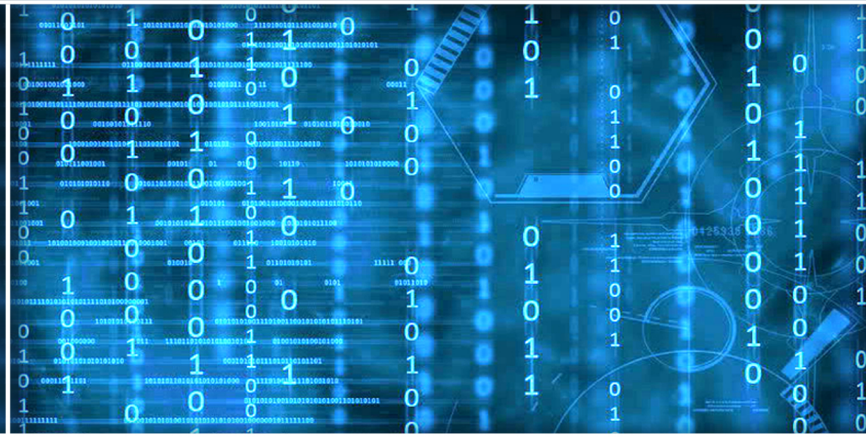


Volume 11 Issue 9

September 2020



ISSN 2156-5570(Online)

ISSN 2158-107X(Print)



# 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!**

**Managing Editor**  
**IJACSA**  
**Volume 11 Issue 9 September 2020**  
**ISSN 2156-5570 (Online)**  
**ISSN 2158-107X (Print)**  
**©2013 The Science and Information (SAI) Organization**

# 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

**Chao-Tung Yang**

**Department of Computer Science, Tunghai University, Taiwan**

*Domain of Research: Software Engineering and Quality, High Performance Computing, Parallel and Distributed Computing, Parallel Computing*

**Elena SCUTELNICU**

**"Dunarea de Jos" University of Galati, Romania**

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

**Krassen Stefanov**

**Professor at Sofia University St. Kliment Ohridski, Bulgaria**

*Domains of Research: e-Learning, Agents and Multi-agent Systems, Artificial Intelligence, Big Data, Cloud Computing, Data Retrieval and Data Mining, Distributed Systems, e-Learning Organisational Issues, e-Learning Tools, Educational Systems Design, Human Computer Interaction, Internet Security, Knowledge Engineering and Mining, Knowledge Representation, Ontology Engineering, Social Computing, Web-based Learning Communities, Wireless/ Mobile Applications*

**Maria-Angeles Grado-Caffaro**

**Scientific Consultant, Italy**

*Domain of Research: Electronics, Sensing and Sensor Networks*

**Mohd Helmy Abd Wahab**

**Universiti Tun Hussein Onn Malaysia**

*Domain of Research: Intelligent Systems, Data Mining, Databases*

**T. V. Prasad**

**Lingaya's University, India**

*Domain of Research: Intelligent Systems, Bioinformatics, Image Processing, Knowledge Representation, Natural Language Processing, Robotics*

# CONTENTS

Paper 1: Efficient GPU Implementation of Multiple-Precision Addition based on Residue Arithmetic

Authors: Konstantin Isupov, Vladimir Knyazkov

PAGE 1 – 8

Paper 2: Classification of Pulmonary Nodule using New Transfer Method Approach

Authors: Syed Waqas Gillani, Bo Ning

PAGE 9 – 13

Paper 3: 6G: Envisioning the Key Technologies, Applications and Challenges

Authors: Syed Agha Hassnain Mohsan, Alireza Mazinani, Warda Malik, Imran Younas, Nawaf Qasem Hamood Othman, Hussain Amjad, Arfan Mahmood

PAGE 14 – 23

Paper 4: Maximum Likelihood Classification based on Classified Result of Boundary Mixed Pixels for High Spatial Resolution of Satellite Images

Authors: Kohei Arai

PAGE 24 – 30

Paper 5: Weather Variability Forecasting Model through Data Mining Techniques

Authors: Sultan Shekana, Addisu Mulugeta, Durga Prasad Sharma

PAGE 31 – 41

Paper 6: A Recommender System for Mobile Applications of Google Play Store

Authors: Ahlam Fuad, Sahar Bayoumi, Hessah Al-Yahya

PAGE 42 – 50

Paper 7: Factored Phrase-based Statistical Machine Pre-training with Extended Transformers

Authors: Vivien L. Beyala, Marcellin J. Nkenlifack, Perrin Li Litet

PAGE 51 – 59

Paper 8: Reward-Based DSM Program for Residential Electrical Loads in Smart Grid

Authors: Muthuselvi G, Saravanan B

PAGE 60 – 68

Paper 9: SQ-Framework for Improving Sustainability and Quality into Software Product and Process

Authors: Kamal Uddin Sarker, Aziz Bin Deraman, Raza Hasan, Ali Abbas

PAGE 69 – 78

Paper 10: Forecasting the Global Horizontal Irradiance based on Boruta Algorithm and Artificial Neural Networks using a Lower Cost

Authors: Abdulatif Aoihan Alresheedi, Mohammed Abdullah Al-Hagery

PAGE 79 – 92

Paper 11: Towards Computational Models to Theme Analysis in Literature

Authors: Abdulfattah Omar

PAGE 93 – 99

Paper 12: Pynq-YOLO-Net: An Embedded Quantized Convolutional Neural Network for Face Mask Detection in COVID-19 Pandemic Era

Authors: Yahia Said

PAGE 100 – 106

Paper 13: Best Path in Mountain Environment based on Parallel Hill Climbing Algorithm

Authors: Raja Masadeh, Ahmad Sharieh, Sanad Jamal, Mais Haj Qasem, Bayan Alsaaidah

PAGE 107 – 116

Paper 14: High-Speed and Secure Elliptic Curve Cryptosystem for Multimedia Applications

Authors: Mohammad Alkhatib

PAGE 117 – 129

Paper 15: A Survey on Privacy Vulnerabilities in Permissionless Blockchains

Authors: Aisha Zahid Junejo, Manzoor Ahmed Hashmani, Abdullah Abdulrehman Alabdulatif

PAGE 130 – 139

Paper 16: Efficient Method for Three Loop MMSE-SIC based Iterative MIMO Systems

Authors: Zuhaibuddin Bhutto, Saleem Ahmed, Syed Muhammad Shehram Shah, Azhar Iqbal, Faraz Mehmood, Imdadullah Thaheem, Ayaz Hussain

PAGE 140 – 145

Paper 17: An Empirical Comparison of Fake News Detection using different Machine Learning Algorithms

Authors: Abdulaziz Albahr, Marwan Albahar

PAGE 146 – 152

Paper 18: Cloud-Based Outsourcing Framework for Efficient IT Project Management Practices

Authors: Mesfin Alemu, Abel Adane, Bhupesh Kumar Singh, Durga Prasad Sharma

PAGE 153 – 164

Paper 19: A Clustering Hybrid Algorithm for Smart Datasets using Machine Learning

Authors: Dar Masroof Amin, Munishwar Rai

PAGE 165 – 172

Paper 20: Grid Search Tuning of Hyperparameters in Random Forest Classifier for Customer Feedback Sentiment Prediction

Authors: Siji George C G, B.Sumathi

PAGE 173 – 178

Paper 21: Aspect-Based Sentiment Analysis and Emotion Detection for Code-Mixed Review

Authors: Andi Suciati, Indra Budi

PAGE 179 – 186

Paper 22: Pseudo Amino Acid Feature-based Protein Function Prediction using Support Vector Machine and K-Nearest Neighbors

Authors: Anjna Jayant Deen, Manasi Gyanchandani

PAGE 187 – 195

Paper 23: Real Time Implementation and Comparison of ESP8266 vs. MSP430F2618 QoS Characteristics for Embedded and IoT Applications

Authors: Krishnaveni Kommuri, Venkata Ratnam Kolluru

PAGE 196 – 202

**Paper 24: A Critical Analysis of IS Governance Frameworks: A Metamodel of the Integrated Use of CobiT Framework**

*Authors: Lamia MOUDOUBAH, Abir EL YAMAMI, Mansouri KHALIFA, Mohammed QBADOU*

**PAGE 203 – 209**

**Paper 25: Prioritization of Software Functional Requirements from Developers Perspective**

*Authors: Muhammad Yaseen, Aida Mustapha, Noraini Ibrahim*

**PAGE 210 – 224**

**Paper 26: Understanding user Emotions Through Interaction with Persuasive Technology**

*Authors: Wan Nooraishya Wan Ahmad, Nazlena Mohamad Ali, Ahmad Rizal Ahmad Rodzuan*

**PAGE 225 – 235**

**Paper 27: Decision-Making Analysis using Arduino-Based Electroencephalography (EEG): An Exploratory Study for Marketing Strategy**

*Authors: Ahmad Faiz Yazid, Siti Munirah Mohd, Abdul Razzak Khan Rustum Ali Khan, Shafinah Kamarudin, Nurhidaya Mohamad Jan*

**PAGE 236 – 243**

**Paper 28: Finger Movement Discrimination of EMG Signals Towards Improved Prosthetic Control using TFD**

*Authors: E.F. Shair, N.A. Jamaluddin, A.R. Abdullah*

**PAGE 244 – 251**

**Paper 29: Autism Spectrum Disorder Diagnosis using Optimal Machine Learning Methods**

*Authors: Maitha Rashid Alfeneiji, Layla Mohammed Alqaydi, Muhammad Usman Tariq*

**PAGE 252 – 260**

**Paper 30: Educational Tool for Generation and Analysis of Multidimensional Modeling on Data Warehouse**

*Authors: Elena Fabiola Ruiz Ledesma, Elizabeth Moreno Galván, Enrique Alfonso Carmona García, Laura Ivoone Garay Jiménez*

**PAGE 261 – 267**

**Paper 31: The Effects of Speed and Altitude on Wireless Air Pollution Measurements using Hexacopter Drone**

*Authors: Rami Noori, Dahlila Putri Dahnil*

**PAGE 268 – 276**

**Paper 32: An Improved Image Retrieval by Using Texture Color Descriptor with Novel Local Textural Patterns**

*Authors: Punit Kumar Johari, Rajendra Kumar Gupta*

**PAGE 277 – 286**

**Paper 33: A Review of Recommender Systems for Choosing Elective Courses**

*Authors: Mfowabo Maphosa, Wesley Doorsamy, Babu Paul*

**PAGE 287 – 295**

**Paper 34: Susceptible, Infectious and Recovered (SIR Model) Predictive Model to Understand the Key Factors of COVID-19 Transmission**

*Authors: DeepaRani Gopagoni, P V Lakshmi*

**PAGE 296 – 302**

**Paper 35: Automated Estrus Detection for Dairy Cattle through Neural Networks and Bounding Box Corner Analysis**

*Authors: Nilo M. Arago, Chris I. Alvarez, Angelita G. Mabale, Charl G. Legista, Nicole E. Repiso, Rodney Rafael A. Robles, Timothy M. Amado, Romeo Jr. L. Jorda, August C. Thio-ac, Jessica S. Velasco, Lean Karlo S. Tolentino*

**PAGE 303 – 311**

**Paper 36: An Efficient Cluster-Based Approach to Thwart Wormhole Attack in Adhoc Networks**

*Authors: Kollu Spurthi, T N Shankar*

**PAGE 312 – 316**

**Paper 37: A Proposed User Requirements Document for Children's Learning Application**

*Authors: Mira Kania Sabariah, Paulus Insap Santosa, Ridi Ferdiana*

**PAGE 317 – 324**

**Paper 38: Design and Implementation of Real Time Data Acquisition System using Reconfigurable SoC**

*Authors: Dharmavaram Asha Devi, Tirumala Satya Savithri, Sai Sugun.L*

**PAGE 325 – 331**

**Paper 39: GAIT based Behavioral Authentication using Hybrid Swarm based Feed Forward Neural Network**

*Authors: Gogineni Krishna Chaitanya, Krovi Raja Sekhar*

**PAGE 332 – 339**

**Paper 40: Electricity Cost Prediction using Autoregressive Integrated Moving Average (ARIMA) in Korea**

*Authors: Safdar Ali, Do-Hyeun Kim*

**PAGE 340 – 344**

**Paper 41: Deep Learning Bidirectional LSTM based Detection of Prolongation and Repetition in Stuttered Speech using Weighted MFCC**

*Authors: Sakshi Gupta, Ravi S. Shukla, Rajesh K. Shukla, Rajesh Verma*

**PAGE 345 – 356**

**Paper 42: Using of Redundant Signed-Digit Numeral System for Accelerating and Improving the Accuracy of Computer Floating-Point Calculations**

*Authors: Otsokov Sh.A, Magomedov Sh.G*

**PAGE 357 – 363**

**Paper 43: Self-Configurable Current-Mirror Technique for Parallel RGB Light-Emitting Diodes (LEDs) Strings**

*Authors: Shaheer Shaida Durrani, Asif Nawaz, Muhamamd Shahzad, Rehan Ali Khan, Abu Zaharin Ahmad, Ahmed Ali Shah, Sheeraz Ahmed, Zeeshan Najam*

**PAGE 364 – 371**

**Paper 44: Implementation of a Clinical Decision Support Systems-Based Neonatal Monitoring System Framework**

*Authors: Sobowale A. A, Olaniyan O. M, Adetan. O, Adanigbo. O, Esan. A, Olusesi. A.T, Wahab. W.B, Adewumi. O. A*

**PAGE 372 – 377**

**Paper 45: Machine Learning-Based Phishing Attack Detection**

*Authors: Sohrab Hossain, Dhiman Sarma, Rana Joyti Chakma*

**PAGE 378 – 388**

**Paper 46: Identifying Critical Success Factors of Financial ERP System in Higher Education Institution using ADVIAN® Method**

*Authors: Ayogeboh Epizitone, Oludayo. O. Olugbara*

**PAGE 389 – 403**

**Paper 47: Machine Learning based Analysis on Human Aggressiveness and Reactions towards Uncertain Decisions**

*Authors: Sohaib Latif, Abdul Kadir Abdullahi Hasan, Abdaziz Omar Hassan*

**PAGE 404 – 408**

**Paper 48: A Cluster-Based Mitigation Strategy Against Security Attacks in Wireless Sensor Networks**

*Authors: Jahangir Khan, Ansar Munir Shah, Babar Nawaz, Khalid Mahmood, Muhammad Kashif Saeed, Mehmood ul Hassan*

**PAGE 409 – 414**

**Paper 49: A Predictive Model for the Determination of Academic Performance in Private Higher Education Institutions**

*Authors: Francis Makombe, Manoj Lall*

**PAGE 415 – 419**

**Paper 50: The Effect of Requirements Quality and Requirements Volatility on the Success of Information Systems Projects**

*Authors: Eman Osama, Ayman Khedr, Mohamed Abdelsalam*

**PAGE 420 – 425**

**Paper 51: Dissemination and Implementation of THK-ANEKA and SAW-Based Stake Model Evaluation Website**

*Authors: Dewa Gede Hendra Divayana, I Putu Wisna Ariawan, Agus Adiarta*

**PAGE 426 – 436**

**Paper 52: Physically-Based Animation in Performing Accuracy Bouncing Simulation**

*Authors: Loh Ngik Hoon*

**PAGE 437 – 444**

**Paper 53: Physiotherapy: Design and Implementation of a Wearable Sleeve using IMU Sensor and VR to Measure Elbow Range of Motion**

*Authors: Anzalna Narejo, Attiya Baqai, Neha Sikandar, Absar Ali, Sanam Narejo*

**PAGE 445 – 455**

**Paper 54: Outlier Detection using Nonparametric Depth-Based Techniques in Hydrology**

*Authors: Insia Hussain*

**PAGE 456 – 462**

**Paper 55: A Hybrid Approach to Enhance Scalability, Reliability and Computational Speed in LoRa Networks**

*Authors: S. Raja Gopal, V S V Prabhakar*

**PAGE 463 – 469**

**Paper 56: A New Online Plagiarism Detection System based on Deep Learning**

*Authors: El Mostafa Hambli, Faouzia Benabbou*

**PAGE 470 – 478**

**Paper 57: Impact of Project-Based Learning on Networking and Communications Competences**

*Authors: Cristian Castro-Vargas, Maritza Cabana-Caceres, Laberiano Andrade-Arenas*

**PAGE 479 – 488**

**Paper 58: Artificial Intelligent Techniques for Palm Date Varieties Classification**

*Authors: Lazhar Khriji, Ahmed Chiheb Ammari, Medhat Awadalla*

**PAGE 489 – 495**

**Paper 59: Computational Analysis of Arabic Cursive Steganography using Complex Edge Detection Techniques**

*Authors: Anwar H. Ibrahim, Abdulrahman S. Alturki*

**PAGE 496 – 500**



**Paper 60: The Most Efficient Classifiers for the Students' Academic Dataset**

*Authors: Ebtehal Ibrahim Al-Fairouz, Mohammed Abdullah Al-Hagery*

**PAGE 501 – 506**

**Paper 61: An Empirical Study of e-Learning Interface Design Elements for Generation Z**

*Authors: Hazwani Nordin, Dalbir Singh, Zulkefli Mansor*

**PAGE 507 – 515**

**Paper 62: Cotton Leaf Image Segmentation using Modified Factorization-Based Active Contour Method**

*Authors: Bhagya M Patil, Basavaraj Amarapur*

**PAGE 516 – 521**

**Paper 63: Trading Saudi Stock Market Shares using Multivariate Recurrent Neural Network with a Long Short-term Memory Layer**

*Authors: Fahd A. Alturki, Abdullah M. Aldughaiyem*

**PAGE 522 – 528**

**Paper 64: Implementing the Behavioral Semantics of Diagrammatic Languages by Co-simulation**

*Authors: Daniel-Cristian Craciunean*

**PAGE 529 – 536**

**Paper 65: A Cluster based Non-Linear Regression Framework for Periodic Multi-Stock Trend Prediction on Real Time Stock Market Data**

*Authors: Lakshmana Phaneendra Maguluri, R. Ragupathy*

**PAGE 537 – 551**

**Paper 66: Development of a Graphic Information System Applied to Quality Statistic Control in Production Processes**

*Authors: Laura Vazquez, Alicia Valdez, Griselda Cortes, Mariana Rosales*

**PAGE 552 – 558**

**Paper 67: Netnography and Text Mining to Understand Perceptions of Indian Travellers using Online Travel Services**

*Authors: Dashrath Mane, Prateek Srivastava*

**PAGE 559 – 569**

**Paper 68: Multi-Dimensional Fraud Detection Metrics in Business Processes and their Application**

*Authors: Badr Omair, Ahmad Alturki*

**PAGE 570 – 586**

**Paper 69: Video Processing for Animation at Key Points of Movement in the Mimosa Pudica**

*Authors: Rodolfo Romero-Herrera, Laura Mendez-Segundo*

**PAGE 587 – 594**

**Paper 70: Population based Optimized and Condensed Fuzzy Deep Belief Network for Credit Card Fraudulent Detection**

*Authors: Jisha M.V, D. Vimal Kumar*

**PAGE 595 – 602**

**Paper 71: Meta-Analysis of Artificial Intelligence Works in Ubiquitous Learning Environments and Technologies**

*Authors: Caitlin Sam, Nalindren Naicker, Mogiveny Rajkoomar*

**PAGE 603 – 613**

**Paper 72: Hate Speech Detection in Twitter using Transformer Methods**

*Authors: Raymond T Mutanga, Nalindren Naicker, Oludayo O Olugbara*

**PAGE 614 – 620**

**Paper 73: VerbNet based Citation Sentiment Class Assignment using Machine Learning**

*Authors: Zainab Amjad, Imran Ihsan*

**PAGE 621 – 627**

**Paper 74: DBSR: A Depth-Based Secure Routing Protocol for Underwater Sensor Networks**

*Authors: Ayman Alharbi*

**PAGE 628 – 634**

**Paper 75: Product Recommendation in Offline Retail Industry by using Collaborative Filtering**

*Authors: Bayu Yudha Pratama, Indra Budi, Arlisa Yuliawati*

**PAGE 635 – 643**

**Paper 76: Using Wearable Sensors for Human Activity Recognition in Logistics: A Comparison of Different Feature Sets and Machine Learning Algorithms**

*Authors: Abbas Shah Syed, Zafi Sherhan Syed, Muhammad Shehram Shah, Salahuddin Saddar*

**PAGE 644 – 649**

**Paper 77: Real-Time Healthcare Monitoring System using Online Machine Learning and Spark Streaming**

*Authors: Fawzya Hassan, Masoud E. Shaheen, Radhya Sahal*

**PAGE 650 – 658**

**Paper 78: Fundamental Capacity Analysis for Identically Independently Distributed Nakagami-q Fading Wireless Communication**

*Authors: Siam Bin Shawkat, Md. Mazid-Ul-Haque, Md. Sohidul Islam, Borshan Sarker Sonok*

**PAGE 659 – 663**

**Paper 79: Unified Approach for White Blood Cell Segmentation, Feature Extraction, and Counting using Max-Tree Data Structure**

*Authors: Bilkis Jamal Ferdosi*

**PAGE 664 – 673**

**Paper 80: DistB-SDoIndustry: Enhancing Security in Industry 4.0 Services based on Distributed Blockchain through Software Defined Networking-IoT Enabled Architecture**

*Authors: Anichur Rahman, Umme Sara, Dipanjali Kundu, Saiful Islam, Md. Jahidul Islam, Mahedi Hasan, Ziaur Rahman, Mostofa Kamal Nasir*

**PAGE 674 – 681**

**Paper 81: Small-LRU: A Hardware Efficient Hybrid Replacement Policy**

*Authors: Purnendu Das, Bishwa Ranjan Roy*

**PAGE 682 – 686**

**Paper 82: Parameter Estimation of the ALBA Autonomous Surface Craft**

*Authors: Melanie M. Valdivia-Fernandez, Brayan A. Monroy-Ochoa, Daniel D. Yanyachi, Juan C. Cutipa-Luque*

**PAGE 687 – 693**

**Paper 83: Mobility-Aware Container Migration in Cloudlet-Enabled IoT Systems using Integrated Muticriteria Decision Making**

*Authors: Mutaz A. B. Al-Tarawneh*

**PAGE 694 – 701**

**Paper 84: Disaster Recovery in Cloud Computing Systems: An Overview**

*Authors: Abedallah Zaid Abualkishik, Ali A. Alwan, Yonis Gulzar*

**PAGE 702 – 710**

**Paper 85: An IoT based Urban Areas Air Quality Monitoring Prototype**

*Authors: Martin M. Soto-Cordova, Martha Medina-De-La-Cruz, Anderson Mujaico-Mariano*

**PAGE 711 – 716**

**Paper 86: Modeling and Interpretation of Covid-19 Infections Data at Peru through the Mitchell's Criteria**

*Authors: Huber Nieto-Chaupis*

**PAGE 717 – 722**

# Efficient GPU Implementation of Multiple-Precision Addition based on Residue Arithmetic

Konstantin Isupov<sup>1</sup>

Department of Electronic Computing Machines  
Vyatka State University  
Kirov, Russia 610000

Vladimir Knyazkov<sup>2</sup>

Research Institute of Fundamental and Applied Studies  
Penza State University  
Penza, Russia 440026

**Abstract**—In this work, the residue number system (RNS) is applied for efficient addition of multiple-precision integers using graphics processing units (GPUs) that support the Compute Unified Device Architecture (CUDA) platform. The RNS allows calculations with the digits of a multiple-precision number to be performed in an element-wise fashion, without the overhead of communication between them, which is especially useful for massively parallel architectures such as the GPU architecture. The paper discusses two multiple-precision integer algorithms. The first algorithm relies on *if-else* statements to test the signs of the operands. In turn, the second algorithm uses radix complement RNS arithmetic to handle negative numbers. While the first algorithm is more straightforward, the second one avoids branch divergence among threads that concurrently compute different elements of a multiple-precision array. As a result, the second algorithm shows significantly better performance compared to the first algorithm. Both algorithms running on an NVIDIA RTX 2080 Ti GPU are faster than the multi-core GNU MP implementation running on an Intel Xeon 4100 processor.

**Keywords**—Multiple-precision algorithm; integer arithmetic; residue number system; GPU; CUDA

## I. INTRODUCTION

Multiple-precision integer arithmetic, which provides operations with numbers that consist of more than 32 or 64 bits, is an important and often indispensable method for solving scientific and engineering problems that are difficult to solve using the standard numerical precision. The most notable application of multiple-precision integer arithmetic is cryptography, where the level of security depends on the length of the keys [1], [2]. Multiple precision is also required in computer algebra (symbolic computation) systems, which operate with mathematical expressions instead of fixed-precision integer and floating-point numbers [3]. The intermediate data produced during a computation may be very large, and multiple-precision arithmetic is required to prevent overflow. Another problem requiring computations with very large integers and of practical interest in polymer physics is counting Hamiltonian cycles on two- and three-dimensional lattices, triangular grid graph, and other structures [4]. Multiple-precision arithmetic is becoming more and more in demand as the scale of computations increases.

There are several approaches for implementing multiple-precision arithmetic. One of them is special software libraries that emulate operations with large numbers using standard fixed-precision operations. Some of the well-known libraries for central processors (CPUs) include the GNU MP Bignum

Library (GMP) [5], the Library for doing Number Theory (NTL) [6], and the Fast Library for Number Theory (FLINT) [7]. There are also works devoted to the implementation of integer arithmetic operations with arbitrary/multiple precision on GPUs [8], [9], [10], [11], [12].

A higher level of arithmetic precision is also supported in a number of programming languages, e.g., Python (the built-in *int* type), Ruby (the built-in *Bignum* type), Perl (*Math::BigInt*), Java (the *BigInteger* class), Haskell (the *Integer* datatype), and C# (*BigInteger*). Another actual approach is to develop hardware accelerators that support integer and floating-point computations with multiple precision [13], [14], [15].

Previous research in [8], [9], [13], and [15] use the traditional way of representing multiple-precision numbers, according to which a number is represented as an array of weighted digits in some base, and the digits themselves are machine-precision numbers [16]. The need for carry propagation under this number representation is one of the main bottleneck of efficient multiple-precision algorithms.

This paper deals with another type of multiple-precision arithmetic, which is based on the residue number system (RNS) [17], [18]. In the RNS, a number is represented by its residues relative to a set of moduli. The moduli are mutually independent, so multiple-precision integer operations such as addition, subtraction, and multiplication are replaced by groups of reduced-precision operations with residues performed in an element-wise fashion and without the overhead of manipulating carry information between the residues.

Recently, a new software library has been developed for efficient residue number system computations on CPU and GPU architectures. The library is called GRNS and is freely available for download at <https://github.com/kisupov/grns>. GRNS is designed for arbitrary moduli sets with large dynamic ranges that far exceed the usual word length of computers, up to several thousand bits. In addition to a number of optimized non-modular RNS operations such as magnitude comparison and division, GRNS implements multiple-precision integer arithmetic. This paper considers two multiple-precision addition algorithms implemented in GRNS. Along with multiplication, addition and subtraction is key operations for many computational algorithms, e.g., fast Fourier transform. Multiple-precision addition is usually considered to be faster and easier than multiplication. However, in the case of RNS, signed addition is more difficult than multiplication as it requires determining the sign of the result, which is a time-consuming

operation for RNS.

Both of our multiple-precision addition algorithms use an interval floating-point evaluation technique for efficient RNS sign determination [19]. However, the first algorithm relies on *if-else* statements to test the signs of the operands, while the second one uses the radix complement RNS notation for negative numbers. It is shown that the second algorithm is better suited for implementation on massively parallel GPU architectures than the first algorithm.

The rest of this paper is organized as follows. Section II provides the background on RNS arithmetic. Section III describes the RNS-based format of multiple-precision integer numbers. Multiple-precision addition algorithms are presented in Section IV. Performance comparison results are given in Section V, and Section VI concludes the paper.

## II. BACKGROUND ON RNS ARITHMETIC

An RNS is specified by a set of  $n$  pairwise prime moduli  $\{m_0, m_1, \dots, m_{n-1}\}$ . The dynamic range of the RNS is  $M = m_0 \cdot m_1 \cdot \dots \cdot m_{n-1}$ . The mapping of an integer  $X$  into the RNS is defined to be the  $n$ -tuple  $(x_0, x_1, \dots, x_{n-1})$ , where  $x_i = |X|_{m_i}$  is the smallest non-negative remainder when  $X$  is divided by  $m_i$ , that is,  $x_i = X \bmod m_i$ . Within the RNS there is a unique representation of all integers in the range from 0 to  $M - 1$ . Namely, the Chinese Remainder Theorem (CRT) states that [18]

$$|X|_M = \left| \sum_{i=0}^{n-1} M_i |x_i w_i|_{m_i} \right|_M, \quad (1)$$

where  $M_i = M/m_i$ , and  $w_i = |M_i^{-1}|_{m_i}$  is the modulo  $m_i$  multiplicative inverse of  $M_i$ .

Since the RNS moduli are independent of each other, arithmetic operations such as addition, subtraction, and multiplication can be computed efficiently. If  $X$ ,  $Y$ , and  $Z$  have RNS representations given by  $(x_0, x_1, \dots, x_{n-1})$ ,  $(y_0, y_1, \dots, y_{n-1})$ ,  $(z_0, z_1, \dots, z_{n-1})$ , then denoting  $\circ$  to represent  $+$ ,  $-$ , or  $\times$ , the RNS version of the  $Z = X \circ Y$ , satisfies

$$Z = (z_0, z_1, \dots, z_{n-1}) = (|x_0 \circ y_0|_{m_0}, |x_1 \circ y_1|_{m_1}, \dots, |x_{n-1} \circ y_{n-1}|_{m_{n-1}}) \quad (2)$$

provided that  $Z \in [0, M - 1]$ . Thus the  $i$ th RNS digit, namely  $z_i$ , is defined in terms of  $|x_i \circ y_i|_{m_i}$  only. That is, no carry information need be communicated between residue digits, and the overhead of manipulating carry information in more traditional, weighted-number systems can be avoided [20].

The disadvantage of RNS is the high complexity of estimating the magnitude of a number, which is required to perform number comparison, sign calculation, overflow checking, division, and some other operations. The classic technique to perform these operations is based on the CRT formula (1) and consists in computing the binary representations of numbers with their subsequent analysis. However, in large dynamic ranges (e.g., a few thousand bits) this technique becomes slow. Other methods for evaluating the magnitude of residue numbers are based on the mixed-radix conversion (MRC) process [21]. But these methods are often also ineffective since they require a lot of arithmetic operations with residues or the use of unacceptably large lookup tables.

An alternative method for implementing time-consuming operations in the RNS is based on computing the floating-point interval evaluation of the fractional representation of an RNS number [19]. This method is designed to be fast on modern general-purpose computing platforms that support efficient finite precision floating-point arithmetic operations such as IEEE 754 operations. For a given RNS number  $X = (x_0, x_1, \dots, x_{n-1})$ , the floating-point interval evaluation is an interval defined by its lower and upper bounds (endpoints)  $\underline{X/M}$  and  $\overline{X/M}$  that are finite precision floating-point numbers satisfying  $\underline{X/M} \leq X/M \leq \overline{X/M}$ . The floating-point interval evaluation is denoted by  $I(X/M) = [\underline{X/M}, \overline{X/M}]$ .

Thus,  $I(X/M)$  provides information about the range of changes in the fractional representation (also called relative value) of an RNS number. This information may not be sufficient to restore the binary representation, but it can be efficiently used to perform other difficult operations in RNS, e.g., magnitude comparison, sign detection, and division.

The most important benefit of this method is that computation of  $I(X/M)$  requires only standard arithmetic operations, and no residue-to-binary conversion is required. For a given RNS representation  $(x_0, x_1, \dots, x_{n-1})$ , the calculation of the bounds of  $I(X/M)$  is performed on average in linear and logarithmic time for sequential and parallel cases, respectively. Furthermore, the following arithmetic operations are defined:

$$\begin{aligned} I(X/M) + I(Y/M) &= [\underline{X/M} \nabla \underline{Y/M}, \overline{X/M} \triangle \overline{Y/M}], \\ I(X/M) - I(Y/M) &= [\underline{X/M} \nabla \overline{Y/M}, \overline{X/M} \triangle \underline{Y/M}], \\ I(X/M) \times I(Y/M) &= [\underline{X/M} \nabla \underline{Y/M} \nabla W, \overline{X/M} \triangle \overline{Y/M} \triangle V], \\ I(X/M) \div I(Y/M) &= [\underline{X/M} \nabla V \nabla \overline{Y/M}, \overline{X/M} \triangle W \triangle \underline{Y/M}]. \end{aligned} \quad (3)$$

In these interval formulas, the following notation are used:

- $\nabla, \nabla, \nabla$  and  $\nabla$  stand for the floating-point operations of addition, subtraction, multiplication, and division, performed with rounding downwards;
- $\triangle, \triangle, \triangle$  and  $\triangle$  stand for the floating-point operations of addition, subtraction, multiplication, and division, performed with rounding upwards;
- $V$  is the greatest floating-point number that is less than or equal to  $1/M$ ;
- $W$  is the least floating-point number greater than or equal to  $1/M$ .

Interval formulas (3) are useful in that they do not limit the possible values of the result interval in the range of  $[0, 1]$ . This allows for easy overflow detection or sign identification despite the cyclical (modulo  $M$ ) nature of RNS arithmetic.

Using interval evaluations, new algorithms have been proposed in [19] to efficiently implement several difficult RNS operations, such as number comparison and general division.

## III. NUMBER REPRESENTATION

The format for multiple-precision integers is shown in Fig. 1. A multiple-precision integer  $x$  consists of a sign  $s$ , a significant  $X$  composed of  $n$  significant digits  $(x_0$  to  $x_{n-1})$ , and an interval floating-point evaluation of the significant  $I(X/M) = [\underline{X/M}, \overline{X/M}]$ . The sign is interpreted in the same way as in two's complement representation: the sign is equal

to zero when  $x$  is positive and one when it is negative. The significand expresses the absolute value of  $x$  and is represented in the RNS with the moduli set  $\{m_0, m_1, \dots, m_{n-1}\}$ . The significand digits (residues) are represented as ordinary two's complement integers.

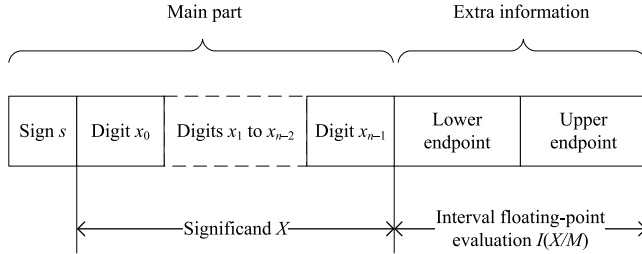


Fig. 1. Multiple-Precision Integer Format.

The size of the moduli set  $n$  specifies the number of digits in the significand. If the product of all RNS moduli is  $M$ , then the precision of  $x$  is equal to  $\lfloor \log_2 M \rfloor$  bits. Thus, changing the size of the moduli set allows one to achieve arbitrary precision.

The following notation is used to denote a multiple-precision integer in the described number format:

$$x = \langle s, X, I(X/M) \rangle. \quad (4)$$

The value of a multiple-precision integer of the form (4) can be computed using the CRT formula:

$$x = (-1)^s \times \left| \sum_{i=0}^{n-1} M_i |x_i w_i|_{m_i} \right|_M. \quad (5)$$

The interval evaluation is included in the number representation as additional information, that is,  $X/M$  and  $\overline{X/M}$  are stored in system memory along with other fields of the multiple-precision integer. This provides efficient comparison, sign computation and overflow detection, allowing one to calculate  $I(X/M)$  in  $O(1)$  time using the formulas (3). Recently, this approach has been successfully used in the context of multiple-precision floating-point arithmetic based on the residue number system [22].

In order to be able to use interval evaluations for virtually any (arbitrarily large) value of  $M$  without worrying about underflow,  $X/M$  and  $\overline{X/M}$  are represented as binary floating-point numbers with an extended exponent range, that is, have the form

$$f \times 2^e, \quad (6)$$

where  $f$  is a regular floating-point number (IEEE 754), and  $e$  is a two's complement integer. This extended-range representation is not intended to improve the level of numerical precision or accuracy, but it does ensure that there is no overflow or underflow when dealing with extremely large or small values.

#### IV. MULTIPLE-PRECISION INTEGER ADDITION

In this section, two algorithms for signed multiple-precision integer addition are presented. A naive implementation is presented first and then an improved one. Step-by-step examples are also provided for both implementations.

##### A. Useful Notation

For given  $a \in \{0, 1\}$  and  $X = (x_0, x_1, \dots, x_{n-1})$ , the paper [22] introduces a function  $B[X, a]$  such that

$$B[X, a] = \begin{cases} X/M, & \text{for } a = 0, \\ \overline{X/M}, & \text{for } a = 1. \end{cases} \quad (7)$$

That is, the lower bound of  $I(X/M)$  is denoted by  $B[X, 0]$ , while the upper one is denoted by  $B[X, 1]$ . Using this notation, we have  $I(X/M) = [B[X, 0], B[X, 1]]$ . This notation is useful in that it allows one to dynamically specify the bound to be accessed. This notation is used in the rest of the present paper.

##### B. Note on Overflow Detection

For the set of RNS moduli  $\{m_0, m_1, \dots, m_{n-1}\}$ , the largest representable integer is  $(M-1)$ , and the result of an arithmetic operation should belong to the interval  $[0, M-1]$  if we want to obtain its valid representation in the RNS. Otherwise, the result will be reduced modulo  $M$ , and this event is classified as an integer overflow. The GRNS library implements efficient overflow detection using the floating-point interval evaluations, but that is beyond the scope of this paper.

---

##### Algorithm 1 Multiple-precision integer addition

---

- 1: **if**  $s_x = s_y$  **then**
  - 2:      $s_z \leftarrow s_x$
  - 3:     **for each**  $i \in \{0, 1, \dots, n-1\}$  **do**
  - 4:          $z_i \leftarrow |x_i + y_i|_{m_i}$
  - 5:     **end for**
  - 6:      $B[Z, 0] \leftarrow B[X, 0] \nabla B[Y, 0]$
  - 7:      $B[Z, 1] \leftarrow B[X, 1] \Delta B[Y, 1]$
  - 8:     **else if**  $B[X, 0] \geq B[Y, 1]$  **then**
  - 9:          $s_z \leftarrow s_x$
  - 10:        **for each**  $i \in \{0, 1, \dots, n-1\}$  **do**
  - 11:             $z_i \leftarrow |x_i - y_i|_{m_i}$
  - 12:        **end for**
  - 13:         $B[Z, 0] \leftarrow B[X, 0] \nabla B[Y, 1]$
  - 14:         $B[Z, 1] \leftarrow B[X, 1] \Delta B[Y, 0]$
  - 15:     **else if**  $B[Y, 0] \geq B[X, 1]$  **then**
  - 16:          $s_z \leftarrow s_y$
  - 17:        **for each**  $i \in \{0, 1, \dots, n-1\}$  **do**
  - 18:             $z_i \leftarrow |y_i - x_i|_{m_i}$
  - 19:        **end for**
  - 20:         $B[Z, 0] \leftarrow B[Y, 0] \nabla B[X, 1]$
  - 21:         $B[Z, 1] \leftarrow B[Y, 1] \Delta B[X, 0]$
  - 22:     **else**
  - 23:         Use mixed-radix conversion to compare the magnitude of  $X$  and  $Y$ . If  $X \geq Y$ , subtract  $Y$  from  $X$  and take  $s_z \leftarrow s_x$ ; otherwise, subtract  $X$  from  $Y$  and take  $s_z \leftarrow s_y$ ; In any case,  $I(Z/M)$  should be recalculated.
  - 24:     **end if**
-

C. Algorithm 1 (Naive Implementation)

1) *Description:* Algorithm 1 takes two multiple-precision integers  $x$  and  $y$  represented as  $x = \langle s_x, X, I(X/M) \rangle$  and  $y = \langle s_y, Y, I(Y/M) \rangle$ , and outputs the sum  $z = x + y$  represented as  $z = \langle s_z, Z, I(Z/M) \rangle$ . This algorithm analyzes the signs of the numbers, and if they are the same, then RNS addition of the significands is performed; otherwise, RNS subtraction is performed. The sign of the result is computed by comparing the magnitude of  $X = (x_0, x_1, \dots, x_{n-1})$  and  $Y = (y_0, y_1, \dots, y_{n-1})$  using the floating-point interval evaluations.

2) *Illustration:* Consider the moduli set  $\{7, 9, 11, 13\}$  with the moduli product  $M = 9009$ . Suppose we are given two integers of the form (3),

$$x = \langle 0, (5, 7, 5, 8), [0.416, 0.420] \rangle,$$

$$y = \langle 1, (3, 7, 6, 4), [0.444, 0.448] \rangle,$$

and we want to find  $z = x + y$ . Since  $B[Y, 0]$  (0.444) is greater than  $B[X, 1]$  (0.420), steps 16 to 21 of the algorithm are performed. They are presented in Table I.

TABLE I. EXAMPLE OF ALGORITHM 1

| Step no. | Calculations                                     |
|----------|--|
| 16       | $s_z = 1$  |
| 17-19    | $Z = (3, 7, 6, 4) - (5, 7, 5, 8) = (5, 0, 1, 9)$ |
| 20       | $B[Z, 0] = 0.444 \nabla 0.420 = 0.024$           |
| 21       | $B[Z, 1] = 0.448 \triangle 0.416 = 0.032$        |

The computed result is  $z = \langle 1, (5, 0, 1, 9), [0.024, 0.032] \rangle$ . We check this result by converting it to decimal:  $z = -243$ . In fact,  $x = 3778$  and  $y = -4021$ .

3) *Drawback:* The main disadvantage of Algorithm 1 is that checking the signs of the operands via conditionals (*if-else* statements) results in branch divergence among threads that concurrently compute different elements of a multiple-precision array. This may be normal for modern multi-core processors with good branch prediction accuracy, but this is a problem for SIMT (single instruction, multiple threads) architectures such as GPUs, where many threads run in lock-step.

For example, a CUDA-compliant GPU consists of an array of streaming multiprocessors (SMs), each of which contains multiple streaming processors. Although each SM can run one or more different instructions, conditionals can greatly decrease performance inside an SM, as each branch of each conditional must be evaluated. Long code paths in a conditional can cause a 2-fold slowdown for each conditional within a warp (a group of 32 threads) and a  $2^N$  slowdown for  $N$  nested conditionals. A maximum 32-time slowdown can occur when each thread in a warp executes a separate condition [23].

This bottleneck is illustrated in Fig. 2, which contains a flowchart of Algorithm 1. In the figure, 14 threads concurrently compute 14 multiple-precision additions on a system that follows the SIMT execution model. The right side of the figure shows threads running at once.

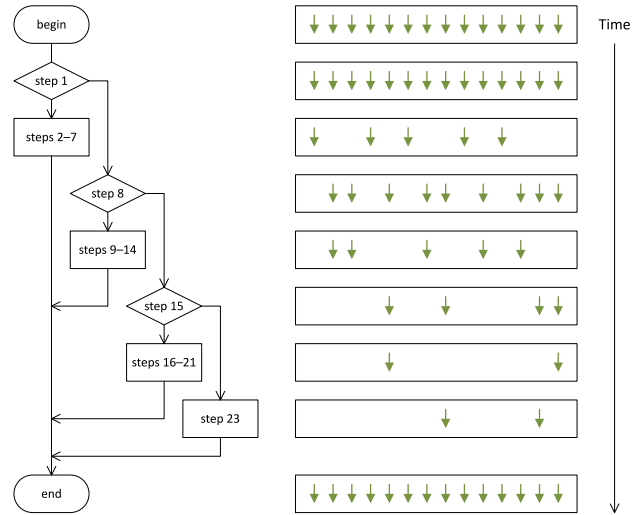


Fig. 2. Flowchart of Algorithm 1.

D. Algorithm 2 (Improved Implementation)

1) *Description:* Algorithm 2 shows how to avoid conditional expressions when adding multiple-precision signed integers. This algorithm is a simplified version of the multiple-precision RNS-based floating-point addition algorithm that was originally proposed in [22]. The main idea is to use the radix-complement representation of a negative number in the RNS. Recall that the precomputed constant  $V$  used in this algorithm is the greatest finite precision floating-point number that is less than or equal to  $1/M$ .

**Algorithm 2** Multiple-precision integer addition using radix complement RNS arithmetic

```

1:  $\alpha \leftarrow (1 - 2s_x)$ 
2:  $\beta \leftarrow (1 - 2s_y)$ 
3: for each  $i \in \{0, 1, \dots, n - 1\}$  do
4:    $z_i \leftarrow (\alpha x_i + \beta y_i) \bmod m_i$ 
5: end for
6:  $B[Z, 0] \leftarrow \alpha B[X, s_x] \nabla \beta B[Y, s_y]$ 
7:  $B[Z, 1] \leftarrow \alpha B[X, (1 - s_x)] \Delta \beta B[Y, (1 - s_y)]$ 
8: if  $B[Z, 0]$  and  $B[Z, 1]$  have the same sign then
9:   Assign the sign of  $B[Z, 0]$  and  $B[Z, 1]$  to  $s_z$ 
10: else
11:   Use mixed-radix conversion to compare  $X$  and  $Y$ :
   • If  $X > Y$ , then assign  $s_z \leftarrow s_x$ .
   • If  $X < Y$ , then assign  $s_z \leftarrow s_y$ .
   • If  $X = Y$ , then assign  $s_z \leftarrow 0$ .
12:    $B[Z, s_z] \leftarrow (1 - 2s_z)V$ 
13: end if
14: if  $s_z = 1$  then
15:   for each  $i \in \{0, 1, \dots, n - 1\}$  do
16:      $z_i \leftarrow (m_i - z_i) \bmod m_i$ 
17:   end for
18:   Swap  $B[Z, 0]$  and  $B[Z, 1]$  with sign inversion, that is,
   set  $B[Z, 0] = -B[Z, 1]$  and  $B[Z, 1] = -B[Z, 0]$ 
19: end if

```

Fig. 3 shows a flowchart of Algorithm 2. The *if-else* statement at steps 8 to 12 cannot be eliminated, since the accuracy of  $B[Z, 0]$  and  $B[Z, 1]$  may be insufficient to unambiguously determine the sign of  $z$ . This ambiguity is possible due to the limited precision arithmetic used in calculating  $B[Z, 0]$  and  $B[Z, 1]$ . However, this is actually a rare case, and it can only occur when the result is too close to zero.

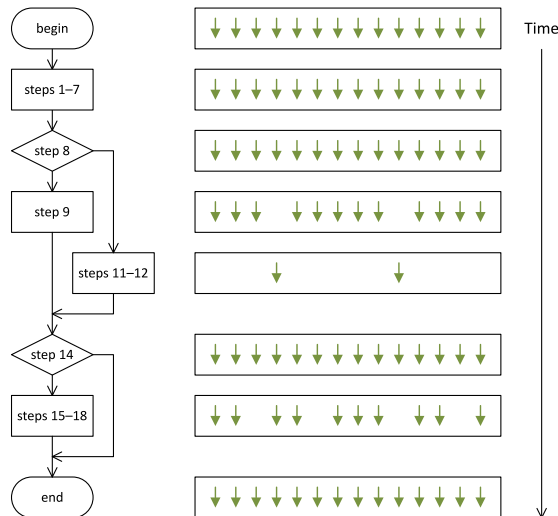


Fig. 3. Flowchart of Algorithm 2.

We note that the *if* statement at step 14 of Algorithm 2 does not cause branch divergence, since there is no the corresponding *else* statement here.

2) *Illustration*: In Table II, Algorithm 2 is used to compute the sum of the numbers from the previous example.

TABLE II. EXAMPLE OF ALGORITHM 2

| Step no. | Calculations  |
|----------|---|
| 1,2      | $\alpha = 1 - 2 \times 0 = 1$ $\beta = 1 - 2 \times 1 = -1$   |
| 3-5      | $z_0 = (5 - 3) \bmod 7 = 2$<br>$z_1 = (7 - 7) \bmod 9 = 0$<br>$z_2 = (5 - 6) \bmod 11 = 10$<br>$z_3 = (8 - 4) \bmod 13 = 4$ |
| 6        | $B[Z, 0] = 0.416 \nabla (-0.448) = -0.032$  |
| 7        | $B[Z, 1] = 0.420 \Delta (-0.444) = -0.024$  |
| 9        | $s_z = 1$   |
| 15-17    | $Z = (7, 9, 11, 13) - (2, 0, 10, 4) = (5, 0, 1, 9)$   |
| 18       | $B[Z, 0] = 0.024, B[Z, 1] = 0.032$  |

Thus, as in the first example, the correct result is computed:  $z = \langle 1, (5, 0, 1, 9), [0.024, 0.032] \rangle$ .

V. PERFORMANCE COMPARISON RESULTS

This section gives comparative results of the presented multiple-precision integer addition algorithms. In the experiments, we used a GeForce RTX 2080 Ti graphics card that has 11 GB of GDDR6 memory, 4352 CUDA cores, and Compute Capability 7.5. This GPU was installed on a machine with an Intel Xeon 4100/8.25M S2066 OEM processor running Ubuntu 18.04.5 LTS, CUDA 10.2 and NVIDIA Driver 450.51.06 were used. The source code was compiled using the nvcc compiler with the *-O3* and *-Xcompiler=-fopenmp* options.

A. Methodology

The parameters of the experiments are shown in Table III. Each dataset was composed of two multiple-precision integer arrays of the same length, and the performance was evaluated for element-by-element addition of the arrays. The performance was measured in the number of multiple-precision arithmetic operations (additions) per second. For comparison purposes, the performance of the GNU MP library was also measured on 4 CPU cores. In the experiments, we considered only the computation time, so the measurements do not include neither the data transfer time nor the time of converting data into internal multiple-precision representations.

TABLE III. EXPERIMENTAL PARAMETERS

| Parameter                       | Value   |
|---------------------------------|---|
| Size of the RNS moduli set, $n$ | from 8 to 256   |
| Bit width of each modulus       | 32  |
| Precision in bits, $p$          | from 128 to 4096  |
| Dataset size                    | 1,000,000   |
| Datasets                        | <i>Dataset-1</i> : pseudo-random integers in the range 0 to $(M - 1)/2$<br><i>Dataset-2</i> : pseudo-random integers in the range $(1 - M)/2$ to 0<br><i>Dataset-3</i> : pseudo-random integers in the range $(1 - M)/2$ to $(M - 1)/2$ |



For each precision  $p$ , a corresponding set of RNS moduli was generated such that

$$\lceil \log_2 M \rceil \geq p, \quad (8)$$

where  $M$  is the product of all the moduli in the set. Table IV shows the relationship between the precision and moduli sets used in the experiments.

TABLE IV. RELATIONSHIP BETWEEN THE PRECISION AND MODULI SETS USED IN THE EXPERIMENTS

| Precision, $p$ | Size of moduli set, $n$ | Dynamic range, $M$ (approx.)  |
|----------------|-------------------------|-------------------------------|
| 128            | 8                       | 3.486474761596273374449E+38   |
| 256            | 16                      | 1.182869237276559892956E+77   |
| 512            | 32                      | 1.381750867498453484869E+154  |
| 1024           | 64                      | 1.834972082650114435387E+308  |
| 2048           | 128                     | 3.267493893788783073405E+616  |
| 4096           | 256                     | 1.113716837551166769174E+1233 |

The moduli sets were generated using Algorithm 3. This algorithm takes as input the smallest odd modulus  $m_0$ , the size of the desired moduli set  $n$ , and produces an increasing sequence of  $n - 1$  consecutive odd integers  $m_1, m_2, \dots, m_{n-1}$  that are coprime to each other and also coprime to  $m_0$ . The value of  $m_0$  is selected by trial and error until the condition (8) is satisfied. The used tool for generating moduli sets is freely available at <https://github.com/kisupov/rns-moduli-generator>.

For the CUDA implementations of the presented multiple-precision addition algorithms, 32 threads per each thread block were used, and the total number of blocks was calculated as follows:

$$nBlocks = \left\lceil \frac{N}{nThreads} \right\rceil + K, \quad (9)$$

where  $N$  is the size of the dataset (1,000,000),  $nThreads = 32$ , and  $K$  is defined as

$$K = \begin{cases} 1, & \text{if } N \bmod nThreads > 0, \\ 0, & \text{otherwise.} \end{cases} \quad (10)$$

The GNU MP library implementation have been accelerated using the OpenMP library.

### B. Results

In the first experiment, the input arrays were filled with pseudo-random non-negative integers ranging from 0 to  $(M - 1)/2$ , where  $M$  is the product of all the RNS moduli. The performance results are shown in Fig. 4.

In the second experiment, the input arrays were filled with pseudo-random non-positive integers ranging from  $(1 - M)/2$  to 0. The results are reported in Fig. 5.

Finally, in the third experiment, the input arrays were filled with pseudo-random positive and negative integers ranging from  $(1 - M)/2$  to  $(M - 1)/2$ . Fig. 6 demonstrates the performance results obtained in this setting.

### Algorithm 3 Moduli set generation

```

1:  $t \leftarrow m_0 + 2$ 
2:  $k \leftarrow 1$ 
3: while  $k < n$  do
4:    $p \leftarrow 1$ 
5:   for  $i \leftarrow 1$  to  $k$  do
6:     if  $\text{gcd}(m_i, t) > 1$  then
7:        $p \leftarrow 0$ 
8:     end if
9:   end for
10:  if  $p = 1$  then
11:     $m_k \leftarrow t$ 
12:     $k \leftarrow k + 1$ 
13:  end if
14:   $t \leftarrow t + 2$ 
15: end while

```

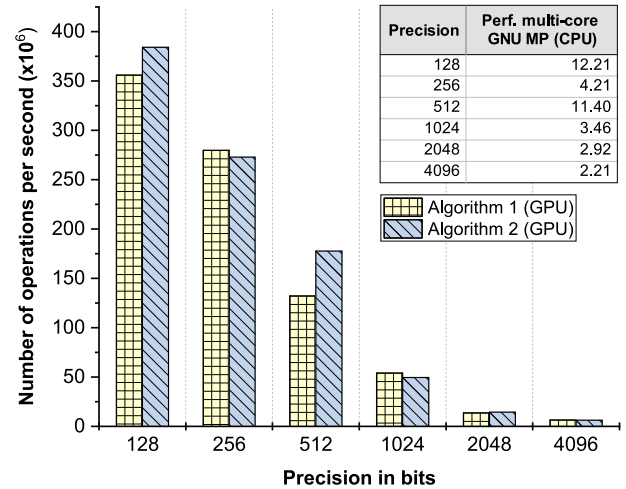


Fig. 4. Performance of Multiple-Precision Integer Addition Implementations with Non-Negative Inputs (Dataset-1).

### C. Discussion

For Dataset-1 (Fig. 4), Algorithm 1 has nearly the same performance as Algorithm 2. This is because in Algorithm 1, all parallel threads follow steps 2–7 and there is no divergent execution paths. The results show that the developed CUDA functions are up to  $65\times$  faster than the parallel CPU implementation using GNU MP.

In the case of Dataset-2 (Fig. 5) the performance of Algorithm 1 remains the same as in the case of Dataset-1, since there are still no branch divergence (all parallel threads follow steps 2–7). In turn, the need to restore negative results reduces the performance of Algorithm 2 by an average of  $1.1\times$  compared to Dataset-1, and this performance degradation does not seem to be significant.

When using Dataset-3 (Fig. 6), branch divergence leads to an average  $1.9$ -fold decrease in the performance of Algorithm

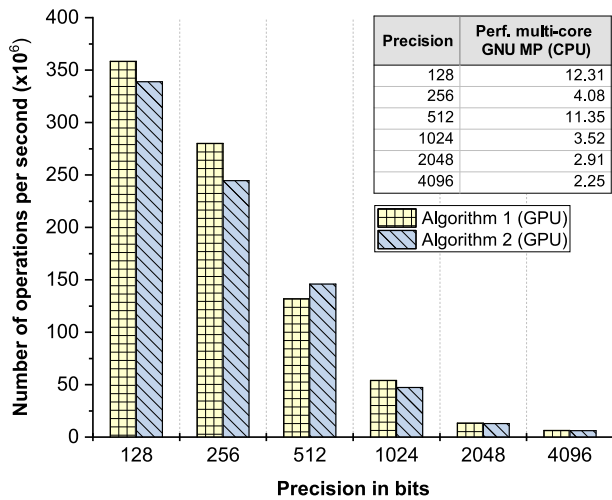


Fig. 5. Performance of Multiple-Precision Integer Addition Implementations with Non-Positive Inputs (Dataset-2).

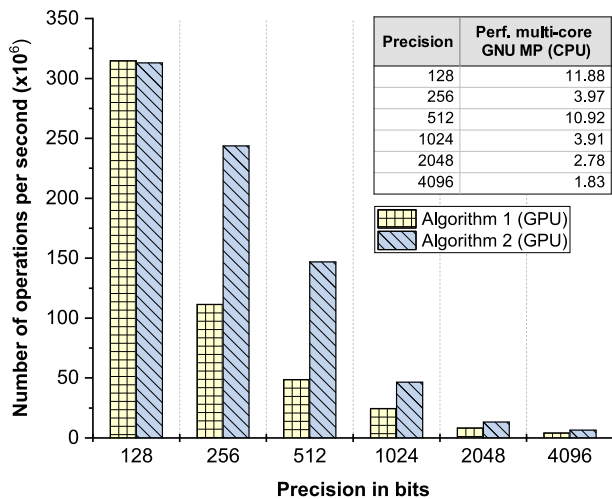


Fig. 6. Performance of Multiple-Precision Integer Addition Implementations with Mixed Positive and Negative Inputs (Dataset-3).

1 compared to Dataset-1 and Dataset-2. With 512-bit precision, the performance of Algorithm 1 is reduced by almost  $3\times$  compared to Dataset-1. In turn, the performance of Algorithm 2 reduced by at most a factor of 1.2 compared to Dataset-1. The net result is that when the operands have different signs, Algorithm 2 outperforms Algorithm 1 by up to  $3\times$ .

A limitation of the proposed CUDA implementations is that the execution time grows linearly with increasing the precision. This happens for the following reasons:

- 1) Each multiple-precision addition is performed as a single thread, that is, the digits of multiple-precision numbers are calculated sequentially.
- 2) As the precision increases, the stride between elements in the input arrays increases accordingly and the effective GPU memory bandwidth decreases.

It should be noted that it is possible to compute all the digits (residues) of multiple-precision significands in parallel across different RNS moduli without worrying about carry

propagation. This parallel arithmetic property of the RNS is employed in [22] to implement GPU-accelerated multiple-precision linear algebra kernels. Furthermore, we note that if all the digits of a multiple-precision number are computed in parallel, then the structure-of-arrays (SoA) layout with a sequential addressing scheme will provide coalesced access to the global GPU memory. Implementing digit-parallel multiple-precision integer addition is a direction for future work.

## VI. CONCLUSION

In this paper, we have considered two multiple-precision integer addition algorithms for graphics processing units. The algorithms are based on the representation of large integers in the residue number system.

The first algorithm uses conditional operators to check the signs of the operands. However, in this case, threads that concurrently compute different elements of a multiple-precision array take divergent execution paths, which leads to an increase in the total computation time. To overcome this disadvantage, the second algorithm uses the radix-complement representation of a negative number in the RNS.

Experiments have shown that when the signs of the operands are different, the second algorithm outperforms the first one by far. In turn, both algorithms running on an NVIDIA RTX 2080 Ti GPU have shown to be faster than the multi-core GNU MP implementation on an Intel Xeon 4100 processor.

The presented implementation is part of GRNS, a library for efficient computations in the residue number system using CUDA-enabled GPUs. In the future, we plan to implement digit-parallel versions of the multiple-precision integer operations to take full advantage of the internal RNS parallelism. Furthermore, we will focus on extending the GRNS functionality and implementing real-world multiple-precision applications using this library.

## ACKNOWLEDGMENT

This research is supported by the Ministry of Science and Higher Education of the Russian Federation, grant id RFMEFI61319X0092.

## REFERENCES

- [1] A. Omondi, *Cryptography Arithmetic*. Springer International Publishing, 2020.
- [2] W. Wang, Y. Hu, L. Chen, X. Huang, and B. Sunar, "Exploring the feasibility of fully homomorphic encryption," *IEEE Transactions on Computers*, vol. 64, no. 3, pp. 698–706, 2015.
- [3] R. Sehgal and V. Nehra, "Symbolic computation of mathematical transforms and its application: A MATLAB computational project-based approach," *IUP Journal of Electrical & Electronics Engineering*, vol. 8, no. 1, pp. 53–76, 2015.
- [4] O. Bodroža-Pantić, H. Kwong, and M. Pantić, "Some new characterizations of Hamiltonian cycles in triangular grid graphs," *Discrete Applied Mathematics*, vol. 201, pp. 1–13, 2016.
- [5] "The GNU multiple precision arithmetic library," 2020. [Online]. Available: <https://gmplib.org/>
- [6] "NTL: A library for doing number theory," 2020. [Online]. Available: <https://shoup.net/ntl/>
- [7] "FLINT: Fast library for number theory," 2020. [Online]. Available: <http://www.flintlib.org/>

- [8] K. Zhao and X. Chu, "GPUMP: A multiple-precision integer library for GPUs," in *Proceedings of the 2010 10th IEEE International Conference on Computer and Information Technology (CIT 2010)*, Bradford, UK, 2010, pp. 1164–1168.
- [9] T. Ewart, A. Hehn, and M. Troyer, "VLI – a library for high precision integer and polynomial arithmetic," in *Supercomputing*, J. M. Kunkel, T. Ludwig, and H. W. Meuer, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2013, pp. 267–278.
- [10] E. Ochoa-Jiménez, L. Rivera-Zamarripa, N. Cruz-Cortés, and F. Rodríguez-Henríquez, "Implementation of RSA signatures on GPU and CPU architectures," *IEEE Access*, vol. 8, pp. 9928–9941, 2020.
- [11] N. Emmart and C. C. Weems, "High precision integer multiplication with a GPU using Strassen's algorithm with multiple FFT sizes," *Parallel Processing Letters*, vol. 21, no. 3, pp. 359–375, 2011.
- [12] B.-C. Chang, B.-M. Goi, R. C.-W. Phan, and W.-K. Lee, "Multiplying very large integer in GPU with Pascal architecture," in *Proceedings of the 2018 IEEE Symposium on Computer Applications Industrial Electronics (ISCAIE)*, Penang, Malaysia, 2018, pp. 401–405.
- [13] K. Rudnicki, T. P. Stefański, and W. Żebrowski, "Open-source coprocessor for integer multiple precision arithmetic," *Electronics*, vol. 9, no. 7, p. article no. 1141, 2020.
- [14] A. Bocco, Y. Durand, and F. De Dinechin, "SMURF: Scalar multiple-precision unum Risc-V floating-point accelerator for scientific computing," in *Proceedings of the Conference for Next Generation Arithmetic 2019*. New York, NY, USA: ACM, 2019.
- [15] M. J. Schulte and E. E. Swartzlander, "A family of variable-precision interval arithmetic processors," *IEEE Transactions on Computers*, vol. 49, no. 5, pp. 387–397, 2000.
- [16] R. Brent and P. Zimmermann, *Modern Computer Arithmetic*. Cambridge: Cambridge University Press, 2010.
- [17] P. V. Ananda Mohan, *Residue Number Systems: Theory and Applications*. Cham: Birkhäuser, 2016.
- [18] A. Omondi and B. Premkumar, *Residue Number Systems: Theory and Implementation*. London, UK: Imperial College Press, 2007.
- [19] K. Isupov, "Using floating-point intervals for non-modular computations in residue number system," *IEEE Access*, vol. 8, pp. 58 603–58 619, 2020.
- [20] F. J. Taylor, "Residue arithmetic a tutorial with examples," *Computer*, vol. 17, no. 5, pp. 50–62, 1984.
- [21] N. S. Szabo and R. I. Tanaka, *Residue Arithmetic and its Application to Computer Technology*. New York, USA: McGraw-Hill, 1967.
- [22] K. Isupov, V. Knyazkov, and A. Kuvaev, "Design and implementation of multiple-precision BLAS level 1 functions for graphics processing units," *Journal of Parallel and Distributed Computing*, vol. 140, pp. 25–36, 2020.
- [23] R. Farber, *CUDA Application Design and Development*. Boston: Morgan Kaufmann, 2011.

# Classification of Pulmonary Nodule using New Transfer Method Approach

Syed Waqas Gillani<sup>1</sup>, Bo Ning<sup>2\*</sup>

College of Information Science and Technology  
Dalian Maritime University, Dalian, Liaoning, China

**Abstract**—Lung cancer is among the world's worst cancers, and accounted for 27% of all cancers in 2018. Despite substantial improvement in recent diagnoses and medications, the five year cure ratio is just 19%. Before even the diagnosis, classification of lung nodule is an essential step, particularly because early detection can help doctors with a highly valued opinion. CT image detection and classification is possible easily and accurately with advanced vision devices and machine-learning technology. This field of work has been extremely successful. Researchers have already attempted to improve the accuracy of CAD structures by computational tomography (CT) in the screening of lung cancer in several deep learning models. In this paper, we proposed a fully automated lung CT system for lung nodule classification, namely, new transfer method (NTM) which has two parts. First features are extracted by applying different VOI and feature extraction techniques. We used intensity, shape, contrast of border and spicula extraction to extract the lung nodule. Then these nodules are transfer to the classification part where we used advance-fully convolution network (A-FCN) to classify the lung nodule between benign and malignant. Our A-FCN network contain three types of layers that helps to enhance the performance and accuracy of NTM network which are convolution layer, pooling layer and fully connected layer. The proposed model is trained on LIDC-IDRI dataset and attained an accuracy of 89.90 % with AUC of 0.9485.

**Keywords**—New transfer method; VOI extraction; feature extraction; classification; LIDC-IDRI dataset

## I. INTRODUCTION

In this modern era of machine learning, doctors are finding some form of support that encourages their ability to analyze and diagnose CT images of patients easily and to identify extremely effective and accurate pathologies. Timely detection and classification of lung nodules enhances clinical results and improves the chances of survival rates [1]. Instead, cancer is now the ultimate common pale pathology that endangers anyone irrespective of age. Between various diseases, lung cancer is the unregulated cell growth in an overt body district [2]. The development of lung nodules is indicated for lung cancer and shows the clinical disease phase [3]. The nodules existing in the lungs equate to variants in lung tissue as from standard, which is round or fit as a 3 millimeter and 30 millimeter wide fiddle [4]. The human body comprises of several cells. The nodule is formed whenever cells develop feral outside of lung [5]. Computed tomography (CT) is mostly used for the detection of lung cancer [6]. Centered on the findings of a nationwide lung screening study performed in the USA [7], scanning for low dose CT scans decreased deaths from lung cancer by 18%. Hence, CT is

considered an appropriate diagnostic technique for the detection of cancer. Computer-aided diagnostic (CAD) methods are designed to help clinicians in analyzing health data and diagnosing the disease. Fig. 1 shows that 24% of CAD's are just used for lung cancer [23]. It means that the medical technology is being rapidly accepted and implemented quickly. CAD could be categorized into two kinds: the detection and diagnostic system (CADE and CADx). CADE's aim is to find the ROI to detect unique abnormalities and CADx offers doctor medical help to discern the form, frequency, level, development and disease deterioration. The CADx could just only use to the diagnosis purpose for example shape, thickness, and texture.

CNN began as an advancement of the deep neural network with the use of convolution algorithms to help interpret the input image. It is built on the basis of the biologically visual field and is therefore quite efficient for difficulties of image recognition, irrespective of size or volume. Ginneken et al. [8] evaluated CNN as well as food over fat as a clinical CAD tool for the identification of lung cancer tumors. LIDC CT-scanning images have been used to locate lesions. Over feat CNN collected characteristics of the lung nodules, and support vector machine methods were used to characterize the tumors. In fact, CAD-systems found the nodules. Analysis indicates that each approach can detect lesions with a sensitivity of 70%. Anthimopoulos et al. [9] Suggested that CNN identify and describe multiple lung tissues of respiratory problems. The CNN included five CNN layers, one layer of pooling, and three FC layers. The algorithm suggested was contrasted with many other CNN systems such as AlexNet and VGG Network. Results showed that the suggested CNN for tissue classification and identification was preferable to the other architectures. The new CNN obtained accuracy 85.61 %. Li et al. [10] suggested a CNN with just one convolution layer for patch detection on CT images with high resolutions. For this reason LIDC lung database is included. In addition, a pair of Support vector machine classifier with three feature extraction techniques their implemented CNN was contrasted to the mixture of a three-function extraction process and Support vector machine classification. Analysis indicates that CNN reached higher values of Sensitivity 0.88 and Precision 0.93 as compared to other methodologies. Shen et al. [11] suggested a multi-scale CNN to distinguish malignant and benign nodules in the lungs. The CNN classification methods have been used as support vector machine and random forest. In [12] a CNN is used to identify the lung CT scan images with an ILD dataset. CNN findings for the classification of good, ground glassy opacities, nano-nodules, reservoirs, of

\*Corresponding Author

## II. METHOD

ILDs were stated to be beneficial. In [13] used the updated edition of ResNet-18 as a classification system to label CT scanning images for the data science bowl and Kaggle lung. For certain of such fields, DNN may also attain near human intelligence [14]. CNN is the popular DNN model, integrates supervised learning methods that can help it to capture high-level features from unpasteurized images, and is promoted to handle lung nodule classifier accuracy improvements [15]. A new work adopts a deeper CNN with a one CNN layer for the classification of nodules and reveals a higher precision in comparison to the feature's traditional extraction methods. Owing to the customization of the CNN model [16], the transfer learning method often provides promising outcomes for the nodules classification [17]. Until the advent of deep learning, manual feature designing assisted by classifiers has been the basic nodule classification method. With the public availability of the LIDC-IDRI dataset [18], deep learning techniques [19] are already the dominant nodule classification system for research.

In order to obtain a suspicious-sensitive classification in CT images, we have to address at minimum two significant obstacles, the complexity of nodule depiction induced by a broad variety of nodule's morphology variants, and the problem raised by analytical models' radiological complexity to identify qualitative features as it is difficult to differentiate benign nodules from malignant nodules. So, to solve these challenges in this paper, we proposed a new transfer method (NTM) approach is used for lung nodule classification as shown in Fig. 2. First, volume of interest (VOI) extraction and feature extraction is used to extract the main features through CT image. Then for lung nodule classification task we proposed advance-fully convolution network (A-FCN).

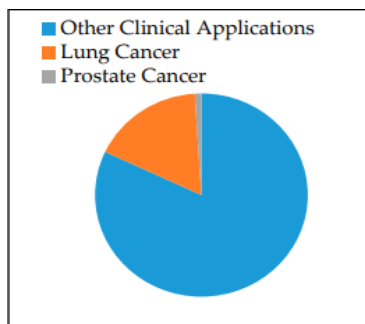


Fig. 1. Different CAD Applications.

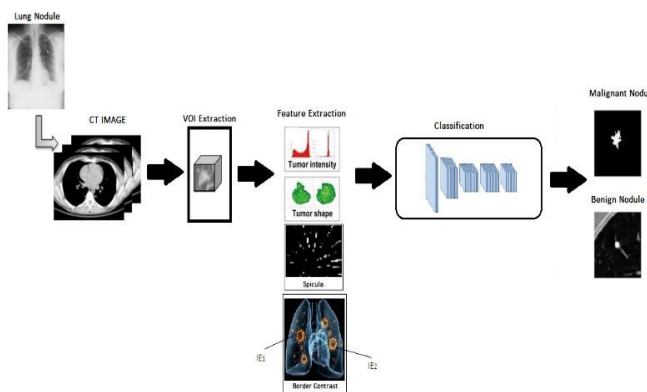


Fig. 2. Proposed New Transfer Method (NTM) Approach.

The proposed fully automated new transfer method (NTM) for lung nodule classification has two parts: (1) volume of interest (VOI) extraction and feature extraction; (2) advance-fully convolution network (A-FCN) which will automatically classified lung nodule between benign and malignant. In the proposed NTM technique the features are extracted by applying different VOI and feature extraction techniques then these features is being used as an input of advance-fully convolution network to classify the nodules among benign and malignant.

### A. Volume of Interest (VOI) Extraction

The physicist has indicated the location and size of the nodule to be examined using standard CT scan. Thus the VOI segmentation across the pulmonary nodule is conducted for examination upon CT images. The center position of a VOIs obtained from the standard CT images are configured manually whereas the CT images were tested for MPR images. Firstly, there is the transaxial's image with the highest nodule region, and manually determined the central coordinates. Therefore the diameter of the image is fixed to the maximum, and denoted by  $MD_{xy}$ , in Nodule path x-y for transaxial plane. Then, even as deforming including its slice in the path of the axis direction of body, the value of a slice wherein the lesion is still visible is being acquired and establish as  $MD_z$ . The VOI is now extract from the original image by using pixel on 3 sides, that is,  $2MD_{xy}$ ,  $2MD_{yx}$  and  $2MD_z$ .

### B. Feature Extraction

Feature extraction is the key of proposed model before the learning stage. One of the strategies for removal of dimensionality in image-processing is feature extraction of images. If the input image being examined is too complicated to process due to its repetitive features, therefore it is preferable to implement feature extraction technique. It effectively turns the immense collection of data into reduced range of feature. We extract CT image pixel intensity, shape, contrast of border and spicula for the lung nodule classification through CT image.

#### 1) Features Explanation

a) *CT Image Pixel Intensity and Shape of Nodule:* Most malignant lesions in CT-images include a larger intensity of a pixel. Standard uptake or silly useless value (SUV) [20] of slightly earlier CT images has therefore been described as ESUV and DSUV. In addition, the gap between the deferred and early stages in SUV was specified as part of the SUV. Two approaches are developed in the measurement of SUV, namely  $SUV_{max}$  and  $SUV_{peak}$ . The CT values in the nodule central and the average CT value within the nodule is determined throughout the CT images. As with the nodules' shape, malignant nodules frequently have a small ball type of form, whereas benign has a line type structure. An approach was suggested to use a Hessian-matrix to compare the ball type and line types [21]. It the Hessian matrix was done by adding the 3D image differentiation in 2nd order.

b) *Contrast of Nodule Border and Spicula:* Sometimes the location of a malignant nodule is uncertain. Consequently

the boundary comparison was measured using the discrepancy between the CT values of the nodules' exterior and interior borderlines. The estimated CT values only at numbers corresponding here to interior edge  $IE_1$  (CTIE<sub>1</sub>) and the exterior area  $IE_2$  (CTIE<sub>2</sub>) are extracted in order to quantify this value, and the distinction among the two values, is established as either the contrast which is shown in Fig. 3 The image is binaries to acquire  $IE_1$  and  $IE_2$ , and the shape was generated by the Sobel system. On the outlining the range of pixels has also been specified as  $IE_2$ . Consequently, morphological erosion with a systemic variable which includes a diameter of one pixel and the shape of the decreased area was derived in the same way as mentioned above; a collection of such pixels is utilized as  $IE_1$ . The appearance of spicula across the nodule raises the chances for malignancy in the nodule. We used Gabor filter used to detect spicula [22]. Using Gabor Filter allows imagining line shapes and the orientation as shown in Fig. 4. Finally we got detected spicula is by using Gabor filter as represent in Fig. 5 and the amount of radial items and the ranges were determined as spicula,  $S_a$  and  $S_b$  characteristics.

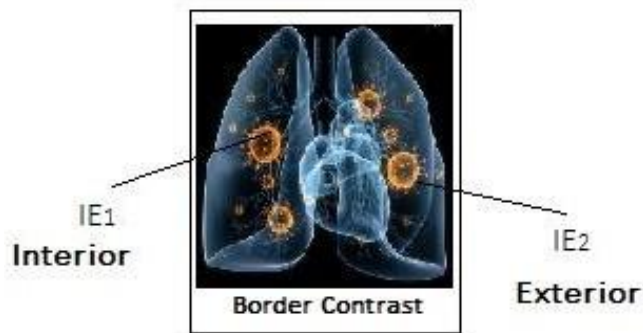


Fig. 3. Nodule Contrast.

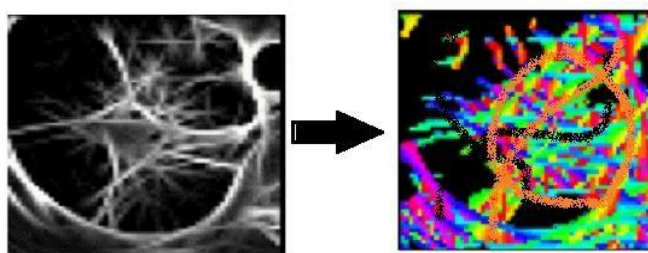


Fig. 4. (Left) Intensity, (Right) Angle using Gabor Filter.



Fig. 5. Detected Spicula.

### C. Advance-Fully Convolutional Network (A-FCN)

We proposed a new method to enhance the classification accuracy of lung nodule. Our A-FCN is composed of three separate kinds' layers which are convolution layer (CL), pooling layer (PL), and fully connected layer (Fc). On every input image CL executes convolution operators. As during training process these all layers can obtain features from the input data. Deep layers could identify high abstractions features. A pooling layer operates on independent feature channel and analyzes the surrounding values into just one. It therefore decreases the number of training parameter and significantly reduces the training time effectively. An fc layer connects every neuron in the existing layer from all neurons through last layer. Fc layers perform less reliably as compared to convolution layers owing to the lack of structural connection in images. It often expands the number of training parameters and thus, elongates time needed in training. Dropout which is a regulation strategy to decreasing amount of neurons and interactions is proposed to fix certain issue. As shown in Fig. 6 there are fifteen layers in which nine are convolution one. We employ three convolution blocks: 3 CL and 1 Pooling layer and max operator pooling layer. Every layer is accompanied by a ReLU activation function. Due to its efficient computing efficiency the last that have been fully connected are aimed at resolving a problem of classification from the extracting features. A ReLU and a softmax activation function are observed, respectively.

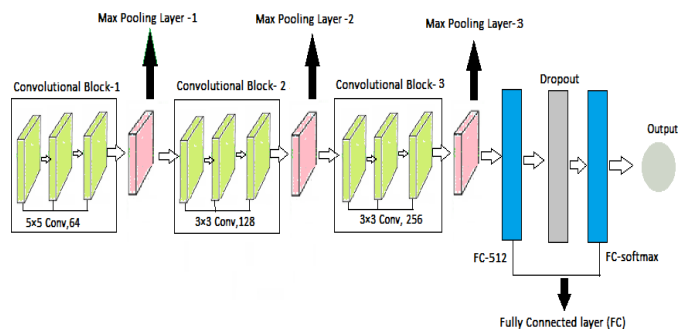


Fig. 6. Proposed Advance-Fully Convolutional Network for Classification of Lung Nodule.

## III. EXPERIMENTS

This section provides descriptions of how the proposed model is implemented and evaluated and it also discusses the experimental findings. Section 3.1 shows about dataset which were described for our experimentation and also training and testing phase. The analysis on the LIDC dataset is shown in Section 3.2. It explains the contrast between our model and several recent methods.

### A. Experimental Setting

1) *LIDC-IDRI Dataset*: The original edition of LIDC is 399 scans. Later, LIDC-IDRI has been extended to 1018 patient LIDC scans. In every CT scan, four professional radiologists used a LIDC-like labeling technique and an XML related file to report the data using a double-stage image classification process. Each radiologist interpreted and labeled every case separately at the first blind reading point. Growing radiologist studied its features and the attributes of any of the

three radiologists. Then they separately interpreted and recorded each case for the final decision. These two-phase marking will accurately classifies all pulmonary nodules, eliminating the need for pressured consensus. There were three types in the region of the nodule which are nodules  $\geq 3$  millimeter, nodules  $< 3$  millimeter, and non-nodules  $\geq 3$  millimeter. If there were more than three annotated nodule meaning by and over two doctors, the nodule would be labeled malignant. Otherwise nodule is considered benign in comparison. Around 198 malignant nodules and 153 benign nodules were available in our dataset. The nodules are ignored by the same votes. To order to reduce computational uncertainty, we separate the key transect for every voxel. We then apply the Data augmentation (DA) approach to expand the data by applying types to a dataset in terms of reducing deep learning over fitting. We dynamically flip and magnify the image through zooming 0.2 in general. The translation step is chosen using just a voxel of  $[-6, 6]$ , and a rotation angle of  $[90^\circ, 180^\circ \text{ and } 270^\circ]$  has been chosen at random. Eventually, 1958 malignant nodules, and 1867 benign nodules are available.

2) *Training and testing phase:* We employ Adam optimizer throughout the training process avoiding exponential decline from a learning rate=0.0001, and  $\beta_1=0.9$ ,  $\beta_2=0.99$  as a standard parameter. The k-cross validation approach is used. The procedure divides the data into the very same size k sections. The calculation is taken in k iterations; one component is used to test, whereas the rest is used for training. Iterations  $K = 10$  is a popular approach of such a validation approach.

### B. Experimental Result and Analysis

We evaluated the efficiency of our proposed new transfer method (NTM) model on the LIDC dataset and the ROC curve is also shown in Fig. 7. Our model achieved accuracy of 89.90% and area under curve (AUC) 0.9485. The efficiency of the suggest model is also contrasted with several other effective models in order to assess the effects of the appropriate methods, as seen in Table I. Liu et al. [24] presented a new hybrid convolution neural network in which LeNet as well as AlexNet were used and they combined both networks layers. Zhao et al. developed a multi-scale VGG-16 learning method to extract exclusion features from substitute stacked layers [25]. The latest 3D multi-scale convolution neural network framework for the classification of lung nodule has been created by Tafti et al. [26]. Yu Gu et al. [27] recommended the usage of a systemic method to identify a Deep convolutional neural network 3D multi-scalar predictor lung nodule. Moreover, it is suggested that an extra tiny nodule be observed using a classification approach with several scale cube clusters. Shen et al. [28] employed multi-scale to consider the complexity of a nodule, utilizing instead layered strata to isolate the discriminatory features.

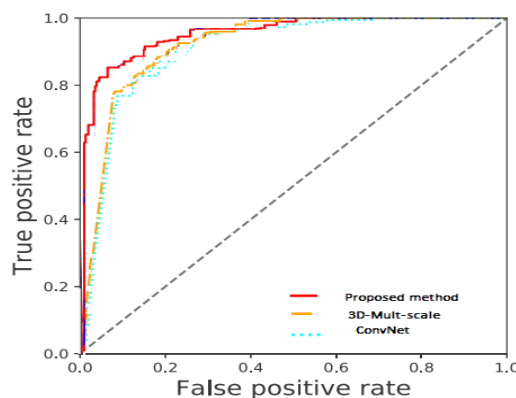


Fig. 7. ROC Curve of Proposed Method.

TABLE I. COMPARISON OF PROPOSED NTM WITH OTHER METHODS

| Research               | Methods  | Sensitivity | Specificity | Accuracy % | AUC    |
|------------------------|--|-------------|-------------|------------|--------|
| Liu et al. [24]        | CNN  | -           | -           | 82.23      | -      |
| Zhao et al. [25]       | ConvNet  | 0.843       | 0.858       | 84.97      | 0.902  |
| Tafti et al. [26]      | 3D-Multi scale CNN   | -           | -           | 83.75      | 0.926  |
| Yu Gu et al. [27]      | Multi scale 3D-DCNN  | 0.832       | 0.847       | 84.66      | -      |
| Shen et al. [28]       | Multi scale CNN  | -           | -           | 86.84      | -      |
| <b>Proposed method</b> | NTM: VOI and feature extraction with Advance-fully Convolutional Network | -           | -           | 89.90      | 0.9485 |

## IV. CONCLUSION

In this paper the aim is to introduce a new architecture named: new transfer method (NTM) for lung nodule classification. The network is divided into two parts: (1) VOI extraction and feature extraction; (2) advance-fully convolutional network (A-FCN). First features are extracted from lung CT image by using two extraction techniques in which we used intensity, shape, contrast of border and spicula extraction to extract the lung nodule. After finding nodule then it transfer to the A-FCN for lung nodule classification. Our A-FCN part used three convolutional blocks in which we have convolutional layer, pooling layer and fully connected layer. We also use ReLU and softmax activation function in our classification phase to solve the complexity of lung nodule classification. Finally, our system is thoroughly equipped for benign and malignant lung nodule classification. The result of LIDC-IDRI indicates that the NTM framework has improved accuracy.

REFERENCES

- [1] S. Lee, A. Kouzani, and E. J. Hu, "Hybrid Classification of Pulmonary Nodules" Communications in Computer and Information Science, vol. 51, pp. 472-481, 2009.
- [2] Ur Rehman MZ, Javaid M, Shah SI, Gilani SO, Jamil M, Butt SI (2018) An appraisal of nodules detection techniques for lung cancer in CT images. Biomedical Signal Processing and Control 1(41):140-151.
- [3] Silva D, Giovanni LF, Thales Levi AV, AristófanésCS ACP, Marcelo G (2018) Convolutional neural network-based PSO for lung nodule false positive reduction on CT images. Comput Methods Prog Biomed 162:109-118.
- [4] Skourt BA, El Hassani A, Majda A (2018) Lung CT image segmentation using deep neural networks. Procedia Computer Science 127:109-113.
- [5] Badura P, Pietka E (2014) Soft computing approach to 3D lung nodule segmentation in CT. Comput Biol Med 53:230-243.
- [6] Sone S, Takashima S, Li F, Yang Z, Honda T, Maruyama Y, et al. Mass screening for lung cancer with mobile spiral computed tomography scanner. Lancet. 1998;351(9111):1242-5.
- [7] National Lung Screening Trial Research Team, Aberle DR, Adams AM, Berg CD, Black WC, Clapp JD, et al. Reduced lung-cancer mortality with low-dose computed tomographic screening. N Engl J Med. 2011;365(5):395-409.
- [8] Ginneken, B.V.; Setio, A.A.A.; Jacobs, C.; Ciompi, F. Off-The-Shelf Convolutional Neural Network Features for Pulmonary Nodule Detection in Computed Tomography Scans. In Proceedings of the IEEE 12th International Symposium on Biomedical Imaging (ISBI), New York, NY, USA, 16-19 April 2015.
- [9] Anthimopoulos, M.; Christodoulidis, S.; Ebner, L.; Christe, A.; Mougiakakou, S. Lung Pattern Classification for Interstitial Lung Diseases Using a Deep Convolutional Neural Network. IEEE Trans.Med. Imaging 2016, 35, 1207-1216.
- [10] Li, Q.; Cai, W.; Wang, X.; Zhou, Y.; Feng, D.D.; Chen, M. Medical Image Classification with Convolutional Neural Network. In Proceedings of the 13th International Conference on Control, Automation, Robotics & Vision, Marina Bay Sands, Singapore, 10-12 December 2014.
- [11] Shen, W.; Zhou, M.; Yang, F.; Yang, C.; Tian, J. Multi-scale Convolutional Neural Networks for Lung Nodule Classification. Inf. Process. Med. Imaging 2015, 24, 588-599.
- [12] Bondfale, N.; Banait, S. Lung Pattern Classification for Interstitial Lung Diseases Using a Deep Convolutional Neural Network. IJARCC 2017, 5, 9851-9856.
- [13] Data Science Bowl 2017. Kaggle. Available online: <https://www.kaggle.com/c/data-science-bowl-2017> (accessed on 5 February 2019).
- [14] Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., van den Driessche, G., et al. (2016). Mastering the game of Go with deep neural networks and tree search. Nature, 529, 484-489.
- [15] Szegegy, C., Vanhoucke, V., Ioffe, S., Shlens, J., & Wojna, Z. (2015b). Rethinking the inception architecture for computer vision. CoRR.
- [16] Tajbakhsh, N., Shin, J. Y., Gurudu, S. R., Hurst, R. T., Kendall, C. B., Gotway, M. B., et al. (2016). Convolutional neural networks for medical image analysis: Full training or fine tuning? IEEE Transactions on Medical Imaging, 35(5), 1299-1312.
- [17] Zhao, X., Liu, L., Qi, S., Teng, Y., Li, J., & Qian, W. (2018). Agile convolutional neural network for pulmonary nodule classification using CT images. International Journal of Computer Assisted Radiology and Surgery, 13(4), 585-595.
- [18] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In MICCAI, 2015.
- [19] S. G. Armato et al. The lung image database consortium (lidc) and image database resource initiative (idri): a completed reference database of lung nodules on ct scans. Medical physics, 38(2):915-931, 2011.
- [20] Keyes JW Jr. SUV: standard uptake or silly useless value? J Nucl Med. 1995;36(10):1836-9.
- [21] LeCun Y, Bengio Y, Hinton G. Deep learning. Nature. 2015;521(7553):436-44.
- [22] Rangayyan RM, Ayres FJ. Gabor filter and phase portraits for the detection of architectural distortion in mammograms. Med Biol Eng Comput. 2006;44(10):883-94.
- [23] Global Computer-Aided Detection (CAD) Market US\$ 2.2 Billion by 2023.
- [24] X., Liu, L., Qi, S., Teng, Y., Li, J., & Qian, W. (2018). Agile convolutional neural network for pulmonary nodule classification using CT images. International Journal of Computer Assisted Radiology and Surgery, 13(4), 585-595.
- [25] Zhao, D., Zhu, D., Lu, J., Luo, Y., & Zhang, G. (2018). Synthetic medical images using F&BGAN for improved lung nodules classification by multi-scale VGG16. Symmetry, 10(10), 519.
- [26] Tafti, A. P., Bashiri, F. S., LaRose, E., & Peissig, P. (2018, June). Diagnostic classification of lung CT images using deep 3d multi-scale convolutional neural network. In 2018 IEEE International Conference on Healthcare Informatics (ICHI) (pp. 412-414). IEEE.
- [27] Gu, Y., Lu, X., Yang, L., Zhang, B., Yu, D., Zhao, Y., ... & Zhou, T. (2018). Automatic lung nodule detection using a 3D deep convolutional neural network combined with a multi-scale prediction strategy in chest CTs. Computers in biology and medicine, 103, 220-231.
- [28] Shen, W., Zhou, M., Yang, F., Yang, C., & Tian, J. (2015, June). Multi-scale convolutional neural networks for lung nodule classification. In International Conference on Information Processing in Medical Imaging (pp. 588-599). Springer, Cham.



# 6G: Envisioning the Key Technologies, Applications and Challenges

Syed Agha Hassnain Mohsan<sup>1</sup>, Alireza Mazinani<sup>2</sup>, Warda Malik<sup>3</sup>, Imran Younas<sup>4</sup>  
Nawaf Qasem Hamood Othman<sup>5</sup>, Hussain Amjad<sup>6</sup>, Arfan Mahmood<sup>7</sup>

Department of Electrical Engineering, COMSATS University Islamabad (CUI), Islamabad, Pakistan<sup>1</sup>

School of Electronic and Information Engineering, Beihang University, Beijing, China<sup>2</sup>

Department of Electrical Engineering, COMSATS University Islamabad (CUI), Wah, Pakistan<sup>3</sup>

Department of Electronics Engineering, Xi'an Jiaotong University (XJTU), Xi'an, China<sup>4</sup>

Department of Information and Communication Engineering, Xi'an Jiaotong University (XJTU), Xi'an, China<sup>5</sup>

Department of Marine Science and Information Technology, Ocean College, Zhejiang University, Zhoushan, China<sup>6</sup>

Complex Networks and Control Lab, Shanghai Jiao Tong University, Shanghai, China<sup>7</sup>

**Abstract**—In 2030, 6G is going to bring remarkable revolution in communication technologies as it will enable Internet of Everything. Still many countries are working over 5G and B5G has yet to be developed, while some research groups have already initiated projects on 6G. 6G will provide high and sophisticated QoS e.g. virtual reality and holographic communication. At this stage, it is impossible to speculate every detail of 6G and which key technologies will mark 6G. The wide applications of ICT, such as IoT, AI, blockchain technology, XR (Extended Reality) and VR (Virtual Reality), has created the emergence of 6G technology. On the basis of 5G technique, 6G will put profound impact over ubiquitous connectivity, holographic connectivity, deep connectivity and intelligent connectivity. Notably, research fraternity should focus on challenges and issues of 6G. They need to explore various alternatives to meet desired parameters of 6G. Thus, there are many potential challenges to be envisioned. This review study outlines some future challenges and issues which can hamper deployment of 6G. We subsequently define key potential features of 6G to provide the state of the art of 6G technology for future research. We have provided a review of extant research on 6G. In this review, technology prospects, challenges, key areas and related issues are briefly discussed. In addition, we have provided technologies breakdown and framework of 6G. We have shed light over future directions, applications and practical considerations of 6G to help researchers for possible breakthroughs. Our aim is to aggregate the efforts and eliminate the technical uncertainties towards breakthrough innovations for 6G.

**Keywords**—IoT; AI; communication technologies; holographic communication; blockchain

## I. INTRODUCTION

Although the era of 5G is not fully developed, the limitations of 5G have created the demand for 6G networks. In 2019, communication synergy around the globe drafted first 6G white paper in world's first 6G summit in Finland. After that, many government organizations and research group from prestigious institutes started introducing their 6G projects. UK government has decided to invest in 6G technology [1], while Academy of Finland has launched "6 Genesis" project.

What is 6G technology? Some people expect more than just a faster version of 5G. For example, there should be no limitation of coverage to ground level. Instead, it must provide undersea and space coverage. It must enable higher Artificial Intelligence (AI) characteristics. In fact, some researchers consider it as an "AI-empowered" network [2]. It should not merely involve AI but it must integrate AI networking functions and tool. In addition, secrecy, privacy and risk mitigation must be a core component of its architecture [3]. In this review, we have investigated privacy and security challenges along with potential applications of 6G network. An overview of different dimensions of 6G networks is shown in Fig. 1.

After commercialization of 5G network, academia and industrial experts have started thinking about next 6G network, services and requirements behind it. If we look at standardization methods of 5G technology, three aspects were investigated as, ultra-reliable and low latency communications (URLLC), massive machine type communications (mMTC) and enhanced mobile broadband (eMBB). Although such scenarios are not fully investigated for 6G networks, however some pioneering works [4-5] forecast the idea to link everything via unlimited, reliable and instantaneous wireless resources. We have shown an overview of 6G coverage in Fig. 2.

To bring this revolution to connect everything worldwide, 6G will require extreme communication techniques such as smart living based wireless brain-computer interactions [6], smart working based on seamless holographic projection [7] and smart design considering real-time digital twins [8]. The evolution from 5G to 6G is summarized in Table I.

We have provided some performance metrics for 6G networks below and compared with conventional 5G requirements.

- **Mobility:** The highest speed to be achieved will be increased from 500 km/h to 1000 km/h.
- **Reliability:** 99.99% reliability will be achieved to support unmanned vehicles including AUVs and collaborative robotics.

- **Latency:** The communication latency will be decreased by 10 times for end-to-end point of view.
- **Throughput:** A maximum throughput of 1 Tb/s will be needed for 6G which is 1000 times speedy than 5G. 100 times advancement is expected.
- **Energy and Spectrum Efficiency:** 100 times energy efficiency and 10 times spectrum efficiency will be achieved.

The above described metrics involve disruptive features in 6G networks to use more flexible frame structure, more frequency bands and more spatial dimensions. Many industrial experts and technologies have discussed to meet these requirements. Such as, Space-Air-Ground integrated network [9] have suggested to enhance the spatial degrees of freedom by incorporating airborne, terrestrial and satellite networks, which extend 2D into 3D space for reliable and efficient connectivity [10]. Under-utilized high frequency bands can be explored through Terahertz (THz). Visible light communication (VLC) is a promising candidate for tens of GHz bandwidth [11] and 1 Tb/s throughput. Meanwhile, AI driven communication [12] with intelligent control will be possible.

TABLE I. EVOLUTION FROM 5G TO 6G

| Key parameter                                | 5G              | 6G              |
|--|-----------------|-----------------|
| Mobility (km/h)                              | 350-500         | 1000            |
| Peak spectral efficiency (b/s/Hz)            | 30              | 60              |
| End-to-end latency (ms)                      | 1               | 0.1             |
| Reliability                                  | 10-5            | 10-9            |
| Connection Density (device/km <sup>2</sup> ) | 10 <sup>6</sup> | 10 <sup>7</sup> |
| Area traffic capacity (Mbps/m <sup>2</sup> ) | 10              | 1000            |
| Channel bandwidth (GHz)                      | 1               | 100             |
| Spectral efficiency (b/s/Hz)                 | 0.3             | 3               |
| Energy Efficiency (Tb/J)                     | NA              | 1               |
| User Data rate (Gbps)                        | 1Gbps           | >10Gbps         |
| Peak data rate                               | 10-20Gbps       | >100Gbps        |
| Receiver sensitivity                         | -120dBm         | <-130dBm        |
| Position precision                           | m               | cm              |
| Coverage                                     | 70%             | >99%            |
| Delay  | ms              | <ms             |

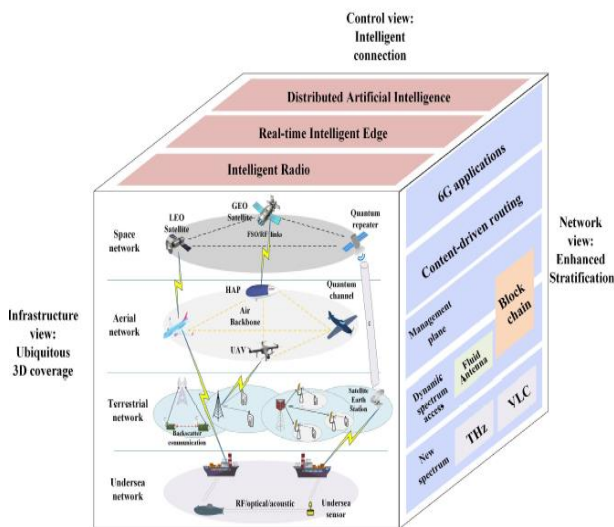


Fig. 1. Different Dimensions of 6G Architecture [17].

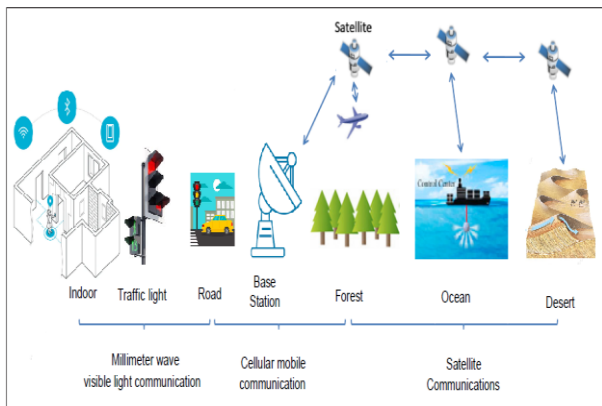


Fig. 2. An Overview 6G Network Coverage.

## II. HISTORICAL OVERVIEW

### A. 1G and 2G - 10 Times Reduction

1G and 2G networks provide the basic service of voice calling. Significant contribution has been made from 1G to 2G realization, such as China Mobile’s annual report revealed 10 times price depletion from 0.1 to 0.01 US dollar/minute [13]. In addition, world’s population using these services also increased from 10% (1G) to above 50% (2G) within 20 years [14].

### B. 3G and 4G - 1000 Times Reduction

3G and 4G networks provide the key service of data transmission. Technical development from 3G to 4G include orthogonal frequency division multiplexing (OFDM) and multiple-input multiple-output (MIMO) and user-sensitive goal of 1000 times price reduction. Initial 3G users are limited to business community to access company resources and emails, while further enhancement occurs only after the deployment of 4G networks.

### C. 5G and 6G - 1000 Times Reduction

An explosive growth of 5g and beyond is found to facilitate human-to-machine and machine-to-machine communications. Although the existing 5G is still based on eMBB with the similar price strategy of 4G networks. However, it will be more reasonable to charge on the basis of connection rather than data traffic. According to FTTH systems, China is charging 100-200 US dollars for each terminal [15]. However, 100 trillion sensors are expected to be manufactured and connect to internet by the end of 2030 to revolutionize 6G. Hence, 1000 times price reduction will be required to develop a sustainable smart society. Table II summarizes different features of 5G and 6G.

A details comparison of 1G to 6G [16] technologies is summarized in Table III.

TABLE II. COMPARISON BETWEEN 5G AND 6G

| Feature     | 5G     | 6G      |
|-------------|--------|---------|
| VLC         | No     | Yes     |
| Reliability | Good   | Extreme |
| AI          | No     | Yes     |
| Centre      | User   | Service |
| Capacity    | 1D /2D | 3D      |
| WPT         | No     | Yes     |
| Core        | IoT    | IoE     |
| Privacy     | Good   | Extreme |
| Real Time   | No     | Yes     |

TABLE III. COMPARISON OF 1G TO 6G TECHNOLOGIES

| Feature           | 1G        | 2G           | 3G                    | 4G                     | 5G                 | 6G                         |
|-------------------|-----------|--------------|-----------------------|------------------------|--------------------|----------------------------|
| Time span         | 1980-1990 | 1990-2000    | 2000-2010             | 2010-2020              | 2020-2030          | 2030-2040                  |
| Highlight         | Mobility  | Digitization | Internet connectivity | Real-time applications | Extreme data rates | Privacy, secrecy, security |
| Core network      | PSTN      | PSTN         | Packet N/W            | Internet               | IoT                | IoE                        |
| Services          | Voice     | Text         | Picture               | Video                  | 3D VR/AR           | Tactile                    |
| Architecture      | SISO      | SISO         | SISO                  | MIMO                   | Massive MIMO       | Intelligent Surface        |
| Multiple xing     | FDMA      | FDMA, TDMA   | CDMA                  | OFDMA                  | OFDMA              | Smart OFDMA plus IM        |
| Maximum Frequency | 894 MHz   | 1900 MHz     | 2100 MHz              | 6 GHz                  | 90 GHz             | 10 THz                     |
| Maximum Data rate | 2.4 kb/s  | 144 kb/s     | 2 Mb/s                | 1 Gb/s                 | 35.46 Gb/s         | 100 Gb/s                   |

### III. CURRENT RESEARCH PROGRESSES TOWARDS 6G

Many research groups have shown the vision of 6G and research fraternity has started advance research activities and projects [18-20]. There is a growing inclination in research publications in this domain. Recently, Yang Lu et al. [21] filtered extant articles about 6G as various institutes have been conducting research on several approaches towards 6G. Publishing trend between 2016 and 2020 is depicted in Fig. 3. X-axis shows the number of publications while Y-axis shows specific year. It can be seen that maximum papers were published in IEEE conferences and journals.

E. Basar et al. [22] have discussed MIMO paradigm for 6G. They focused on research activities related to device manufacturing capabilities. S.M. Bohloul et al. [23] have made a good discussion about trends, opportunities and developments in 6G. They have outlined communication technologies e.g. tactile internet, flying networks and holographic calls for future networks in 2030. In [24] and [25], future trends and applications enabling 6G technology have

been summarized. Blockchain technology, human centric services and key performance indicators of 6g are investigated in these studies. 6G prospect, challenges and key performance indicators are defined. Authors have illustrated the role of OWC [26] in 6G technology. Some recent articles have provided detailed discussions about green 6G network architecture [27], 6G spectrum management [28], security challenges [29], potential solutions [33], machine learning technologies for 6G [30-31] and performance evolution of terahertz [32] communications. Some publications have discussed data center connectivity [34] and practical implementation of multiple access [35] for 6G networks. Network patterns for 6G are highlighted in some studies [36-37]. 6G based AI applications [38-39] which will unlock the full potential of radio signals are outlined in some studies. Hardware foundation of AI [40] is proposed in an article. Zhao et al. [41] have provided a survey on intelligent reflective surfaces for 6G networks. These promising materials can enhance the spectral efficiency [42] in 6G networks. In addition, several countries have started research projects to initiate, develop, define and reshape framework of 6G networks. Table IV summarizes country wise research initiatives in 6G networks.

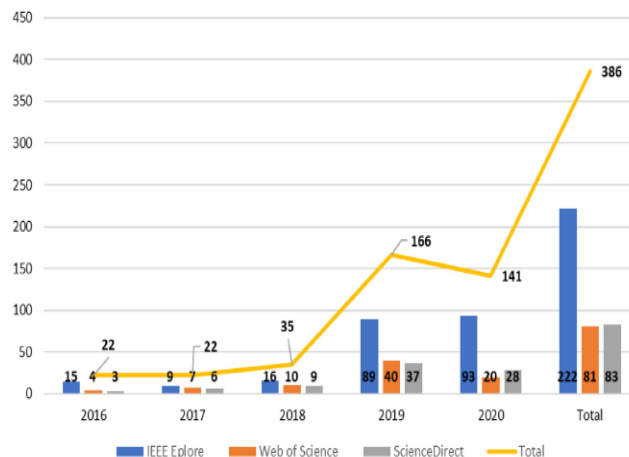


Fig. 3. A Trend of Publications on 6G [21].

TABLE IV. 6G PROJECTS IN DIFFERENT COUNTRIES

| Country   | Year         | Research Initiative  |
|-----------|--------------|--|
| 2018      | Finland      | 6G initiative was launched in University of Oulu.  |
| 2019      | China        | 37 research institutes have started focusing on 6G research.   |
| 2019      | USA          | Spectrum between 95 GHz and 3 THz has been opened.   |
| 2019      | South Korea  | KAIST and LE Electronics have established a 6G research center with collaboration.   |
| 2020      | Japan        | Sony, Intel and NTT have collaborated to work on 6G technology. Japan has planned to spend \$US 2 billion on 6G industrial research. |
| 2020      | Saudi Arabia | Researchers from KAUST have started working on 6G technology.  |
| 2021-2026 | South Korea  | Government of Korea will invest \$169 million to secure 6G and planning to launch 6G pilot project in 2026.                          |

#### IV. TECHNOLOGY BREAKDOWN

We have discussed each generation in the aspects of frequency, spatial and time domains as given below. Technology breakdown from 1G to 6G is also displayed in Fig. 4.

##### A. Spatial - 10 Times

The purview of the Space-Air-Ground integrated network enfold an extensive range of terminals, satellite communications, flying drones, which proffers two times cost reduction with low number of base stations. Ultra-scale MIMO can improve 50% throughput without extra costs; thereby 1.5 times cost reduction can be achieved. Intelligent adaptation of beam eventually brings three to four times reduction, while 10 times reduction is possible through different network architectures.

##### B. Frequency - 10 Times

In frequency domain, the cost reduction is dependent on utilization of low cost spectrum. Although mmWave, VLC and THz are capable to offer significant bandwidth for wireless transfer, the befitting scenario is indoor users with pedestrian mobility, which is 70% of the overall traffics. Thus, higher frequency bands can facilitate with 3 times reduction. Moreover, another 3-4 times reduction is possible by flexible usage of multiple frequency bands.

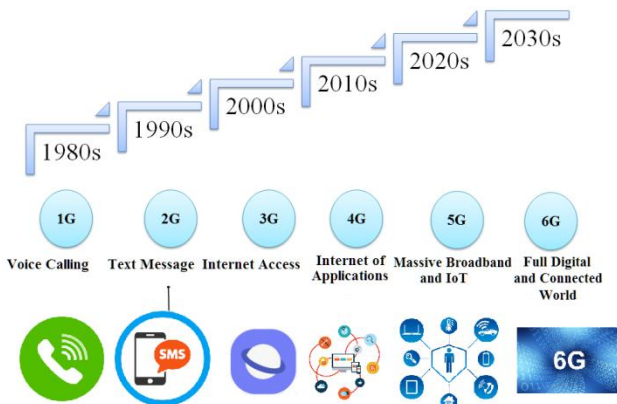


Fig. 4. An Overview of 1G-6G Devices and the Corresponding Technology Breakdown.

##### C. Time - 10 Times

Another prominent alternative is to profoundly impact the resolution of time-frequency resource to feature flexible frame structure and integrate modulation scheme like index modulation. A fast mode adaptation can enhance the performance with a massive combination of duplex schemes, modulation techniques and frame structures. By incorporating several techniques, we expect 1000 times reduction can be achieved. The core element is AI-assisted intelligent communication which can reduce cost up to 20-50 times.

#### V. 6G REQUIRES A NEW PARADIGM

Next generation 6G network requires wide bandwidth for high resolution and high carrier frequencies for small antennas. A potential issue is to analyze and process radio systems over wide bandwidth without prior information of signal,

modulation and carrier frequency. An idea option is photonics defined system as it can provide high spectrum capacity with extreme bandwidth. It is an extended version of microwave photonics through coherent optics, optical computing and photonics DSP. A paradigm shift and hyper-S curve [43] presenting a revolution of mobile of communication technologies is shown in Fig. 5.

Open loop control, reduced feedbacks, software defined systems and interference cancellation have developed this system. A radical innovation is expected in case of 6G which will result into a new S curve. The logical start of 6G is shown in Fig. 6.

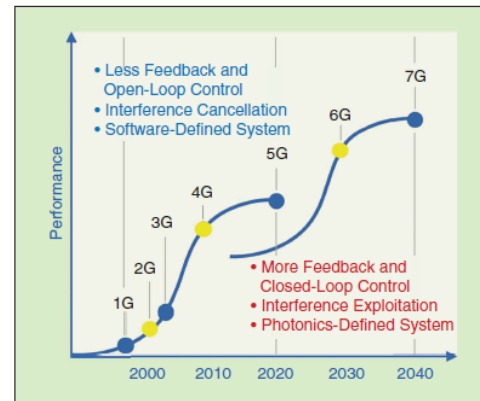


Fig. 5. Hyper-S Curve and Paradigm Shift [43].



Fig. 6. New Logical Start of 6G.

#### VI. KEY AREAS IN 6G NETWORKS

We will discuss key areas in 6G networks and we have also investigated privacy and security issues in these areas.

##### A. Real-Time Intelligent Edge

It is not fully possible to implement Unmanned Aerial Vehicle (UAV) networks with existing technologies as it needs real time intelligence and extremely low latency to control the network. Although 5G technology has supported autonomous driving, however prediction, self-adaption and self-awareness for network entities is not supported [44]. Thus, a new technology is required to overcome these issues. It will be possible through 6G technology to enable AI-powered services. As AI will be incorporated in vehicle networks, it will support several security mechanisms. However, it will cause new privacy and security issues. Tang et al. [45]

investigated that both network and physical environments should be considered for a vehicle network as it can reduce malicious activities.

**B. Distributed AI**

6G networks will support Internet of Everything (IoE). It will make 6G network advance enough to take intelligent decisions [27]. In addition, IoT needs to support various requirements. 1) The edge device must compute and store data. 2) It should have the capability to clean and abstract data [46]. This approach can improve the privacy and security of the network. Machine learning algorithms can be integrated with 6G to ensure security [47] and data integrity.

**C. 3D Intercoms**

In 6G network, network optimization and designing will move from 2D to 3D [48]. 6G technologies will be capable of supporting 3D communication to enable undersea, UAVs and satellite communication. A 3D intercom can facilitate this feature with accurate time and location. In addition, resource management, routing and mobility characteristics also require network optimization in 3D intercom. Currently, THz bands are being experienced. With this band, some new technologies e.g. quantum and molecular communications can be applied for remote communication [49]. Wei et al. [50] highlighted some security risks for authentication process. In addition, performance of 6G networks in undersea environment is still unpredicted. Once 6G network operations in undersea environment are possible, more opportunities and challenges will emerge in near future. Fig. 7 illustrates some application scenarios supported by 6G technologies.

**D. Intelligent Radio**

The transceiver devices can be separated in 6G while they were designed together in earlier generations. Hence, it has the capability to update itself. Some operating systems are developed on the basis of hardware information and AI technology. Researchers have investigated signal jamming and suspicious activities in data transmission. Thus, 6G will enable intelligent and secure data transmission.

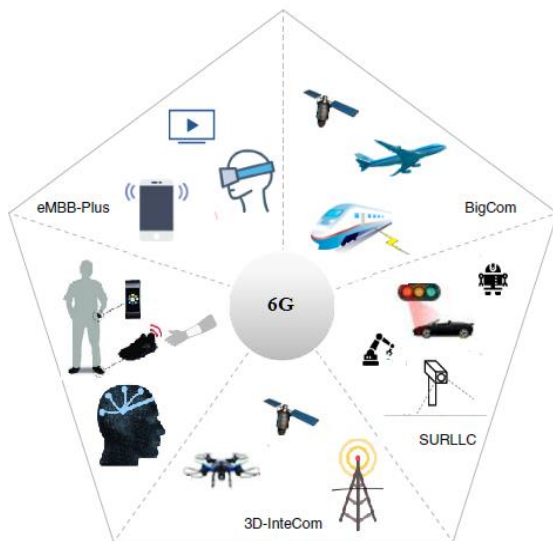


Fig. 7. Applications Supported by 6G.

**VII. 6G TECHNOLOGIES**

In this section, we have discussed 6G technologies and associated privacy and security concerns. Table V presents an overview of 6G technologies and security issues. While Fig. 8 illustrates potential key technologies of 6G networks.

TABLE V. EVOLUTION FROM 5G TO 6G

| Technology                   | Reference | Privacy and security issue |
|------------------------------|-----------|----------------------------|
| AI                           | [48]      | Malicious attack           |
| AI                           | [51]      | Communication              |
| AI and quantum communication | [52]      | Encryption                 |
| Blockchain                   | [53]      | Communication              |
| Blockchain                   | [54]      | Access control             |
| Blockchain                   | [55]      | Authentication             |
| VLC                          | [56]      | Malicious attack           |
| VLC                          | [57]      | Communication              |
| THz                          | [58]      | Malicious attack           |
| THz                          | [59]      | Authentication             |
| Quantum communication        | [60]      | Encryption                 |

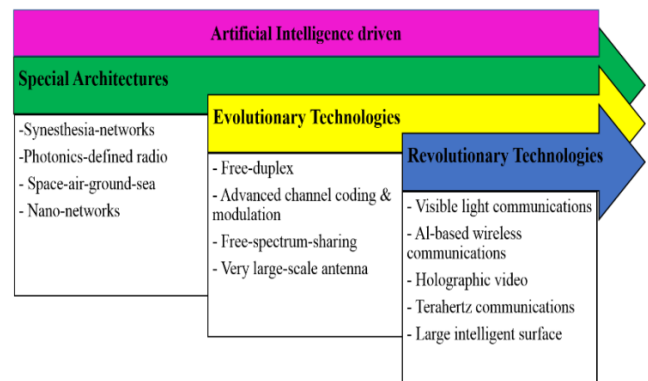


Fig. 8. Key Technologies of 6G Network.

**A. AI**

AI plays the most important role in future network infrastructures as compared to other technologies. AI has gained a lot of attention from various research groups. With this growing research, various privacy and security problems are also emerging. Although AI is also integrated in 5G technology, however it is considered as the key component of future 6G. AI technologies are subdivided into physical layer consisting of network infrastructure, architecture layer, computing layer which contains software defined networks, edge/cloud computing and network function virtualization.

**B. Quantum Communication**

Another promising technology in 6G network is quantum communication. It can significantly increase reliability and security of data transmission. Quantum state is affected with any adverse eavesdrop. Quantum communication offers security with essential breakthroughs. It can provide solutions and elevate communication which is not possible to achieve

through traditional communication techniques [61]. However, it is not the only panacea for all security threats. Although research has been carried out to develop quantum cryptography, but fiber attenuation is a serve issue in long distance quantum communication. Zhang et al. [62] and Nawaz et al. [52] have presented quantum mechanism for secure communication through quantum key distribution models.

### C. Blockchain

Another prominent technology is 6G network is blockchain. It has several used such as spectrum sharing, distributed ledger technology and network decentralization. S. Dang et al. [48] used network decentralization to enhance network performance. Strinati et al. [63] also increased authentication security through distributed ledger technology. Blockchain technology can also overcome spectrum monopoly and low spectrum utilization [64]. Blockchain privacy concerns are related to communication, authentication and access control. X. Ling et al. [65] have illustrated authentication and secure network access features through blockchain technology.

### D. Visible Light Communication (VLC)

VLC is a promising technology to meet the rapidly growing needs of wireless connectivity [66]. VLC has been deployed in vehicular Ad Hoc networks and indoor positioning systems. J. Luo et al. [67] have presented an indoor positioning system based on VLC. It is noticed that VLC limits EM interference. Some research studies have demonstrated high speed data transmission by using LEDs. Some deficiencies exist which affect the performance of VLC communication. In particular, VLC technology mainly supports indoor scenario as it is severely affected by natural light. The security issues of VLC technology include communication problem and malicious activities. A SecVLC protocol [57] is developed for secure data transmission in a vehicular network. Fig. 9 presents an overview of OWC in 6G technologies. We have provided a detailed discussion of OWC and 6G our recent systematic study [68]. 6G is expected revolution in UWPT [69] and UWOC.

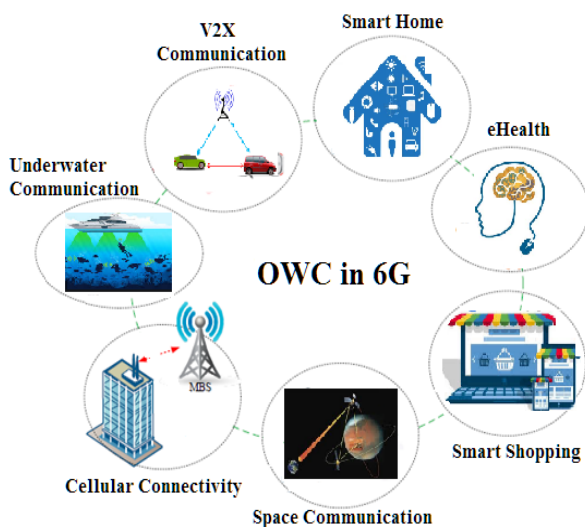


Fig. 9. OWC and 6G.

### E. Terahertz Technology (THz)

Existing RF band cannot be utilized for future 6G technologies [70]. It has spurred the demand for THz technology. THz communication technology used 0.1-10 THz band. Moreover, it exploits optical signals and EM waves. Huang et al. [27] have highlighted several benefits of THz band including 100 Gbps data rate, high security and limited eavesdropping. THz can significantly minimize intercell impact [64]. Strianti et al. [63] have investigated energy consumption problem in THz communication. THz faces security risks of authentication and malicious attack. We have summarized a comparison between VLC and THz communication in Table VI.

TABLE VI. COMPARISON BETWEEN VLC AND THZ

| Feature                 | VLC        | THz       |
|-------------------------|------------|-----------|
| Cost                    | Cheap      | Expensive |
| Data rate               | 10 Gbps    | 100 Gbps  |
| EM radiation            | no         | yes       |
| Transmission            | LOS        | NLOS      |
| Transmission power      | Low        | High      |
| Spectrum regulation     | Unlicensed | Licensed  |
| Inter cell interference | No         | Yes       |
| Bandwidth               | 10-100 GHz | 100 THz   |

## VIII. POTENTIAL CHALLENGES

There are several critical challenges which can affect future 6G technology. In this section, we have discussed big data, power, latency and hardware design challenges.

### A. Wireless Big Data

AI technology has proven its great stature in computer vision tasks. It has potential application in ImageNet big data sets. Such supervised learning method can solve complex optimization challenges in wireless communication. However, there exist many serious concerns for developing public wireless data sets for research purpose. As big data is processed and stored through cloud computing. The DIOE will cause new challenges to manage this data.

### B. Portable and Low-Latency Algorithm

The current AI technologies are developed to meet certain requirements; however, it has limited migration capability. However, an important performance metric is to design portable and low latency algorithms. In addition, latency trade-off and accuracy is highly required as compared to than traditional computer vision tasks.

### C. Hardware Co-Design

High density parallel computing methods are required in AI-assisted technologies. Wireless network architecture requires certain parameters to support AI-assisted communication. Moreover, computer performance can face degradation in case of advance materials e.g. graphene transistors and high temperature superconductors.

#### D. Power Supply

6G technology can make an efficient connection between mobile devices. Energy-efficient algorithms and strategies must be adopted in such cases. 6G will introduce new power control mechanism such as advance wireless power transfer (WPT) for smart devices. It will enable energy harvesting and optimization technique for efficient performance in harsh environment such as undersea environment.

#### E. Network Security Issue

Researchers need to focus on privacy concerns in future 6G technology. They must investigate new security approaches for secure data transmission. A significant extension in 5G security methods can also enable 6G security. Researchers can find new techniques to efficiently integrate THz with mmWaves. It can put profound impact on 6G privacy and security mechanism.

### IX. 6G POTENTIAL APPLICATIONS

Every new epoch of network technology introduces new services. In this section, we have outlined some potential applications for future 6G technology.

#### A. Multi-Sensory XR Applications

The low latency and high bandwidth of 5G technology has extended the VR/AR experience for 5G users. Nonetheless, some existing challenges should be removed in 6G network. The VR/AR experience will be enhanced in 6G network. Multiple sensors can be allocated to gather sensory data. Hence, the XR in 6G network will be formulated from URLLC and eMBB. The security concerns of eMBB and URLLC include internal communication, access control and malicious attack. Chen et al. [71] have investigated security problems in URLLC applications. J.M. Hamamreh et al. [72] have suggested a technique to improve security against URLLC attacks. Similarly, Yamakami et al. [73] have proposed a 3D model for secrecy risks in XR applications.

#### B. Connected Robotics and Autonomous Systems

Another promising application of 6G technology is the connected robotics and autonomous systems. A comprehensive autonomous system is required in 6G network as compared to 5G. This system should be based on a multi-dimensional network. In addition, the system must be capable to embed AI across the network. This feature will support automatic controlling of internal components. Strianti et al. [63] have envisioned resource control, caching and automatic handling in network. They developed an automated factory which contains cloud services, database and UAV networks to make it a complete autonomous system. 6G will be helpful for underwater robotic tasks such as security, imaging and rescue. 6G will enable efficient surveillance, navigation and robotic communication. It will develop a reliable, secure and smooth communication channel for real-time applications. Low latency and high speed data transmission of 6G will be helpful to obtain video data.

#### C. Wireless Brain-Computer Interactions

The concept of wireless BCI is to develop a link between device and human brain. This device can be placed inside or outside the body. The key application of wireless BCI is to

control auxiliary equipment for disabled people. It is expected that BCI will have more applications with involving 6G technology. In 2015, Chen et al. [74] developed a brain-computer interface to speed up spelling. The security risks of wireless BCI contain encryption and malicious behavior. Several research studies [75-76] have discussed security issues, protection techniques and hacking applications to mitigate security issues.

#### D. Accurate Indoor Positioning

With evolving GPS, outdoor positioning systems have been developed accurately. However, indoor position systems need research attention to cope up with complicated indoor EM propagation. New aspects of full-fledge applications are expected with reliable and accurate indoor positioning services. However, alone RF communication cannot achieve accurate indoor positioning. Such crucial application can only be realized with 6G technology.

#### E. Holographic Communications

6G will make it possible to realize virtual in-person meeting than traditional video conferencing. It can be achieved through a realistic projection of real time mobility in short time. It is not sufficient to transmit 3D image with voice to realize in-person presence. However, it requires a stereo audio incorporated in 3D video. We can state that user interacts with holographic data and can carry out possible modifications as needed. This scenario can be captured by reliable communication networks with extremely large bandwidth.

#### F. Tactile Communications

After realizing holographic communication for virtual in-person meeting, it is advantageous to carry out tactile communication to transfer the physical interaction remotely. Specifically, it includes interpersonal communication, cooperative automated driving and teleoperation. Stringent demands or these applications can be met through reliable cross-layer communication-system. Moreover, delay can be mitigated by carefully handling handover, scheduling, queuing and buffering.

#### G. Internet of Nano-Things

Nanotechnology is providing remarkable opportunities to design advance materials. It has developed nanodevices like nanosensors. It has the capability to perform simple task and enable internet connectivity. IoNT [77] is developed by integrating nanotechnology with IoT. It has the ability to sense and transmit information. IoNT can be deployed with allied technologies such as big data, cloud computing, WSN, UWSN [78-79] and IoT. However, IoNT faces limited memory space issue for real-time implementation as data storage depends on the size of nano memory. Another potential issue is high biological noise and congestion control in nanodevices. An overview of 6G applications is summarized in figure 10.

#### H. Intelligent Internet of Medical Things (IIoMT)

IIoMT will remove space and time hurdles to perform surgical operations. 6G will provide high speed communication for efficient performance of telesurgery beyond boundaries. IIoMT will make use of holographic communication, tactile communication and AR/VR to assist remote doctors. Thus, 6G

technology will ensure intelligent healthcare. It is expected to bring mobile hospital technology which can remove ambulance services. In future, it can offer medical devices to perform special medical tasks which can greatly reduce the possibility of medical staff in contacting with viruses. An overview of 6G applications is summarized in Fig. 10.

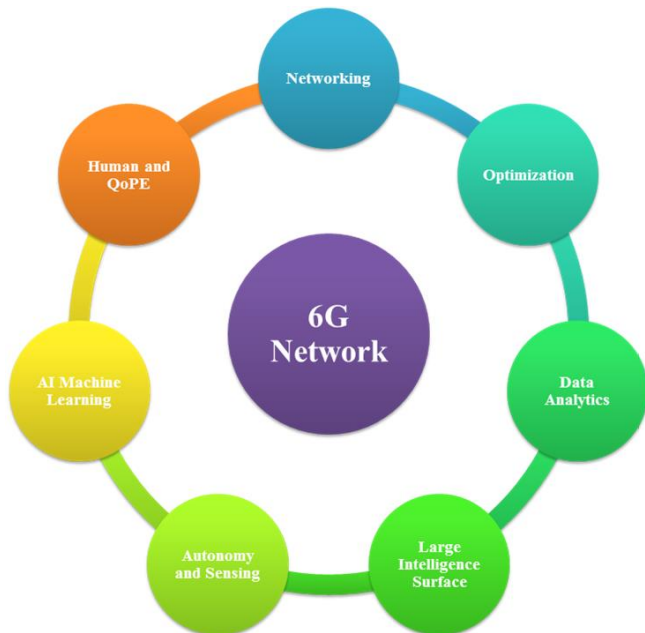


Fig. 10. Potential Applications of 6G.

## X. CONCLUSIONS

During the worldwide deployment of 5G, academia and industrial experts have started conceptualizing 6G. Unlike 5G networks, the next generation 6G will focus on communication among users, industries and multiple objects. Network transmission performance is no longer only important parameter; blockchain technology, IoT and AI have become important components. 6G network will keep penetrating into virtual society, human-perceived actions and ubiquitous spaces. It will provide a secure, reliable, intelligent, deep, seamless and holographic network infrastructure. 6G network will fulfil the growing demands of industries with continuous innovations of AI. We outlined research activities in different countries which aim to create a vision of 6G. 6G will enable many new technologies such as VLC, tactile and holographic communication. In conclusion, we expect that this review article will pave the way to identify 6G roadmap. This paper reviews the key technologies and areas of 6G networks and highlights a prospective on future research. We have presented a vision of 6G network as a research guide for readers. We have also addressed key features, security challenges and explained potential applications which will be supported in 6G. We have presented an overview of 1G to 6G. We then examine the key areas of 6G network. This review article started by highlighting the historical overview of communication technologies and their pivotal elements aiming at fostering future 6G in various dimensions. Then, we discussed technology breakdown, potential challenges associated with future 6G technology and possible solutions to foster 6G. In addition, we have profoundly examined research activities in

different countries including industries and research institutes. Finally, this study concludes with potential applications of future 6G. The key contribution of our study is that it clarifies the promising solution for potential issues and challenges in 6G technology. Thus, this review will open new horizons for future research directions.

## REFERENCES

- [1] P. Yang, Y. Xiao, M. Xiao, S. Li, 6g wireless communications: vision and potential techniques, *IEEE Network* 33 (4) (2019) 70–75.
- [2] K.B. Letaief, W. Chen, Y. Shi, J. Zhang, Y.-J.A. Zhang, The roadmap to 6g: Ai empowered wireless networks, *IEEE Commun. Mag.* 57 (8) (2019) 84–90.
- [3] T. Zhu, P. Xiong, G. Li, W. Zhou, S. Y. Philip, Differentially private model publishing in cyber physical systems, *Future Generat. Comput. Syst.*
- [4] K. David and H. Berndt, “6G vision and requirements: Is there any need for beyond 5G?” *IEEE Veh. Technol. Mag.*, vol. 13, no. 3, pp. 72–80, Sep. 2018.
- [5] T. Chen, S. Barbarossa, X. Wang, G. B. Giannakis, and Z.-L. Zhang, “Learning and management for internet of things: Accounting for adaptivity and scalability,” *Proc. IEEE*, vol. 107, no. 4, pp. 778–796, Apr. 2019.
- [6] W. Saad, M. Bennis, and M. Chen, “A vision of 6G wireless systems: Applications, trends, technologies, and open research problems,” *IEEE Netw.*, to be published, doi: 10.1109/MNET.001.1900287.
- [7] K. Wakunami et al., “Projection-type see-through holographic three-dimensional display,” *Nature Commun.*, vol. 7, no. 1, pp. 1–7, Oct. 2016.
- [8] E. Ahmed, I. Yaqoob, A. Gani, M. Imran, and M. Guizani, “Internet-of-things- based smart environments: State of the art, taxonomy, and open research challenges,” *IEEE Wireless Commun.*, vol. 23, no. 5, pp. 10–16, Oct. 2016.
- [9] N. Zhang, S. Zhang, P. Yang, O. Alhussein, W. Zhuang, and X. Shen, “Software defined space-air-ground integrated vehicular networks: Challenges and solutions,” *IEEE Commun. Mag.*, vol. 55, no. 7, pp. 101–109, Jul. 2017.
- [10] S. Zhang, S. Xu, G. Y. Li, and E. Ayanoglu, “First 20 years of green radios,” *IEEE Trans. Green Commun. Netw.*, vol. 4, no. 1, pp. 1–15, Mar. 2020.
- [11] P. Yang, Y. Xiao, M. Xiao, and S. Li, “6G wireless communications: Vision and potential techniques,” *IEEE Netw.*, vol. 33, no. 4, pp. 70–75, Jul. 2019.
- [12] K. B. Letaief, W. Chen, Y. Shi, J. Zhang, and Y.-J. A. Zhang, “The roadmap to 6G: AI empowered wireless networks,” *IEEE Commun. Mag.*, vol. 57, no. 8, pp. 84–90, Aug. 2019.
- [13] China Mobile Limited, “China mobile limited (2015) annual report 2015,” China Mobile Limited, Tech. Rep., 2016. [Online]. Available: <http://www.chinamobileltd.com/en/ir/reports/ar2015.pdf>
- [14] X. Gong and C. Cortese, “A socialist market economy with Chinese characteristics: The accounting annual report of China Mobile,” *Accounting Forum*, vol. 41, no. 3, pp. 206–220, 2017.
- [15] TTH Council Asia-Pacific, “APAC FTTH market panorama report 2019,” FTTH Council Asia-Pacific, Tech. Rep., 2019. [Online]. Available: [http://www.ftthcouncilap.org/wp-content/uploads/2019/04/FTTH-APAC-Panorama-Report-2019\\_Low.pdf](http://www.ftthcouncilap.org/wp-content/uploads/2019/04/FTTH-APAC-Panorama-Report-2019_Low.pdf)
- [16] Dang, Shuping, et al. “What should 6G be?.” *Nature Electronics* 3.1 (2020): 20-29.
- [17] Huang, Tongyi, et al. “A survey on green 6G network: Architecture and technologies.” *IEEE Access* 7 (2019): 175758-175768.
- [18] Raghavan, V. & Li, J. Evolution of physical-layer communications research in the post-5G era. *IEEE Access* 7, 10392-10401 (2019).
- [19] Calvanese Strinati, E. et al. 6G: the next frontier: from holographic messaging to artificial intelligence using subterahertz and visible light communication. *IEEE Veh. Technol. Mag.* 14, 42-50 (2019).



- [20] Saad, W., Bennis, M. & Chen, M. A vision of 6G wireless systems: applications, trends, technologies, and open research problems. *IEEE Netw.* <https://doi.org/10.1109/MNET.001.1900287> (2019).
- [21] Lu, Yang, and Xianrong Zheng. "6G: A survey on technologies, scenarios, challenges, and the related issues." *Journal of Industrial Information Integration* (2020): 100158.
- [22] E. Basar, Reconfigurable intelligent surface-based index modulation: a new beyond MIMO paradigm for 6G, *IEEE Trans. Commun.* 68 (5) (2020) 3187-3196.
- [23] S.M. Bohloul, Smart cities: a survey on new developments, trends, and opportunities, *J. Ind. Integr. Manag.* (2020) Early Access.
- [24] E. Carter, P. Adam, D. Tsakis, S. Shaw, R. Watson, P. Ryan, Enhancing pedestrian mobility in Smart Cities using Big Data, *J. Manag. Anal.* (2020) 1-16.
- [25] S. Chen, S. Sun, G. Xu, X. Su, Y. Cai, Beam-space multiplexing: practice, theory, and trends, from 4G TD-LTE, 5G, to 6G and Beyond, *IEEE Wirel. Commun.* 27 (2) (2020) 162-172.
- [26] M.Z. Chowdhury, M. Shahjalal, M. Hasan, Y.M. Jang, The role of optical wireless communication technologies in 5G/6G and IoT solutions: prospects, directions, and challenges, *Appl. Sci.* 9 (20) (2019) 4367.
- [27] T. Huang, W. Yang, J. Wu, J. Ma, X. Zhang, D. Zhang, A survey on green 6G network: architecture and technologies, *IEEE Access* 7 (2019) 175758-175768.
- [28] M. Matinmikko-Blue, S. Yrjölä, P. Ahokangas, Spectrum management in the 6G Era: the role of regulation and spectrum sharing, 2020 2nd Wireless Summit (6G SUMMIT), IEEE, 2020, March, pp. 1-5.
- [29] G. Gui, M. Liu, F. Tang, N. Kato, F. Adachi, 6G: opening new horizons for integration of comfort, security and intelligence, *IEEE Wireless Commun* (2020).
- [30] N. Kato, B. Mao, F. Tang, Y. Kawamoto, J. Liu, Ten challenges in advancing machine learning technologies toward 6G, *IEEE Wireless Commun.* (2020).
- [31] S.J. Nawaz, et al., Quantum machine learning for 6G communication networks: state-of-the-art and vision for the future, *IEEE Access* 7 (2019) 46317-46350.
- [32] L. Yan, C. Han, J. Yuan, Hybrid precoding for 6G terahertz communications: performance evaluation and open problems, 2020 2nd 6G Wireless Summit (6G SUMMIT), IEEE, 2020, pp. 1-5.
- [33] Alsharif, Mohammed H., et al. "Sixth Generation (6G) Wireless Networks: Vision, Research Activities, Challenges and Potential Solutions." *Symmetry* 12.4 (2020): 676.
- [34] Rommel, S., Raddo, T. R. & Monroy, I. T. Data center connectivity by 6G wireless systems. In *Proc. IEEE PSC* <https://doi.org/10.1109/PS.2018.8751363> (IEEE, 2018).
- [35] Clazzer, F. et al. From 5G to 6G: has the time for modern random access come? Preprint at <https://arxiv.org/abs/1903.03063> (2019).
- [36] Yaacoub, E. & Alouini, M.-S. A key 6G challenge and opportunity—connecting the remaining 4 billions: a survey on rural connectivity. Preprint at <https://arxiv.org/abs/1906.11541> (2019).
- [37] Giordani, M., Polese, M., Mezzavilla, M., Rangan, S. & Zorzi, M. Towards 6G networks: use cases and technologies. Preprint at <https://arxiv.org/abs/1903.12216> (2019).
- [38] Stoica, R.-A. & de Abreu, G. T. F. 6G: the wireless communications network for collaborative and AI applications. Preprint at <https://arxiv.org/abs/1904.03413> (2019).
- [39] Letaief, K. B., Chen, W., Shi, Y., Zhang, J. & Zhang, Y. A. The roadmap to 6G: AI empowered wireless networks. *IEEE Commun. Mag.* 57, 84-90 (2019).
- [40] Renzo, D. et al. Smart radio environments empowered by reconfigurable AI meta-surfaces: an idea whose time has come. *EURASIP J. Wireless Commun. Netw.* 2019, 129 (2019).
- [41] Zhao, J. A Survey of intelligent reflecting surfaces (IRSs): towards 6G wireless communication networks. Preprint at <https://arxiv.org/abs/1907.04789v3> (2019).
- [42] Basar, E. Reconfigurable intelligent surface-based index modulation: a new beyond MIMO paradigm for 6G. Preprint at <https://arxiv.org/abs/1904.06704v2> (2019).
- [43] Zong, Baiqing, et al. "6G technologies: Key drivers, core requirements, system architectures, and enabling technologies." *IEEE Vehicular Technology Magazine* 14.3 (2019): 18-27.
- [44] M.G. Kibria, K. Nguyen, G.P. Villardi, O. Zhao, K. Ishizu, F. Kojima, Big data analytics, machine learning, and artificial intelligence in next-generation wireless networks, *IEEE Access* 6 (2018) 32328-32338.
- [45] F. Tang, Y. Kawamoto, N. Kato, J. Liu, Future intelligent and secure vehicular network toward 6g: machine-learning approaches, *Proc. IEEE*.
- [46] I. Tomkos, D. Klonidis, E. Pikasis, S. Theodoridis, Toward the 6g network era: opportunities and challenges, *IT Prof.* 22 (1) (2020) 34-38.
- [47] L. Loven, T. Leppänen, E. Peltonen, J. Partala, E. Harjula, P. Porambage, M. Ylianttila, J. Riekk, Edge AI: A Vision for Distributed, Edge-Native Artificial Intelligence in Future 6g Networks, *The 1st 6G Wireless Summit*, 2019, pp. 1-2.
- [48] S. Dang, O. Amin, B. Shihada, M.-S. Alouini, What should 6g be? *Nat. Electron.* 3 (1) (2020) 20-29.
- [49] M. Katz, P. Pirinen, H. Posti, Towards 6g: getting ready for the next decade, in: 2019 16th International Symposium on Wireless Communication Systems (ISWCS), IEEE, 2019, pp. 714-718.
- [50] Y. Wei, H. Liu, J. Ma, Y. Zhao, H. Lu, G. He, Global voice chat over short message service of beidou navigation system, in: 2019 14th IEEE Conference on Industrial Electronics and Applications (ICIEA), IEEE, 2019, pp. 1994-1997.
- [51] T. Hong, C. Liu, M. Kadoch, Machine learning based antenna design for physical layer security in ambient backscatter communications, *Wireless Commun. Mobile Comput.* (2019).
- [52] S.J. Nawaz, S.K. Sharma, S. Wyne, M.N. Patwary, M. Asaduzzaman, Quantum machine learning for 6g communication networks: state-of-the-art and vision for the future, *IEEE Access* 7 (2019) 46317-46350.
- [53] P. Ferraro, C. King, R. Shorten, Distributed ledger technology for smart cities, the sharing economy, and social compliance, *IEEE Access* 6 (2018) 62728-62746.
- [54] K. Kotobi, S.G. Bilen, Secure blockchains for dynamic spectrum access: a decentralized database in moving cognitive radio networks enhances security and user access, *IEEE Veh. Technol. Mag.* 13 (1) (2018) 32-39.
- [55] S. Kiyomoto, A. Basu, M.S. Rahman, S. Ruj, On blockchain-based authorization architecture for beyond-5g mobile services, in: 2017 12th International Conference for Internet Technology and Secured Transactions (ICITST), IEEE, 2017, pp. 136-141.
- [56] S. Cho, G. Chen, J.P. Coon, Enhancement of physical layer security with simultaneous beamforming and jamming for visible light communication systems, *IEEE Trans. Inf. Forensics Secur.* 14 (10) (2019) 2633-2648.
- [57] S. Ucar, S. Coleri Ergen, O. Ozkasap, D. Tsonev, H. Burchardt, Secure visible light communication for military vehicular networks, in: *Proceedings of the 14th ACM International Symposium on Mobility Management and Wireless Access*, 2016, pp. 123-129.
- [58] J. Ma, R. Shrestha, J. Adelberg, C.-Y. Yeh, Z. Hossain, E. Knightly, J.M. Jornet, D.M. Mittleman, Security and eavesdropping in terahertz wireless links, *Nature* 563 (7729) (2018) 89-93.
- [59] I.F. Akyildiz, J.M. Jornet, C. Han, Terahertz band: next frontier for wireless communications, *Phys. Commun.* 12 (2014) 16-32.
- [60] J.-Y. Hu, B. Yu, M.-Y. Jing, L.-T. Xiao, S.-T. Jia, G.-Q. Qin, G.-L. Long, Experimental quantum secure direct communication with single photons, *Light Sci. Appl.* 5 (9) (2016), e16144.
- [61] L. Gyongyosi, S. Imre, H.V. Nguyen, A survey on quantum channel capacities, *IEEE Commun. Surv. Tutorials* 20 (2) (2018) 1149-1205.
- [62] W. Zhang, D.-S. Ding, Y.-B. Sheng, L. Zhou, B.-S. Shi, G.-C. Guo, Quantum secure direct communication with quantum memory, *Phys. Rev. Lett.* 118 (22) (2017), 220501.
- [63] E. C. Strinati, S. Barbarossa, J. L. Gonzalez-Jimenez, D. Ktenas, N. Cassiau, C. Dehos, 6g: the Next Frontier, *arXiv Preprint arXiv:1901.03239*.

- [64] Z. Zhang, Y. Xiao, Z. Ma, M. Xiao, Z. Ding, X. Lei, G.K. Karagiannidis, P. Fan, 6g wireless networks: vision, requirements, architecture, and key technologies, *IEEE Veh. Technol. Mag.* 14 (3) (2019) 28–41.
- [65] X. Ling, J. Wang, T. Bouchoucha, B.C. Levy, Z. Ding, Blockchain radio access network (b-ran): towards decentralized secure radio access paradigm, *IEEE Access* 7 (2019) 9714–9723.
- [66] M.S. Islim, R.X. Ferreira, X. He, E. Xie, S. Videv, S. Viola, S. Watson, N. Bamiedakis, R.V. Penty, I.H. White, et al., Towards 10 gb/s orthogonal frequency division multiplexingbased visible light communication using a gan violet microled, *Photon. Res.* 5 (2) (2017) A35–A43.
- [67] J. Luo, L. Fan, H. Li, Indoor positioning systems based on visible light communication: state of the art, *IEEE Commun. Surv. Tutorials* 19 (4) (2017) 2871–2893.
- [68] Mehedi, Syed Agha Hassnain Mohsan Md, et al. "A Systematic Review on Practical Considerations, Recent Advances and Research Challenges in Underwater Optical Wireless Communication." (2020): 11-7.
- [69] Mohsan, Syed Agha Hassnain, et al. "A Review on Research Challenges, Limitations and Practical Solutions for Underwater Wireless Power Transfer." *environment* 11.8 (2020).
- [70] S. Elmeadawy, R.M. Shubair, 6g wireless communications: future technologies and research challenges, in: 2019 International Conference on Electrical and Computing Technologies and Applications (ICECTA), IEEE, 2019, pp. 1–5.
- [71] R. Chen, C. Li, S. Yan, R. Malaney, J. Yuan, Physical layer security for ultra-reliable and low-latency communications, *IEEE Wireless Commun.* 26 (5) (2019) 6–11.
- [72] J.M. Hamamreh, E. Basar, H. Arslan, Odfm-subcarrier index selection for enhancing security and reliability of 5g urllc services, *IEEE Access* 5 (2017) 25863–25875.
- [73] Yamakami, A privacy threat model in xr applications, in: International Conference on Emerging Internetworking, Data & Web Technologies, Springer, 2020, pp. 384–394.
- [74] X. Chen, Y. Wang, M. Nakanishi, X. Gao, T.-P. Jung, S. Gao, High-speed spelling with a noninvasive brain–computer interface, *Proc. Natl. Acad. Sci. Unit. States Am.* 112 (44) (2015) E6058–E6067.
- [75] P. McCullagh, G. Lightbody, J. Zygierewicz, W.G. Kernohan, Ethical challenges associated with the development and deployment of brain computer interface technology, *Neuroethics* 7 (2) (2014) 109–122.
- [76] R.A. Ramadan, A.V. Vasilakos, Brain computer interface: control signals review, *Neurocomputing* 223 (2017) 26–44.
- [77] Pramanik, Pijush Kanti Dutta, et al. "Advancing Modern Healthcare With Nanotechnology, Nanobiosensors, and Internet of Nano Things: Taxonomies, Applications, Architecture, and Challenges." *IEEE Access* 8 (2020): 65230-65266.
- [78] Mohsan, Syed Agha Hassnain, et al. "Investigating Transmission Power Control Strategy for Underwater Wireless Sensor Networks."
- [79] Mohsan, Syed Agha Hassnain, et al. "Impact of Circular Field in Underwater Wireless Sensor Networks."

# Maximum Likelihood Classification based on Classified Result of Boundary Mixed Pixels for High Spatial Resolution of Satellite Images

Kohei Arai

Faculty of Science and Engineering  
Saga University, Saga City  
Japan

**Abstract**—Maximum Likelihood Classification: MLC based on classified result of boundary Mixed Pixels (Mixel) for high spatial resolution of remote sensing satellite images is proposed and evaluated with Landsat Thematic Mapper: TM images. Optimum threshold indicates different results for TM and Multi Spectral Scanner: MSS data. This may since the TM spatial resolution is 2.7 times finer than MSS, and consequently, TM imagery has more spectral variability for a class. The increase of the spectral heterogeneity in a class and the higher number of channels being used in the classification process may play significant role. For example, the optimum threshold for classifying an agricultural scene using MSS data is about 2.5 standard deviations, while that for TM corresponds to more than four standard deviations. This paper compares the optimum threshold between MSS and TM and suggests a method of using unassigned boundary pixels to determine the optimum threshold. Further, it describes the relationship of the optimum threshold to the class variance with a full illustration of TM data. The experimental conclusions suggest to the user some systematic methods for obtaining an optimal classification with MLC.

**Keywords**—Maximum likelihood classification; optimum threshold; Landsat TM; MSS; Mixed Pixel; spatial resolution

## I. INTRODUCTION

Now a day, high spatial resolution of remote sensing satellite imagery data is available. The finest resolution as of now is 30cm provided by very high-resolution commercial satellites. For instance, 50cm high-resolution image from Pleiades satellite (Airbus); it allows to clearly see the buildings, small boats, narrow streets. There is another example, a 50cm resolution image from SuperView-1 and 40cm (Komsat-3A), etc.

One of problems of utilization of these high spatial resolution of satellite images is comparatively poor classification performance. In general, variance of a designated class in a feature space is increased in accordance with the spatial resolution which result in poor classification performance. Another disadvantage of the high spatial resolution of satellite image classification is determination of optimum threshold for the discrimination between classes in the well-known Maximum Likelihood classification or some other classification methods such as Support Vector Machine, Deep Learning Based Method, etc. This paper is intended to solve the latter problem.

Supervised maximum likelihood classification based on multidimensional normal distribution for each pixel is widely used as a classification method for remote sensing data. In this method, the training area of each class of interest is first set by the analyst, and each class is defined by the distribution of the classes in the probability space from the samples in the area. Then, using the likelihood between the distribution and the pixel to be classified as an index, the pixel is classified into a distribution class showing the highest likelihood.

However, if the likelihood is not so large, it is determined that it is better not to be classified into some forcibly set classes, so that the pixel is set as an unset class. At this time, the value of the likelihood for deciding whether to make the class unset is defined as a threshold here. In general, the threshold is determined empirically by repeated trial and error using the classified image according to the subjectivity of the analyst. Therefore, it takes a lot of time to find the optimal threshold and lacks general bias.

The optimal threshold decision method proposed in this paper extracts boundary pixels in different class regions, performs maximum likelihood classification with a certain threshold on those pixels, and determines the number of pixels classified into unset classes. It is optimized with a threshold value that is half the number of pixels. In this method, the target region to be classified for obtaining the optimum threshold is limited to the boundary pixels, and therefore the optimum threshold can be determined objectively in a short time.

This paper first describes related research works followed by the definition of the optimal threshold together with research background. Then, the optimal thresholds of Landsat TM and MSS data with different Instantaneous Field of View: IFOV are subjectively obtained by the conventional method, and the two elderly are compared to clarify the relationship between the optimal threshold and the IFOV. Finally, the proposed method is explained, and it is shown using TM data that it matches the optimal threshold value obtained by the conventional method, and the validity of the proposed method is shown.

## II. RELATED RESEARCH WORKS

Classification by re-estimating statistical parameters based on auto-regressive model is proposed for purification of

training samples [1]. Meanwhile, multi-temporal texture analysis in TM classification is proposed for high spatial resolution of optical sensor images [2]. On the other hand, Maximum Likelihood (MLH) TM classification considering pixel-to-pixel correlation is proposed [3].

Supervised TM classification with a purification of training samples is proposed [4] together with TM classification using local spectral variability is proposed [5]. A classification method with spatial spectral variability is also proposed [6] together with TM classification using local spectral variability [7].

Application of inversion theory for image analysis and classification is proposed [8]. Meanwhile, polarimetric Synthetic Aperture Radar: SAR image classification with maximum curvature of the trajectory in eigen space domain on the polarization signature is proposed [9]. On the other hand, A hybrid supervised classification method for multi-dimensional images using color and textural features is proposed [10].

Polarimetric SAR image classification with high frequency component derived from wavelet multi resolution analysis: MRA is proposed [11]. Comparative study of polarimetric SAR classification methods including proposed method with maximum curvature of trajectory of backscattering cross section in ellipticity and orientation angle space is conducted and well reported [12].

Human gait gender classification using 2D discrete wavelet transforms energy is proposed [13] together with human gait gender classification in spatial and temporal reasoning [14].

Comparative study on discrimination methods for identifying dangerous red tide species based on wavelet utilized classification methods is conducted [15]. On the other hand, multi spectral image classification method with selection of independent spectral features through correlation analysis is proposed [16].

Image retrieval and classification method based on Euclidian distance between normalized features including wavelet descriptor is proposed [17]. Meanwhile, gender classification method based on gait energy motion derived from silhouettes through wavelet analysis of human gait moving pictures is proposed [18] together with human gait skeleton model acquired with single side video camera and its application and implementation for gender classification [19].

Gender classification method based on gait energy motion derived from silhouette through wavelet analysis of human gait moving pictures is proposed [20] together with human gait gender classification using 3D discrete wavelet transformation feature extraction [21].

Wavelet Multi-Resolution Analysis: MRA and its application to polarimetric SAR classification is proposed [22]. On the other hand, object classification using a deep convolutional neural network and its application to myoelectric hand control is proposed and evaluated effectiveness [23]. On the other hand, image classification considering probability density function based on Simplified

beta distribution is proposed and validated with remote sensing satellite imagery data [24].

### III. RESEARCH BACKGROUND

#### A. Definition of Optimum Thershold for Maximum Likelihood Classification

The log likelihood function  $g_i(x)$  of the observation vector  $x$  with respect to the class  $i$  in the maximum likelihood classification [25] can be expressed by Eq. (1).

$$g_i(x) = \frac{1}{2 \ln \det \Sigma_i^{-1}} - \frac{1}{2} (x - \mu_i)^t \Sigma_i^{-1} (x - \mu_i) \quad (1)$$

where,  $\mu_i$  and  $\Sigma_i$  are the mean vector and the covariance matrix of class  $i$ ,  $\det$ . Is the determinant, and  $t$  and  $-1$  are the transpose and inverse matrix, respectively.  $x$  is classified into a class  $i$  in which  $g_i(x)$  is maximum. At this time, if  $g_i(x)$  does not exceed the threshold value  $\theta$ ,  $x$  is classified into an unset class. No, usually  $\theta$  is defined by likelihood.

Suppose now that the classes  $i$  and  $j$  and other unset classes are considered as shown in Fig. 1. In the figure,  $P_{u1}$  to  $P_u$  are the probabilities of being classified into the unset class,  $P_{ci}$   $P_{cj}$  is the probability of taking pixels that originally belong to the set class into classes  $i$  and  $j$ , and  $P_{ci}$ ,  $P_{cj}$  is the pixel that originally belongs to class  $j$ .

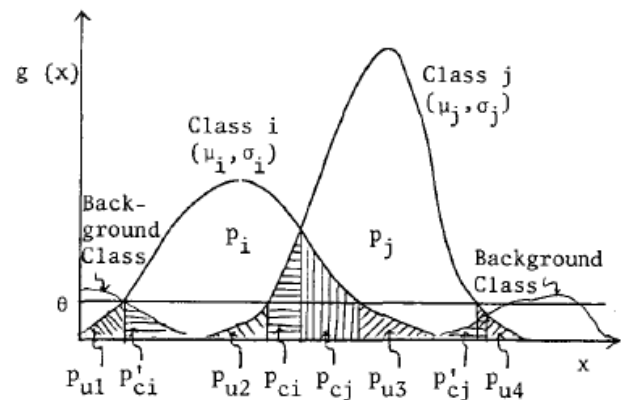


Fig. 1. Explanatory Illustration for the Determination of Optimum Threshold of the Pointwise Maximum Likelihood Classification in the case of Two Designated Classes and Background Class.

The probability of classifying pixels of class  $i$  into class  $j$ .  $P_i$  and  $P_j$  are the correct answer rates for classifying pixels belonging to classes  $i$  and  $j$  into classes  $i$  and  $j$ , respectively. Here, the correct answer rate  $P_c$ , the misclassification probability  $P_e$ , and the classification into unset classes, that is, the probability  $P_u$  not classified because the confidence level of the classification result is low are defined by the following equations.

$$P_c = P_i + P_j \quad (2)$$

$$P_e = P_{ci}' + P_{ci} + P_{cj}' + P_{cj} \quad (3)$$

$$P_u = P_{u1} + P_{u2} + P_{u3} + P_{u4} \quad (4)$$

When the threshold value  $\theta$  is increased,  $P_c$  and  $P_e$  decrease, and  $P_u$  increases. The amount of the increase or decrease depends on the probability density function of the classes  $i$  and  $j$  and the distance between classes. The threshold

that maximizes (aPc-bPecPu) is defined as the optimal threshold  $\theta_{opt}$ . Here, a, b, and c are weighting factors for each probability, and are determined based on the subjectivity of the analyst. Therefore, according to the subjective opinion of the analyst,  $\theta_{opt}$  is obtained by repeating trial and error using the classification result with several threshold values (this is called the conventional method here).

### B. Example of Optimum Thershold

1) *Analysis image:* Acquired by Landsat 5 TM on August 15, 1984. Using the data of Cranbrook, British Columbia, Canada (Path: 43, Row: 25), the optimal threshold was determined by the conventional method [26]. This area is mostly covered with forests, including logging areas, alpine meadows, roads, rivers, etc. The following four classes were set from these. i) New logging area, ii) Old logging area, iii) Alpine grassland, iv) Forest where, a new felling area was defined as a forest that was harvested within 5 years, an old logging area was harvested between 5 and 40 years, and an area with a tree age of more than 40 years was defined as a forest.

2) *Optimal threshold:* Fig. 2(a) and (b) show examples of classified images when a part of the original image and the

threshold  $\theta$  are changed to  $-40$ ,  $-30$ , and  $-20$ , respectively. In the case of  $\theta = -20$ , the probability of being classified into an unset class is high, and it is difficult to say that the threshold is optimal. Looking at the classified image of  $b = -40$ , the shaded forest on the upper left and the exposed part of the rock on the left are misclassified into the class of "new logging area" with similar spectral characteristics. On the other hand, in the classified image of  $\theta = -30$ , those portions are classified into the unset class, and it turns out that the threshold value is the optimal threshold after all.

### C. Resolution and Optimal Threshold

The optimal threshold for MSS data is said to be the likelihood corresponding to two to three times the standard deviation of the band showing the maximum variance. What is the optimal threshold for TM data with improved quantization bit rate and instantaneous visual field? With improved resolution, the intra-class variance increases [27] and the optimal threshold generally decreases. To illustrate this, the optimal thresholds were compared using both TM and MSS data acquired simultaneously by Landsat 5 on July 29, 1984. The target area is an agricultural area in Melfort, Saskatchewan, Canada.

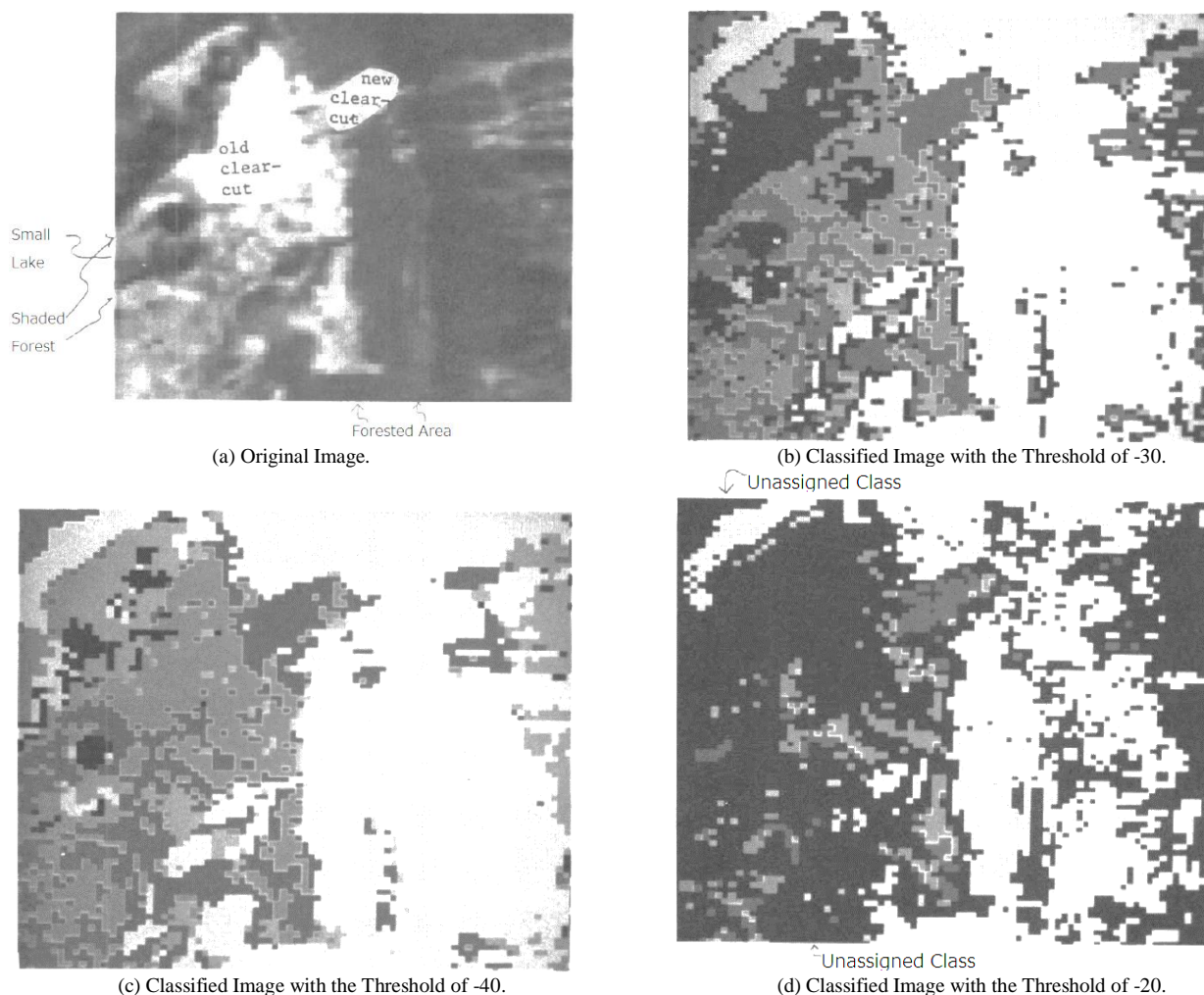


Fig. 2. Original Image and Examples of Classified Images with Several Threshold.

The following classes were set from the main components.

- 1) Urban area,
- 2) Barley field,
- 3) Wheat field,
- 4) Bean field,
- 5) Fallow field

Comparing the logarithm of  $\det.\Sigma$  (generalized variance), which indicates the degree of variance of each class, with TM and MSS data, TM is larger as shown in Table I. In addition, using the average vector of each class and the covariance matrix, the average of the Mahalanobis distance<sup>1</sup> with the observation vector in the training area is obtained, and when comparing TM and MSS, it is found that TM is smaller as shown in the table. That is, it is understood that the likelihood is large.

The maximum likelihood method with various thresholds was applied to both data, and the optimal threshold was determined by the conventional method of subjectively evaluating the classified images. As a result, it was found that the optimum value is about 18 for MSS and about 30 for TM. Examining this in Fig. 3 showing the relationship between the threshold and the classification accuracy Pc of each class, the threshold at which all Pc is almost saturated is optimal.

TABLE I. GENERALIZED VARIANCE AND MAHALANOBIS DISTANCE ON A PER CLASS BASIS

| Class Name | $\log_{10} \det.\Sigma$ | $(x-\mu)^T(\Sigma)^{-1}(x-\mu)$ |
|------------|-------------------------|---------------------------------|
|            | TM_MSS                  | TM_MSS                          |
| Urban      | 10.9_5.91               | 12.09_19.01                     |
| Barley     | 6.62_3.28               | 14.72_23.38                     |
| Wheat      | 3.71_2.91               | 15.09_26.27                     |
| Canola     | 4.46_4.99               | 13.01_25.54                     |
| Fallow     | 7.22_5.12               | 12.88_22.78                     |

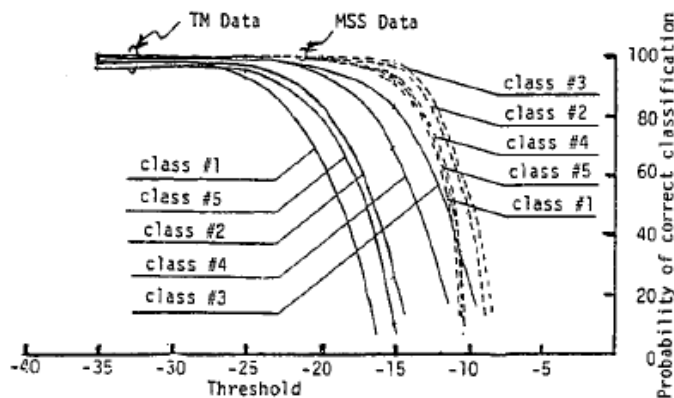


Fig. 3. A Comparison between Optimum Thresholds for TM and MSS Data (Optimum Threshold for TM is around -30, while that for MSS is about -18) TM 1-5, 7, MSS 1-4 where the Classes are: 1) Urban, 2) Barley Field, 3) Wheat Field, 4) Canola Field and 5) Fallow Field. Fallow Field.

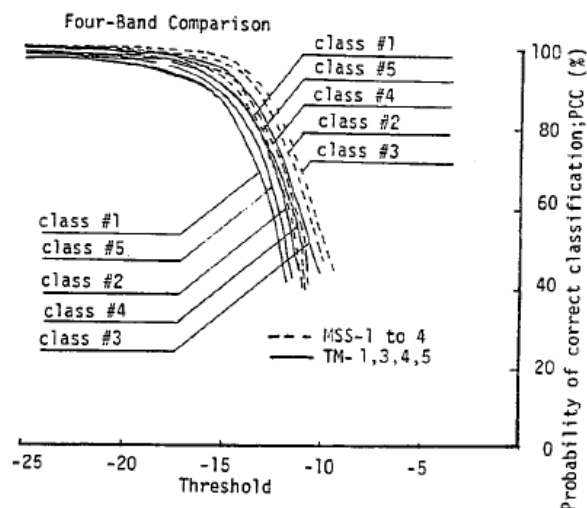


Fig. 4. Changes in the Optimum Threshold Due to differences in Both Radiometric and Spatial Resolutions. where the Classes are: 1) Urban, 2) Barley Field, 3) Wheat Field, 4) Canola Field and 5) Fallow Field.

However, the number of TM bands in this comparison is 6, and that of MSS is different from 4. Since the likelihood comparison cannot be performed without the dimensions of the variables, the TM dimension is further reduced to 4. investigated. As a method of dimensionality reduction, the method using only the degree of separation as an index is considered, and selected TM bands 1, 3, 4, and 5. As a result, as shown in Fig. 4, the optimum threshold value for 4-Pand TM data is about -22, which is lower than that of MSS. This difference is thought to be due to the difference of class variance in feature space based on the difference of spatial resolution.

#### IV. PROPOSED METHOD (MAXIMUM LIKELIHOOD CLASSIFICATION WITH CLASSIFIED BOUNDARY PIXEL

Looking at the classified image at the optimal threshold value obtained by the analyst's subjective view, and Fig. 2(b), about half of the boundary pixels of different setting classes are classified as unset classes. These boundary pixels are also called Mixels (Mixed Pixels) and have at least two or more types of spectral classes having complex spectral characteristics weighted by the area ratio of pixels in each class. Focusing on a certain class  $i$ , consider a boundary pixel set in contact with this class area. To simplify the discussion, consider a one-dimensional pixel array, and consider the classification of boundary pixels at the sampling phase of In-Phase and Out of Phase shown in Fig. 5.

Since the sampling phase is uncorrelated between the radiometer scanning time and the boundary position of the surface object, it is rare in the case of In-Phase, almost out of phase, and independent of the position within the boundary pixel. It is a uniform probability. Now, when the likelihood is classified based on 50%, the probability of being classified into an unset class is 0, and the boundary pixels are classified into either  $i$  or  $j$  class.

<sup>1</sup> [https://en.wikipedia.org/wiki/Mahalanobis\\_distance](https://en.wikipedia.org/wiki/Mahalanobis_distance)

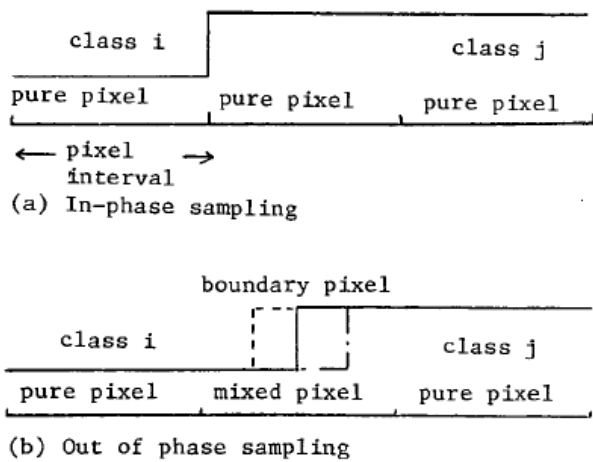


Fig. 5. In-Phase and Out of Phase Samplings.

If the likelihood is classified based on 100%, all boundary pixels are classified into an unset class. The criterion of the likelihood that half of the boundary pixels are classified into the unset class is 75% shown by the broken line and the dashed line in Fig. 5(b). The proposed scheme is based on this. Although the discussion so far is valid for one class, finding the optimal threshold common to all classes is more complicated and difficult to derive theoretically.

Therefore, the proposed method considers a state where the probability of being classified into the unset class of the boundary pixel set is 0.5 as the first approximation of the optimal threshold. Based on this assumption, using the extracted boundary pixel classification results, the threshold value such that half of them are classified into the unset class was considered optimal. Fig. 6 shows a flowchart of the optimum threshold value determination method. In the figure, the mask pattern is a binary image in which the boundary pixel is "1" and the other pixels are "0".

### V. EXPERIMENT

Numerous methods for detecting boundary pixels have been proposed depending on the application and the definition of the boundary [28]. In this paper, a method [29] that uses the variance of local spectral information emphasized in high-resolution images as an index in used. Specifically, the boundary pixel was extracted by binarizing the variation coefficient of the pixel value in a  $2 \times 2$  pixels cell by a predetermined threshold.

Next, an example of extraction using TM data used in the analysis is shown in 22. Table II shows the average, variance, and coefficient of variation of each class for Band 4 (TM-4), which has the largest variance of the set class. The minimum value of the difference in the coefficient of variation between the classes is 0.176. If this value is set as a threshold and the image representing the variation coefficient is binarized, the boundary pixels between the classes can be extracted. To

confirm this, boundary pixels were extracted by changing the threshold to 0.15, 0.2, 0.3. As a result, as shown in Fig. 7, the threshold is best when the threshold is 0.15. Confirmed that it represents the world.

A maximum likelihood classification with a threshold of -10 to -40 was applied to the above extracted boundary pixels to obtain a classified image. After that, the ratio of the pixels classified into the unset class in the boundary pixels was calculated, and the relationship with the threshold was examined.

As a result, as shown in Fig. 8, the threshold value corresponding to 50% of the pixels classified into the unset class among the boundary pixels is about 130, and this value is 2.2, which is subjectively calculated by the conventional method. The validity of the proposed method was evident from the fact that it almost coincided with the above.

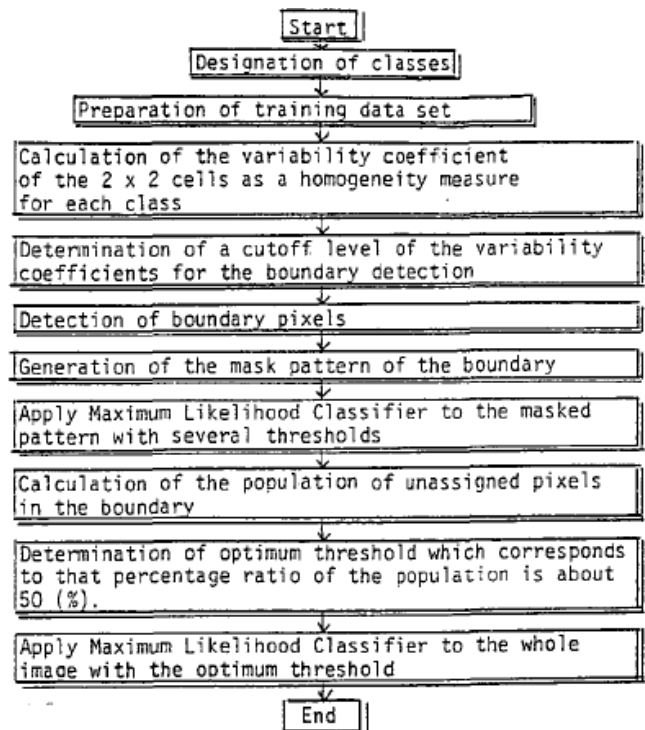


Fig. 6. Schematic Diagram for the Determination of the Optimum Threshold of Maximum Likelihood Classifier.

TABLE II. AVERAGE AND VARIANCE OF THE 2X2 CELLS OF TM4 (TM BAND 4)

| Class         | Sample | Mean  | Variance | Variability_Coefficient |
|---------------|--------|-------|----------|-------------------------|
| New_Cleacut   | 956    | 5.09  | 45.7     | 1.33                    |
| Old_Clearcut  | 270    | 4.99  | 109      | 2.09                    |
| Alpine_Meadow | 435    | 19.27 | 2528.32  | 2.61                    |
| Forest        | 1377   | 1.4   | 11.6     | 2.43                    |

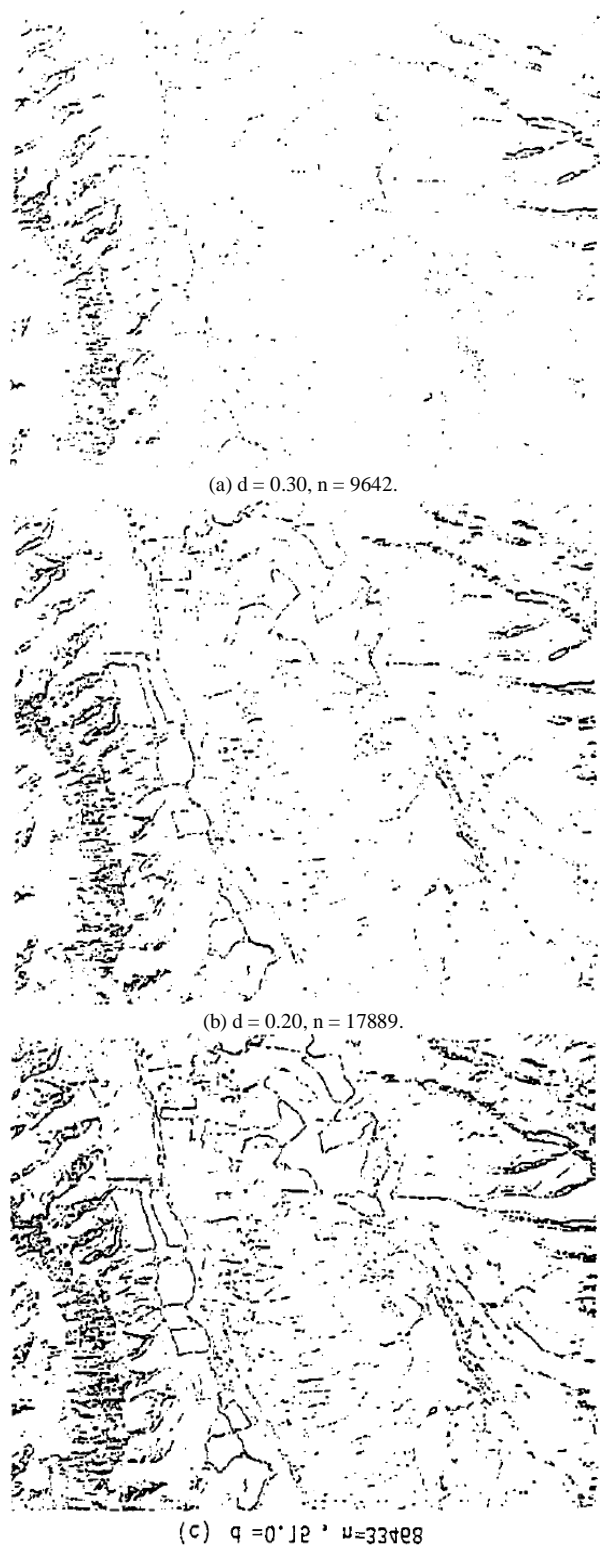


Fig. 7. Examples of the Boundary Images with the Variability Coefficient Cutoff(s) of 0.15, 0.20 and 0.30 for TM-4 of B.C. Forest where n is the Number of Boundary Pixels.

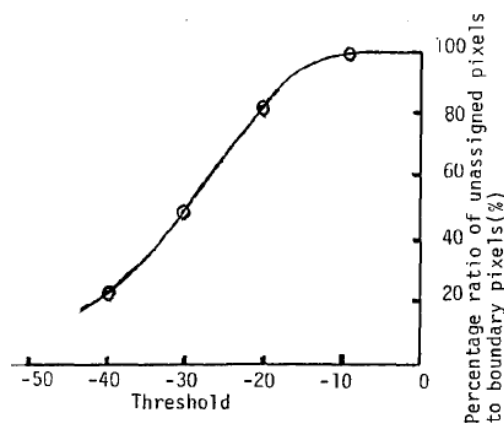


Fig. 8. Relationship between Threshold and the Population of unassigned Pixels in the Boundary Pixels, around 50(%) of the Population Corresponds to the Optimum Threshold (-30).

## VI. CONCLUSION

When applying the maximum likelihood classification to high-resolution images, the following points need to be considered to determine the optimal threshold. That is, the high-resolution image has a relatively large variation in the spectral information, so that the threshold value needs to be selected low. In addition, the optimal threshold decision method proposed in this paper, which has a threshold value at which the occupation probability of a pixel classified into an unset class at the boundary pixel becomes 0.5, has the following characteristics.

1) According to the conventional method, it is determined by the subjectivity of the analyst, and it is difficult to compare the classification results. However, according to the present method, it can be obtained objectively.

2) In the conventional method, several maximum likelihood methods with different thresholds are applied to the entire analysis image and determined by trial and error. However, according to this method, classification is performed only on extracted boundary pixels. Since the optimum threshold value is obtained, the processing time is short.

3) This method is based on a completely new concept of determining the optimal threshold of the maximum likelihood method using the probability of being classified into an unset class of boundary pixels.

## VII. FUTURE RESEARCH WORKS

The proposed method must be validated with a variety of high spatial resolution of optical sensors onboard satellites.

## ACKNOWLEDGMENT

This study was conducted by the author as a Canadian Government Scholarship Student (Postdoctoral Fellow) at the Canadian National Remote Sensing Center (CCRS). The author would like to thank Dr. D.G. Goodenough, Dr. J. Iisaka Mr. K. B. Fung, Mr. M. A. Robson and others who participated in the discussion.



The author, also, would like to thank Professor Dr. Hiroshi Okumura and Professor Dr. Osamu Fukuda for their valuable discussions.

#### REFERENCES

- [1] Kohei Arai, Classification by Re-Estimating Statistical Parameters Based on Auto-Regressive Model, Canadian Journal of Remote Sensing, Vol.16, No.3, pp.42-47, Jul.1990.
- [2] Kohei Arai, Multi-Temporal Texture Analysis in TM Classification, Canadian Journal of Remote Sensing, Vol.17, No.3, pp.263-270, Jul.1991.
- [3] Kohei Arai, Maximum Likelihood TM Classification Taking into account Pixel-to-Pixel Correlation, Journal of International GEOCARTO, Vol.7, pp.33-39, Jun.1992.
- [4] Kohei Arai, A Supervised TM Classification with a Purification of Training Samples, International Journal of Remote Sensing, Vol.13, No.11, pp.2039-2049, Aug.1992.
- [5] Kohei Arai, TM Classification Using Local Spectral Variability, Journal of International GEOCARTO, Vol.7, No.4, pp.1-9, Oct.1992.
- [6] Kohei Arai, A Classification Method with Spatial Spectral Variability, International Journal of Remote Sensing, Vol.13, No.12, pp.699-709, Oct.1992.
- [7] Kohei Arai, TM Classification Using Local Spectral Variability, International Journal of Remote Sensing, Vol.14, No.4, pp.699-709, 1993.
- [8] Kohei Arai, Application of Inversion Theory for Image Analysis and Classification, Advances in Space Research, Vol.21, 3, 429-432, 1998.
- [9] Kohei Arai and J.Wang, Polarimetric SAR image classification with maximum curvature of the trajectory in eigen space domain on the polarization signature, Advances in Space Research, 39, 1, 149-154, 2007.
- [10] Hiroshi Okumura, Makoto Yamaura and Kohei Arai, A hybrid supervised classification method for multi-dimensional images using color and textural features, Journal of the Japanese Society of Image Electronics Engineering, 38, 6, 872-882, 2009.
- [11] Kohei Arai, Polarimetric SAR image classification with high frequency component derived from wavelet multi resolution analysis: MRA, International Journal of Advanced Computer Science and Applications, 2, 9, 37-42, 2011.
- [12] Kohei Arai Comparative study of polarimetric SAR classification methods including proposed method with maximum curvature of trajectory of backscattering cross section in ellipticity and orientation angle space, International Journal of Research and Reviews on Computer Science, 2, 4, 1005-1009, 2011.
- [13] Kohei Arai, Rosa Andrie Asmara, Human gait gender classification using 2D discrete wavelet transforms energy, International Journal of Computer Science and Network Security, 11, 12, 62-68, 2011.
- [14] Kohei Arai, R.A.Asunara, Human gait gender classification in spatial and temporal reasoning, International Journal of Advanced Research in Artificial Intelligence, 1, 6, 1-6, 2012.
- [15] Kohei Arai, Comparative study on discrimination methods for identifying dangerous red tide species based on wavelet utilized classification methods, International Journal of Advanced Computer Science and Applications, 4, 1, 95-102, 2013.
- [16] Kohei Arai, Multi spectral image classification method with selection of independent spectral features through correlation analysis, International Journal of Advanced Research in Artificial Intelligence, 2, 8, 21-27, 2013.
- [17] Kohei Arai, Image retrieval and classification method based on Euclidian distance between normalized features including wavelet descriptor, International Journal of Advanced Research in Artificial Intelligence, 2, 10, 19-25, 2013.
- [18] Kohei Arai, Rosa Andrie Asmara, Gender classification method based on gait energy motion derived from silhouettes through wavelet analysis of human gait moving pictures, International Journal of Information Technology and Computer Science, 6, 3, 1-11, 2014.
- [19] Kohei Arai, Rosa Andrie Asmara, Human gait skeleton model acquired with single side video camera and its application and implementation for gender classification, Journal of the Image Electronics and Engineering Society of Japan, Transaction of Image Electronics and Visual Computing, 1, 1, 78-87, 2014.
- [20] Kohei Arai, Rosa Andrie Asmara, Gender classification method based on gait energy motion derived from silhouette through wavelet analysis of human gait moving pictures, International Journal of Information technology and Computer Science, 5, 5, 12-17, 2013.
- [21] Kohei Arai, Rosa Andrie Asmara, Human gait gender classification using 3D discrete wavelet transformation feature extraction, International Journal of Advanced Research in Artificial Intelligence, 3, 2, 12-17, 2014.
- [22] Kohei Arai, Wavelet Multi-Resolution Analysis and Its Application to Polarimetric SAR Classification, Proceeding of the SAI Computing Conference 2016,
- [23] Yoshinori Bando, Nan Bu, Osamu Fukuda, Hiroshi Okumura, Kohei Arai, Object classification using a deep convolutional neural network and its application to myoelectric hand control, Proceedings of the International Symposium on Artificial Life and Robotics (AROB2017), GS12, 2017.
- [24] Kohei Arai, Image classification considering probability density function based on Simplified beta distribution, International Journal of Advanced Computer Science and Applications IJACSA, 11, 4, 481-486, 2020.
- [25] Colwell, R.N.Edt., Manual of Remote Sensing, The American Society of Photogrammetry, 2nd Edition, Vol.1, p.800-801, 1983.
- [26] Kohei Arai, et al, On optimum threshold of Maximum Likelihood Classifier for TM classification, The 10th Canadian Symposium on Remote Sensing, p.8., 1986.
- [27] Kohei Arai, Multispectral image classification using spatial information, The Remote Sensing Society of Japan, Vol. 7, No. 4, p. 17-24, 1987.
- [28] Kast, J.L. et al, ECHO user's guide, LARS Publication No. 083077, Purdue Univ., P.72, 1977.
- [29] Kohei Arai., Methodologies for TM classification, The 6th Image Systems and Artificial Intelligence Working Group Meeting, p.24, 1987.

#### AUTHOR'S PROFILE

**Kohei Arai**, He received BS, MS and PhD degrees in 1972, 1974 and 1982, respectively. He was with The Institute for Industrial Science and Technology of the University of Tokyo from April 1974 to December 1978 also was with National Space Development Agency of Japan from January, 1979 to March, 1990. During from 1985 to 1987, he was with Canada Centre for Remote Sensing as a Post Doctoral Fellow of National Science and Engineering Research Council of Canada. He moved to Saga University as a Professor in Department of Information Science on April 1990. He was a councilor for the Aeronautics and Space related to the Technology Committee of the Ministry of Science and Technology during from 1998 to 2000. He was a councilor of Saga University for 2002 and 2003. He also was an executive councilor for the Remote Sensing Society of Japan for 2003 to 2005. He is a Science Council of Japan Special Member since 2012. He is an Adjunct Professor of University of Arizona, USA since 1998. He also is Vice Chairman of the Science Commission "A" of ICSU/COSPAR since 2008 then he is now award committee member of ICSU/COSPAR. He wrote 55 books and published 620 journal papers as well as 450 conference papers. He received 66 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>.

# Weather Variability Forecasting Model through Data Mining Techniques

Sultan Shekana<sup>1</sup>, Addisu Mulugeta<sup>2</sup>

Lecturer, AMiT,  
Arba Minch University  
Arba Minch, Ethiopia

Durga Prasad Sharma<sup>3</sup>

Professor, AMUIT, MOEFDRE under UNDP  
Expat International Consultant  
ILO (United Nations)

**Abstract**—Climate and weather variability are thought-provoking for world communities. In this apprehension, weather variability imposes a broad impact on the economy and the survival of the living entities. In relation to the African continent country Ethiopia, it is desirable to have great attention on the weather variability. The Ethiopian Dodota Woreda region is continuously affected by repeated droughts. It gives a great alarm to investigate and analyze the factors which are major causes of the frequent occurrence of drought. Although the weather scientists and domain experts are overwhelmed with meteorological data but lacking in analyzing and revealing the hidden knowledge or patterns about weather variability. This paper is an effort to design an enhanced predictive model for weather variability forecasting through Data Mining Techniques. The parameters used in this research are temperature, dew point, sunshine, rainfall, wind speed, maximum temperature, minimum temperature, and relative humidity to enhance the accuracy of forecasting. To improve the accuracy, we used the Multilayer-perceptron (MLP), Naïve Bayes, and multinomial logistic regression algorithms to design a proposed Predictive Model. The knowledge discovery in database (KDD) process model was used as a framework for the modeling purpose. The research findings revealed that the aforementioned parameters have a strong positive relationship with weather forecasting in meteorology sectors. The MLP model with selected parameters presents an interesting predictive accuracy result i.e. 98.3908% as correctly classified instances. The most performing algorithm, MLP was chosen and used to generate interesting patterns. The domain experts (meteorologists) validated the discovered patterns for the improved accuracy of weather variability forecasting.

**Keywords**—*Meteorological data; weather forecasting; multilayer-perceptron; Naïve Bayes; multinomial logistic regression algorithms*

## I. INTRODUCTION

The prime objective of this study is to design an enhanced predictive model for weather variability forecasting of Dodota Woreda (a region of Ethiopia) site using data mining techniques. Primarily it was aimed to rigorously review the existing models to find out the deficiencies in weather variability forecasting. This review was done in special reference to the site (Dodota Woreda) selected for the study and primary data collection. The determinant parameters that affect the weather variability were targeted to identify and select suitable data mining techniques for designing a predictive model. The performance of the model was validated

by domain experts i.e. meteorologists. It was also marked to identify which data mining algorithm is better performing in weather variability forecasting.

The data set for this study was collected from the national meteorology agency of the Dodota Woreda region and Awash-Melkasa station. 10 years of the daily dataset from 2006 to 2016 were collected to analyze the objectives specified in the study. The meteorological dataset selected for this research contains 7 parameters and 5282 records. The major limitation of this research study is the unavailability of sufficient datasets for some of the selected parameters of the climate.

The significance of this research study is to support the meteorology department for accurate weather variability forecasting. For this purpose, the research efforts were made to design a predictive model with enhanced performance. The interesting patterns/rules extracted from this research can be used to support the decision-making processes in the domain. Timely information about the weather variability allows people to make better precautions and preparations to alleviate the disastrous challenges. The enhanced forecasting can also help in improved decision-making processes such as: 1) weather threatens to life and property, 2) daily planning for outdoor activities, and 3) routine weather dependent economic activities.

In general, the research can help in improving the current forecasting methods/mechanism especially in minimizing the forecasting errors. In addition, the Govt. institutions like the National Meteorology Agency, Agriculture Department, Dodota Woreda Administration, Community, and Non-Govt. Institution can also use the outcomes of this research. The major challenges faced by the world Meteorologist are the accuracy of the weather variability and its predictive analysis [1]. Scientists have tried to forecast meteorological characteristics using numerous methods. Some of these methods being used provide better accuracy than others [2]. The current practices of the weather variability forecasting include ground observation, observation from the ships, aircraft, Doppler radars and satellites. We need accurate judgments of temperature at a particular time for various reasons such as the planning for the individuals' daily activity, the farmers need for planting and harvesting their crops, agricultural and technical systems need the assessment of the natural hazards, and to design the solar energy systems [3]. The necessity for accurate weather variability prediction is not questionable when the benefits are higher than the expenses.

Accuracy of the weather variability is based on the selected algorithms and parameters. Numerous techniques such as linear regression, auto-regression, Multi-Layer Perceptron, and back-propagation neural network is being applied to predict the weather variability using parameters such as temperature, wind speed, rainfall, dew point, and meteorological pollution etc. [4]. The most effective way to minimize the disastrous damage due to weather variability is forewarning. Though the accurate prediction may not stop a famine or flood, they can help people to prepare in advance.

The meteorological office of the forecasting of Ethiopia uses a model 24\*7\*365 of the year, using one of the world's fastest supercomputers to predict the weather variability for hours, days, weeks, seasons, and even years ahead. Despite all such efforts, weather forecasting still needs a significant improvement in forecasting accuracy. This is due to the complexity of the atmosphere and lack of observational data with limited consideration of weather parameters [5]. Even though, a huge amount of the meteorological dataset is available at the meteorology department, but it is not adequately, and properly analyzed for weather variability forecasting. The discovered new knowledge and hidden insights can be useful support in decision making and strategic planning processes.

The world Climate and Weather are the two critical phenomena that affect human lives. Weather is a short-term phenomenon and its variability has a broad and far-reaching set of impact by imposing significant loss of lives and cause illness. It also affects the economy in terms of transportation blockage, the decline in agricultural production, and land erosion. Although some of the research studies and projects have been undertaken in the Dodota Woreda region covering limited knowledge with scientific observations on weather forecasting and its impacts. These studies used different sets of parameters for predicting weather variability in terms of sunny, cloudy, and rainy which are used for the precautionary measures for the community lives, economic development, and to save the community from drought, food security, unpredictable rain, flood, and soil erosion.

In the year 2008 (ETC) many people at Dodota Woreda faced weather variability challenges and died due to great famine and drought shocks. The food shortage and drought were tightly coupled with the rainfall variability and caused a situation of high dependency on national and international food aids. During the preliminary assessment, it was found that the number of farmers at Dodota Woreda was displaced due to disastrous weather and soil erosions. With the researcher's point of view, the impact of weather variability is still continued and considered to be one of the important study focus areas for the Dodota Woreda region. The extremes of damages by climatic disasters and weather variability cannot be avoided completely, but a forewarning or an advance warning can minimize or alleviate such disastrous damages by advanced precautions and decision makings. This could certainly be a noble help to minimize the adverse effect of weather variability. Hence it was observed that an accurate and enhanced forecasting method for weather variability can be an important research initiative.

In data science analytics, there may be numerous hidden rules/patterns available in a massive amount of meteorological dataset. These patterns/rules can be very evident but cannot easily be discovered by domain experts. For this reason, numerous researchers tried to predict weather variability with different meteorological parameters through different data mining techniques. When we use different predictive models with low performance in decision-making processes; definitely decisions cannot be effective on-ground in reality. From the researcher's point of view, we need to have improved models so as to enhance the accuracy of the weather variability forecasting. These enhanced or improved models can provide better and accurate weather variability forecasting to support decision-making processes.

The Data Mining techniques have proved encouraging results in the prediction of weather variability but the variability of the parameters could lead to a deficiency in the results. In this study, efforts are made to enhance the accuracy of weather variability forecasting with extended parameters. The different types of weather parameters such as rainfall, wind speed, humidity, dew point, and sunshine, etc. were added as additional parameters along with the parameters used by the prior researches. It was done to check the impact on the accuracy of the weather variability forecasting. This the research applies "Multilayer perceptron or Back Propagation, Naïve Bayes, and Multinomial logistic regression algorithms". In this research, we used WEKA software for data analysis.

Initially, it was observed that the Dodota Woreda of Ethiopia is being affected by serious weather variability disasters such as socio-economic crisis, transport blockage, cattle deaths, illness, and food shortage because of lack of accurate weather forecasting. This results in a displacement of the people from one region to another. This phenomenon needs serious attention. The significant research question raised to initiate this research study were: 1) What are the major challenges in the accurate forecasting of the weather variability in general and at Dodota Woreda region as a case? 2) What are the forecasting accuracy gaps in the prior research studies in the 1) region? and 3) How a better research contribution can enhance the forecasting accuracy of the weather variability so as to support the decision-makers for better preparation to alleviate the adverse impacts. The main motivation of this research was to design an enhanced predictive model that can accurately forecast the weather variability using data mining techniques. Finally, the designed predictive model and the hidden knowledge/pattern discovered from this research are presented. This knowledge as rules/patterns discovered from the massive meteorology dataset can help agencies to provide better weather variability information to the public and decision-makers.

The study expected a significant contribution to the new knowledge domain of the weather forecasting in the region. Hence, this study is an attempt to analyze the massive amount of meteorological data to design an enhanced predictive model that can support the forecasting of weather variability using data mining techniques. The study was focused on the Dodota Woreda region located in the Arsi zone in the Oromia Region of Ethiopia.

## II. REVIEW OF LITERATURE

A massive amount of data is essential to generate new information and knowledge. The data retrieval is not enough to hold and process relevant data stored in databases and other repositories. It requires alternative techniques such as data mining. These techniques facilitate several powerful tools for classification, analysis, interpretation, and modeling of the massive amount of data that could aid or support in decision-making processes. Different data mining techniques and algorithms are being used to predict or forecast the weather variability. This can help the organizations, agencies, and individuals to take knowledge-driven proactive precautionary decisions.

The weather is the current state of the atmosphere around us and characterized by temperature, atmospheric pressure, humidity, wind speed, cloudiness, and others. The Climate is the state of the atmosphere over a long period of time, such as over years, decades, centuries, or greater. Weather variability forecasting is the application of data science and technology to estimate the future whether the conditions will be the same or not such as in an hour, tomorrow, next week, or next month. Accurate weather variability forecasts are important for planning day-to-day activities. Farmers need the information to help them plan for the planting and harvesting of their crops. Airlines need to know about the local weather conditions in order to schedule flights. Weather forecasting helps to make a more informed daily decision and may even help (keep) us from danger [6] [7].

Forewarnings and forecasting services that provide information about weather impacts are expressed in contributions focusing on very different national and user contexts [8]. The weather forecasts are divided into the many categories. Now-casting, Short-range forecast, Medium range forecasts (4 to 10 days), and Long-range forecasts (more than 10 days to season): there is no rigid definition for long-range forecasting, which may range from a monthly to a seasonal forecast. This type of forecast is good for the long decision-making process and in the agricultural product that is based on seasonal weather conditions. Especially for the farmers to plant and harvest their crops, they need to know the weather variability of more than 15 days to a month or maybe the season [9].

Weather variability forecasting is the application of science and technology to predict the state of the atmosphere for a given location. Ancient weather forecasting methods usually relied on observed patterns of events, also termed pattern recognition [10]. Ancient times it might be observed that if the sunset is particularly red, the following day often brought fair weather. However, not all of these predictions were proved reliable. The data mining techniques developed recently can successfully be applied for accurate prediction of the weather variability. Data mining is a process that uses a variety of data analysis tools to discover patterns and relationships in data that may be used to make valid predictions. Some of the applications of data mining include the discovery of interesting patterns, clustering of data based on parameters, and prediction of results using the existing data [11]. Data are raw facts or measurements that can be recorded about events, assets, and

have no significance beyond its existence [12]. The data needs to be collected, organized, summarized and analyzed for decision-making purposes. Information is the processed data that has a certain meaning, and the knowledge is the appropriate discovery of hidden facts in the collected and processed data or information in a contextualized form. Usages of computers and storage technologies in different sectors have dramatically increased the availability of digital data in every sector. Like other sectors data captured in the Meteorology, sector has also increased dramatically, and thus huge amounts of datasets are the basic input that administrators and decision-makers use to build theories and models. But the analysis capability to build patterns and model from these available datasets are very slow when compared with the availability of the datasets increasing exponentially. Due to the huge sum of data collected/ recorded by the meteorology offices, the application of traditional tools become unreasonable to discover the hidden pattern in terms of new knowledge that can be helpful in the decision-making processes. Data does not replace skilled business analysts or managers, but rather gives them powerful tools and techniques to improve the job they do.

Data mining is defined as the application of algorithms for discovering hidden patterns and relationships in the variables of data using a variety of data analysis tools to make valid predictions. The data mining objective is to provide accurate knowledge in the form of useful rules, techniques, visual graphs, and models for weather variability forecasting [13] [14]. The interesting patterns are selected based on subjective or objective problems. This knowledge can be used to support decision-making processes in specified sectors such as agriculture, weather, and irrigation. Data mining itself is not able to automatically derive useful knowledge from a vast amount of dataset without machine-oriented guidance. This implies that collecting, processing, exploring, and selecting a suitable tool and technique is critically an important issue. The Predictive model works for making predictions from datasets. The predictive model reveals the unknown patterns as a result of different datasets to predict the future. [15] [16].

The classification is one of the predictive data mining tasks and used to find a model that describes and distinguishes classes or concepts.

The time series is a sequence of events where the next event is determined by one or more of the preceding events. The weather parameters such as wind speed, rainfall, relative humidity, sunshine, dew point, maximum temperature, minimum temperature, and months are used to predict the weather condition such as sunny, cloudy, and rainy and depicted in Fig. 1.

Data mining is a complex process which needs a proper combination of various data mining tools, techniques, and human experts. The process model can help experts to have a common reference and will increase the understanding of complex data mining issues. Modeling of data mining is used to select the best data mining techniques and methodology that fits with the expected patterns and features in data mining applications. Different data mining process models available are: CRISP-DM (Cross-industry Standard Process for Data Mining), SEMMA (Sample, Explore, Modify, Model, and

Assess), KDD (Knowledge Discovery in Database), and Hybrid model.

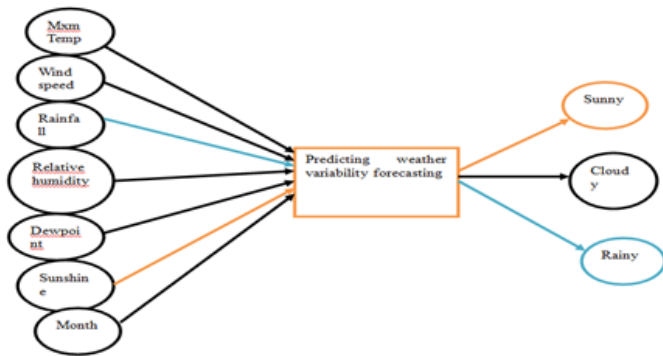


Fig. 1. Conceptual Frameworks of the Research Study.

Some of the studies accomplished prior to this research are reviewed and critically analyzed in this paper. A research conducted by Nishchala et al. [17] for Classification and Forecasting of weather using ANN, k-NN, and Naïve Bayes Algorithms is reviewed and analyzed in this paper. This study conducted a comparison between MLP, k-NN, and Naïve Bayes algorithms to predict and classify the future conditions of the weather. It shows that k-NN provides better accuracy in classification and also in terms of the execution time required. The numeric prediction results also showed that the Naïve Bayes gives better results compared to k-NN and MLP. From the above results, a hybrid system can be developed through Naïve Bayes for numeric prediction and k-NN for classification. This model can provide better accuracy because a single system cannot satisfy all the constraints.

Research of Abhishek Saxena et al. [18] presented the review of weather prediction using Artificial Neural Networks. It yields better results and can be considered as an alternative to the traditional meteorological approach. The study expressed the capability of artificial neural network in predicting various weather phenomena such as temperature, thunderstorms, rainfall, wind speed and concluded that major architecture like BP, MLP is suitable to predict the weather phenomenon.

Another similar research was carried out in Dire Dawa, Dereje region of Ethiopia on Meteorological Data Analysis for designing a Predictive Model to Support Weather Forecasting using Data Mining techniques [5]. This research recommended filling the gaps in the cumulative knowledge in the area. This study used only three parameters i.e. maximum temperature, minimum temperature, and relative humidity which is not sufficient for the promising accuracy of the forecasting results. The research also recommends enhancing the forecasting accuracy by adding some additional parameters like rainfall, wind speed, and dew point. This research used Artificial Neural Network and decision tree algorithms to achieve encouraging results. Another research [19] concluded that the decision tree has limitations with a continuous variable or with complex variables. The decision tree uses the “divide and conquer” method, and therefore it can perform better if the relevant parameters exist in a limited number. If too many complex interactions are present then it can perform poorly. This research recommended testing the alternative

classification algorithms. These algorithms can be tested to investigate their applicability to the problem domain. Therefore, this research proposed to focus on the four parameters i.e. rainfall, wind speed, sunshine, and dew point with three different algorithms like Naïve Bayes, MLP, and multinomial logistic regression.

### III. RESEARCH METHODOLOGY

This research study is a design science, experimental research which mainly focuses on the quantitative approach. To achieve the desired goal, the study followed KDD (Knowledge Discovery in Databases) modeling approach.

CRISP is an industry-standard process while KDD is most preferable for academic purposes. KDD is selected for three main reasons: 1) KDD is the best suited for academic purposes; 2) KDD reduces the skill required for knowledge discovery to the non-experts [20], and 3) KDD is independent of any tool and technique and therefore any technique can be used. Hence, the steps of the KDD process followed are data understanding and data selection, preprocessing, transformation, data mining, evaluation, and interpretation of the mined data to discover the new knowledge or pattern. The five steps involved in the KDD process are [21] presented in Fig. 2.

#### A. Data Sampling, Collection, and Understanding

In this study, datasets were collected from primary and secondary sources as presented in Fig. 3. After the interpretation, a weather variability forecasting model was designed and demonstrated to the domain experts. The potential source of data was the national meteorology agency at Dodota Woreda and Awash Melkasa station. To design a model for weather variability forecasting, 10 years of daily datasets from 2006 to 2016 were collected from national meteorology agencies. To extract new knowledge about weather variability from these large datasets, the researchers used data mining techniques and tools.

The obtained data record includes ‘Ghid’, ‘Name’, ‘latitude’, ‘longitude’, year, month, date, time, wind speed, Relative Humidity, Maximum temperature, Minimum temperature, sunshine, rainfall & dew point, evaporation, and precipitation. This step of the study contains, what research strategies should be set to utilize, the selection of the study outline and strategy of data collection, management, and analysis.

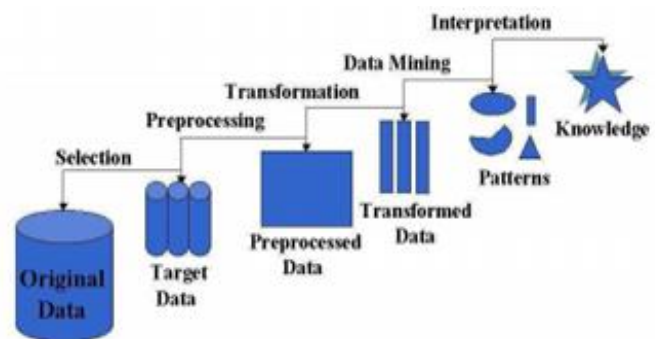


Fig. 2. Steps in KDD and Methodology of the Research.

To understand the nature of the data, descriptive statistics techniques were applied. This step includes collecting sample data and deciding which data, including format and size, will be needed. Background knowledge can be used to guide these efforts. Data were checked for completeness, redundancy, missing values, and the plausibility of attribute values, etc. To discover the profound knowledge from the data, the researchers worked closely with the domain experts. The original meteorological datasets contain 20 parameters and 8967 records. From these collected datasets, the study used only 7 parameters and 5282 records. The data collection methods and procedures are presented in Fig. 3.

### B. Preprocessing

Data preprocessing is an important and time-consuming phase in the knowledge discovery process. It must be taken into consideration with care in the mining of data streams. Poor quality of data is the main challenge for the knowledge discovery process. Thus, to get accurate results, the researcher preprocessed the data carefully to remove the noisy and unwanted climatic parameters.

### C. Attribute/Feature selection

Feature selection is the process of identifying and removing the irrelevant and redundant data or information [22]. Not all parameters can be relevant therefore some of the irrelevant parameters which are not significant to the study's objectives are excluded. For instance, the station name, the latitude, and longitude of the area, Elevation, year, day, time, station name and etc. are not important and do not provide any importance in the results of the data mining process. The parameters selected using the ranker tool are month (mz), \_dew point (dp), \_sunshine (ss), \_maximum temperature (mxt), \_relative humidity (rh), \_wind speed (ws) and rainfall (rf). Fig. 4 illustrates the ranking order of attribute-based on the relevance of the parameters.

It evaluates the worth of parameters by measuring the information gain with respect to the class. The result of parameters, selection shows that the dew point (dp) has gained the highest information gain with 1.5342, and the lowest minimum temperature (mint) 0 information gain, so the minimum temperature is removed from the dataset since it has the least information gain. Fig. 4 presents the order generated by WEKA for the parameters based on their information gain. This implies that how the given parameters are related to the predicted classes (sunny, rainy, and cloudy) or used to determine the predicted class based on their information about the classes. And also, we discussed with domain experts to identify which parameters are more determinant in weather variability. The domain experts also validated that these parameters are the major parameters that can influence weather variability forecasting.

Data cleaning process is used in this paper which refers to the pre-processing of data to remove or reduce noises and the treatment of missing values. Handling missing values by appropriate method does not affect the quality of the data. The researchers worked with weather data was in a form of time series, so they must preserve the series smoothness and consistency. In this study; the average mean value carried forward (manually), k-mean fuzzy logic was used to handle

these missing values. Ignoring instance techniques was also used in this research. Carried forward is an effective method to fill the missing values in the case of time series where the missing value is strongly related to its previous and next value. Replacing missing values with an interpolated estimate or Mean of nearby points (stations) replaces missing values with the mean of surrounding stations of weather datasets. This structure can be exploited by interpolating the missing value. This approach is very effective when it is appropriate, usually with time-series data. This method for handling missing values is used based on the need when filling missing values in datasets. In this research, the most appropriate and effective method is an average mean to analyze the numeric datasets. The statistical summary of the parameters has shown that there are missing values and an unbalanced occurrence of instances. Missing values of the selected parameters can be handled by different methods to fill the missing values, but for evaporation with 66% and Precipitation with 78% is difficult to handle its missing values. For this reason, the researcher excluded these parameters.

Once the data is cleaned, it is needed to transform or consolidate data into suitable forms for applying to mine strategies and transferred into a data mining capable format such as attribute construction, aggregation, and discretization. Therefore, the values of parameters are changed to a new set of replacement values to ease data mining.

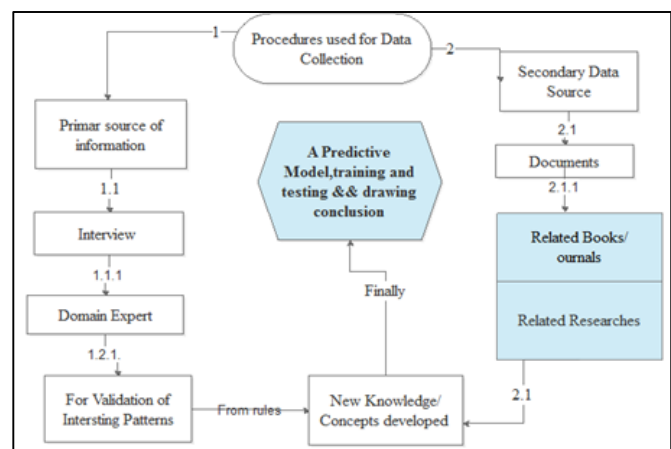


Fig. 3. Data Collection Methods and Procedures.

```
=== Attribute Selection on all input data ===  
Search Method:  
Attribute ranking.  
Attribute Evaluator (supervised, Class (nominal): 9 wc):  
Information Gain Ranking Filter  
Ranked attributes:  
1.5342 8 dp  
1.0439 4 ss  
0.8641 1 mz  
0.0869 7 ws  
0.0517 2 rf  
0.0213 5 mxt  
0.0137 3 Rh  
0 6 mint  
Selected attributes: 8,4,1,7,2,5,3,6 : 8
```

Fig. 4. Parameters Ranking.

Discretization was used to obtain a reduced representation of the data while minimizing the loss of information content [23]. Most of the time in developing a predictive model, several researchers preferred discrete values than continued values. Data mining phase is engaged in searching for patterns of interest in a particular representational form, depending on the DM techniques. From DM techniques, the classification was selected in this study. However, the classification is used because of the nature of parameters and the aim is predicting classification labels. Since the goal of the study is to classify the weather variability forecasting to sunny, cloudy, and rainy target classes based on the weather parameters. In this research study, data mining is used to extract the hidden knowledge from the massive amount of datasets collected from the meteorology agency of Dodota Woreda and Awash Melkasa station and used to generate patterns using the most outperforming algorithm. Interesting patterns are identified by the domain experts from the generated rules to design weather variability, forecasting models. Starting from data collection up to designing and developing a model with a demo through prototype, the research study used a step by step procedure. The detailed flow of basic steps for designing the Model is presented in Fig. 5.

#### D. Interpretation /Evaluation of the Discovered Knowledge

After mining the required pattern; the interpretation and evaluation of the mined patterns were done. The interpretation is concerned that whether the discovered pattern is interesting or not and verifies as knowledge or not. Comparatively measuring the performance of each classifier and representing the results in a suitable model. The performance of the classifiers adopted in the study are measured and evaluated based on their accuracy, TP rate, recall, and precision. The model which has the highest accuracy rate is selected. The rules generated from the selected outperforming algorithm are validated by the domain experts. Finally, visualization and knowledge representation are used to present the mined knowledge to the users and stored as new knowledge in the knowledge base of the weather forecasting division.

#### E. Data Mining Tool Selection

WEKA 3.9.0 machine learning software was selected based on selected suitability parameters after critical assessment as presented in Fig. 6. The predictive data mining task is used to predict the weather variability for the future based on some climatic parameters. This research study used long-range forecasting as it is used to forecast weather variability for a month.

Selection of the most suitable data mining algorithm depends on the structure of the datasets available, performance of the algorithm and objective of the study. Having logical observation from the available algorithms in the WEKA machine learning software tool, the MLP, Naive Bayes, and multinomial logistic regression algorithms were selected as classification algorithms for this research. The back-propagation learning algorithm is one of the most important developments in neural networks. This network is still the most popular and most effective model for complex, multi-layered networks. The typical back-propagation network contains an input layer, an output layer, and at least one hidden layer. The

number of neurons at each layer and the number of hidden layers determine the network's ability on producing accurate outputs for a particular dataset [24]. Back propagation learns by iteratively processing a dataset of training tuples, comparing the network's prediction for each tuple with the actual known target value. The target value may be the known class label of the training tuple (for classification problems) or a continuous value (for numeric prediction). The neural network refers to the set of connected input or output units in which every connection has a weight associated with it. The idea of the back-propagation algorithm is to reduce the error until the ANN learns the training data. The training begins with random weights, and the goal is to adjust them so that the error can be minimized. Artificial Neural network with Back propagation algorithm seems to be the most appropriate method for forecasting weather variability with better accuracy [25].

The multilayer perceptron is one of the most widely used problem-solving architectures in a great variety of areas. In addition, it is easy to use and apply. Multilayer Perceptron is known as back propagation (backward error propagation), is an extension to networks with intermediate layers (multilayer networks). MLP is a classifier that uses back propagation to classify instances. Like any other learning scheme, a multilayer perceptron trained with back propagation may suffer from over fitting [26]. Especially if the network is much larger than what is actually necessary to represent the structure of the underlying learning problem.

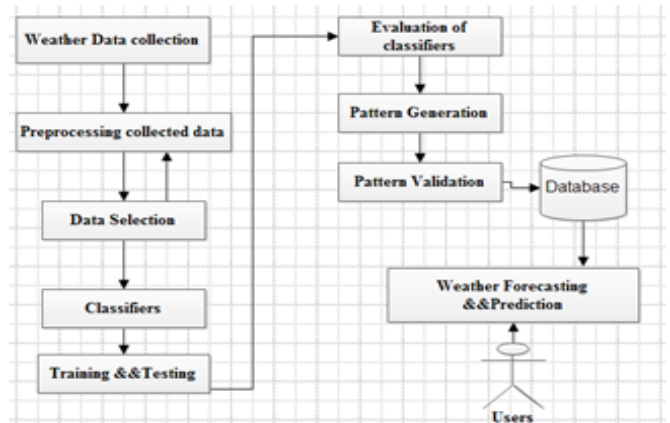


Fig. 5. Basic Flow for Designing Weather Variability Forecasting Model.

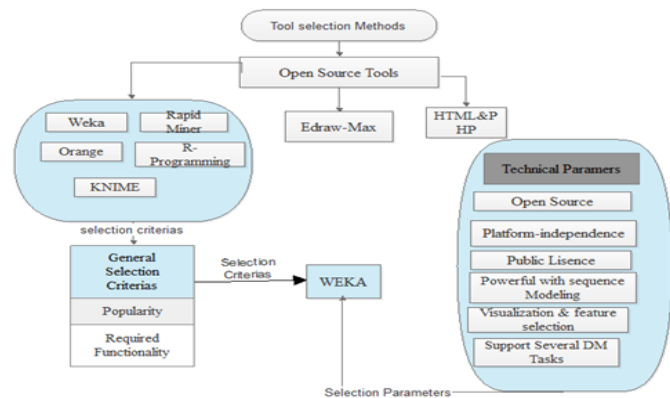


Fig. 6. Tool Selections for Data Analysis, Modeling, Testing and Training.

An advantage of the naive Bayes classifier is that it requires a small amount of training time to estimate the parameters necessary for classification [27].

It performs better in many complex real-world situations like Spam Classification, Medical Diagnosis, and Weather forecasting [27]. This algorithm is important for several reasons. It is very easy to construct, not needing any complicated iterative parameter estimation schemes. This means it may be readily applied to huge datasets. It is easy to interpret, so users unskilled in classifier technology can easily understand why it is making the classification. For the prediction of future events, Naïve Bayes uses knowledge of prior events. Suppose we have more than one evidence for building our Naïve Bayes model, we could run into a problem of dependencies, i.e., some evidence may depend on one or more of other evidence. For instance, the evidence “dark cloud” directly depends on the evidence “high humidity”. However, including dependencies into the model will make it very complicated. This is because one evidence could depend on many other pieces of evidence.

Logistic regression analysis studies the association between a categorical dependent variable and a set of independent (explanatory) variables. The name logistic regression is used when the dependent variable has only two values, such as 0 and 1 or Yes and No. Multinomial logistic regression is used to predict categorical placement in or the probability of category membership on a dependent variable based on multiple independent variables. The binary logistic regression, multinomial logistic regression uses maximum likelihood estimation to evaluate the probability of categorical membership. Logistic regression competes with discriminant analysis as a method for analyzing categorical-response variables. This program computes binary logistic regression and multinomial logistic regression on both numeric and categorical independent variables. It reports on the regression equation as well as the goodness of fit, odds ratios, confidence limits, likelihood, and deviance. It can perform an independent variable subset selection search, looking for the best regression model with the fewest independent variables. It provides confidence intervals on predicted values and provides ROC curves to help determine the best cutoff point for classification. It allows validating the results by automatically classifying rows that are not used during the analysis.

The classifiers were evaluated by cross-validation using the number of folds. K-fold is a natural number used to check the performance of the model through k-times. In this paper 20-folds, 15-folds and 10-fold cross-validation was recommended for estimating accuracy and achieved the highest accuracy in experimentations. The classifiers were tested for 40%, 50%, 66% (the default), 76%, 86%, and 96% for training split percentages and the remaining for testing and different values of learning rates and a number of hidden layers were also used.

#### F. Methods of Analysis and Evaluation of System Performance

Once a model is built using training data, one can be curious to know how the model will perform in the future or compare the forecasting accuracy of multiple models for the same forecasting problem and to decide on the real-world

decision making. To do this we need to measure the performance and accuracy of the model. Commonly, predictive models with higher accuracy are viewed as better. The following evaluation metric was used.

1) *Precision and Recall*- precision is the fraction of true positive from the predicted instance, while recall is the fraction of relevant instances that are retrieved. Precision can be thought of as a measure of exactness or quality, whereas recall is a measure of completeness or quantity [28].

2) *ROC graphs* are two-dimensional graphs in which true positive (TP) rate is plotted on the Y-axis and false positive (FP) rate is plotted on the X-axis.

To test which classifier is highly significant for a given subject is determined by ROC curve analysis. The following Fig. 7 shows the performance of the three classifiers (MLP, Multinomial logistic regression).

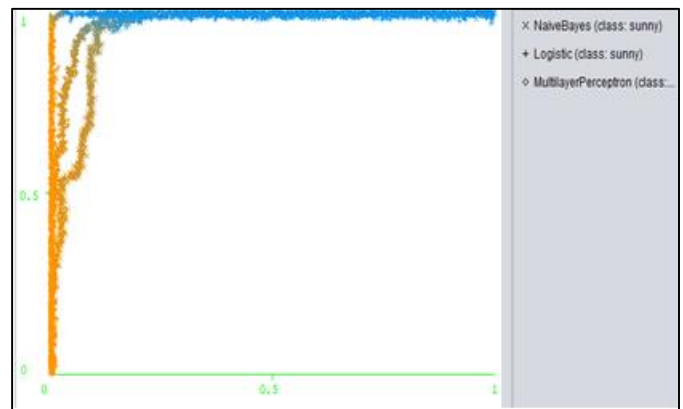


Fig. 7. ROC Area Curves (Performance of the Three Classifiers).

For perfect forecasts (better classifier performance) accuracy gets a value of 1, the maximum possible value. Forecast with little or no skill will obtain a ROC score of approximately 0.5, the area under the diagonal.

#### IV. EXPERIMENT, DISCUSSION AND EVALUATION

This research used different testing methods with different values to check that which testing value provides the highest forecasting accuracy. This was done for all the selected numbers of tests and to select the best models from two or more trained models. The accuracy of the classifier was measured by true positive rate, false-positive rate, F-measure, Recall, Precision, and additionally by the ROC curve.

To avoid the effect of data imbalance on the model created, WEKA based resampling technique was applied. These methods can be grouped into two categories: data perspective and algorithm perspective. The re-sampling method is attractive under the most imbalanced circumstances. This is just because the re-sampling adjusts only the original training datasets instead of modifying the learning algorithms; therefore, it provides a convenient and effective way to deal with imbalanced learning problems using standard classifiers [29]. To test the performance of the models, the researcher used six number of the percentage split starting from 40 to 96. Since it ranges from 1 to 100 excluding the two extreme



borders, different numbers of k-folds from 2-folds to 20-folds, hidden layers from single to multiple hidden layers, and learning rates from (0.01 to 0.9).

**A. Model Designing using Naïve Bayes Classifier**

The Naïve Bayes algorithm used in this experiment was tested with all the attributes. The selected parameters were used to find a better- classification algorithm for the datasets. From the experiment output, the accuracy of the classifier achieved was 84.0212% with an error rate of 15.9788 %. This implies that there are no differences in accuracy for both selected and with all parameters. The Naïve Bayes tested with different split percentages and achieved the highest accuracy rate (84.57%) at 76% split percentage and depicted in Fig. 8.

As presented in Fig. 9, even though the performance variations among different k-values are minimal i.e. nearly 84.0212% successes, the highest performance was observed in 20-fold.

**B. Model Designing using Multilayer Perceptron**

The second algorithm used to classify and discover patterns of factors affecting weather variability was the Multilayer Perceptron. In the Multilayer Perceptron case also, the different testing configurations were tested. In this experiment, multilayer-perceptron used all the parameters with cross-validation techniques with default parameters. The classifier correctly classified 5194 instances (98.3908%) and the true positive rate was 98.4%.

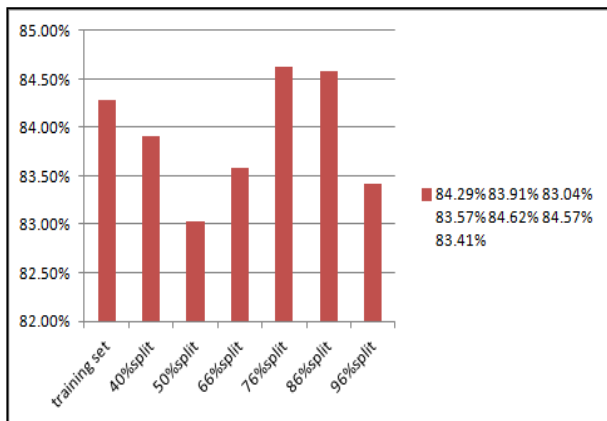


Fig. 8. Output of Naïve Bayes Classifier with different Test Split Percentage Naïve Bayes with different K-Folds.

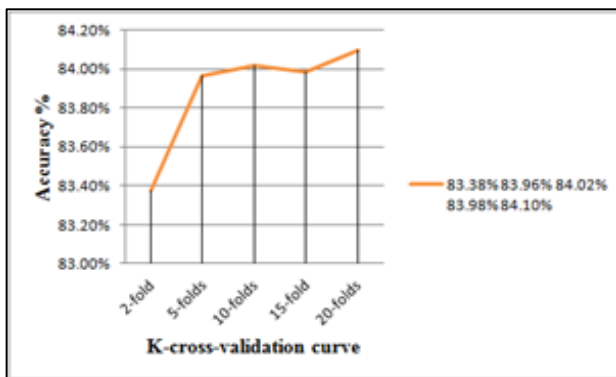


Fig. 9. Learning Progress Curve.

As indicated in Fig. 10 and 11, the performance variations was observed among different models. The accuracy of the Multilayer Perceptron model was found 98.4286%, the highest performance was observed in 15-folds. On the other hand, the performance of MLP was found 99.21% with a percentage split at 76%.

As it is observed from Table I, the highest accuracy rate is achieved for a single hidden layer at value 10 configuration (98.1257 % of instances). The worst accuracy rate is obtained for value 10, 9, 8, and 7 (multiple hidden layer) which is 93.7145% accuracy rate. For learning rates, the highest accuracy was observed at 0.1(98.334% accuracy rate) and the worst accuracy rate was observed at 0.9(95.6835% accuracy rate). In general, the accuracy rate is decreased as the number of hidden layers and the learning rate is increased.

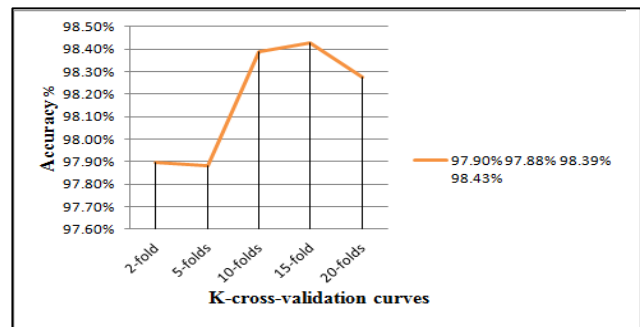


Fig. 10. Output of MLP with different K-Folds-Cross-Validation.

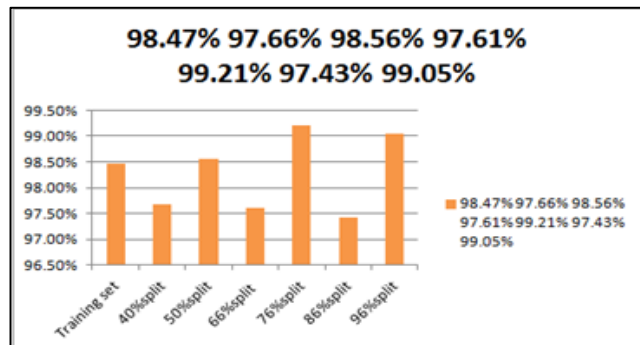


Fig. 11. Output of MLP with different Split Percentage.

TABLE I. THE OUTPUT OF MULTILAYER PERCEPTRON WITH DIFFERENT HIDDEN LAYERS AND LEARNING RATES

| No. | Number of the hidden layers | Accuracy % |
|-----|-----------------------------|------------|
| 1   | 10(single hidden layer) (1) | 98.1257 %  |
| 2   | 10,9(2)                     | 97.8228 %  |
| 3   | 10,9,8(3)                   | 96.4786 %  |
| 4   | 10,9,8,7(4)                 | 93.7145 %  |
|     | Learning rate               |            |
| 1   | 0.01                        | 97.7849 %  |
| 2   | 0.05                        | 98.2204 %  |
| 3   | 0.1                         | 98.334 %   |
| 4   | 0.3                         | 98.1257 %  |
| 5   | 0.5                         | 97.6903 %  |
| 6   | 0.9                         | 95.6835 %  |

C. Model Designing using Logistic Regression Classifier

The third multinomial regression algorithm used in this experiment was tested with all parameters, selected parameters, different split percentage, and K-folds for finding better classification algorithm for the datasets as presented in Fig. 12, and 13.

D. Comparison of the Models

Table II indicates the performance summary of the experimental results for all experiments. As it was observed in Table II, the Naïve Bayes, Logistic Regression, and MLP algorithms compared, and MLP algorithm is selected to generate patterns/rules that classify the weather condition on a particular day such as sunny, rainy, or cloudy. One of the research questions of this study was to find the most suitable data mining algorithm that outperformed in classification. We compared the algorithms used in this study and selected the one which performs the best. In all the experiments, the same datasets were used for all the algorithms (Naïve Bayes, Logistics, and MLP).



Fig. 14. ROC curves for all algorithms

In the comparison of all algorithms using ROC area, the MLP provides the best result in terms of classification accuracy and high ROC area. This is presented in Fig. 14. Thus because of this MLP is selected as a best-fit model for this research study.

E. Expert Validation

The results validation was done using expert interview based on the revealed results. The subjective questions were asked to the experts to understand and forward their expert opinions as a new knowledge contribution with enhanced accuracy in forecasting weather conditions. The experts ensured the weather variability results through comparisons of weather conditions in reality and the forecasted. The first question raised was: 1) which types of additional and parameters can be considered for enhancing the forecasting performance of the model? The expert response, research outcome, and the generated rules clearly indicated that the major parameters that can influence the weather variability are dew point, month, relative humidity, rainfall, maximum temperature, sunshine, and wind speed. The second subjective question raised was: 2) what types of most interesting rules or patterns can be generated for subjective and objective measures? The response of the experts using the selected determinant parameters; the interesting patterns were explicitly identified to predict weather variability. Hence, the explicit interesting pattern can be delivered to the domain experts. These interesting rules for subjective and objective measures are as follows:

If month=m5 and 9<dp>=17 and rh=h and mxt=ht and ws=lir and rf=lr and ss=6 then wc=rainy.

If month=m5 and 9<dp>=17 and rh=h and mxt=ht and ws=lir and rf=lr and ss=6 then wc=rainy.

If month=m5 and 9<dp>=17 and rh=h and mxt=ht and ws=lir and rf=lr and ss=6 then wc=rainy.

If month=m5 and 9<dp>=17 and rh=h and mxt=ht and ws=lir and rf=lr and ss=6 then wc=rainy.

Another subjective question was raised was: 3) which data mining algorithm can be a most suited and provides the enhanced/improved results? The responses of the subject experts were confined to the varied algorithms such as classification algorithms: Naïve Bayes, Multinomial Logistic Regression and MLP. From these algorithms, MLP classifier

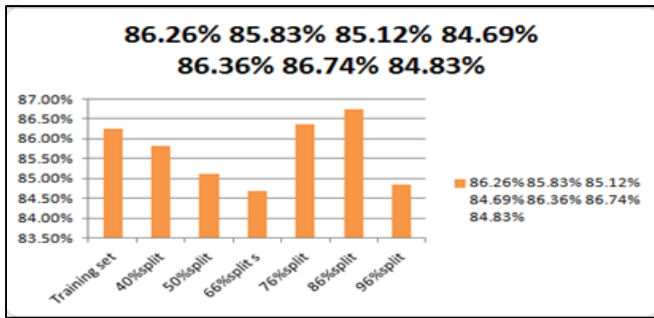


Fig. 12. Output of logistic algorithm with different percentage split

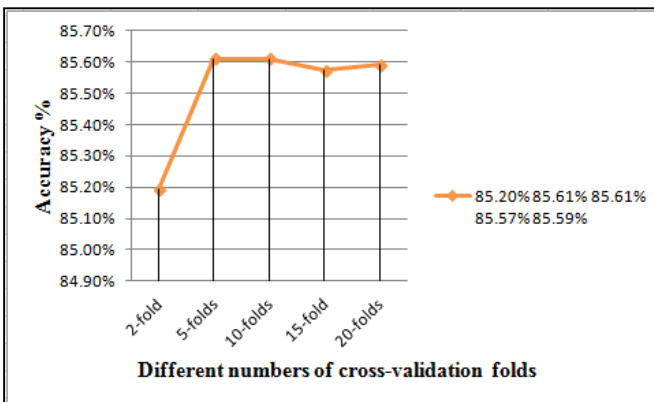


Fig. 13. Output of logistic algorithm with different k-folds cross-validation

TABLE II. THE PERFORMANCE SUMMARY OF THE EXPERIMENTAL RESULTS FOR ALL EXPERIMENTS

| Model    | Accuracy  | TP rate | FP rate | precision | Recall | F-measure | ROC area |
|----------|-----------|---------|---------|-----------|--------|-----------|----------|
| NB       | 84.6215 % | 0.846   | 0.082   | 0.845     | 0.846  | 0.842     | 0.916    |
| MLP      | 99.2114 % | 0.992   | 0.006   | 0.992     | 0.992  | 0.992     | 1.000    |
| Logistic | 86.7388 % | 0.867   | 0.066   | 0.866     | 0.867  | 0.866     | 0.948    |

algorithm with selected parameters achieved relatively better classification accuracy as compared to other two algorithms for weather variability forecasting.

## V. CONCLUSION

Applications of data mining techniques have been increasingly getting popularity and proved to be relevant for the sectors like meteorology, health care, telecommunications, and banking etc. In particular, the meteorology sector has significant possibilities where data mining can be applied for weather variability forecasting to improve the accuracy of weather forecasting and to support decision-makers towards better disaster management. This research study is an attempt to design a forecasting model for the weather variability using data mining techniques in general and Dodota Woreda region as a case. At the end of the study, data mining model was developed using three data mining classification algorithms. In this model, the classifier extracts the hidden knowledge from ten years of massive amounts of meteorological dataset. The algorithm was trained, tested, evaluated against a test dataset. The model was also validated using cost-sensitive evaluation method, ROC curve of the classifier, precision, recall and f-measure methods. In this model was evaluated for the effectiveness. KDD process model was used in this study. WEKA was used as a tool for preprocessing and analyzing the weather variability datasets for the forecasting. Three different models were designed using Naïve Bayes, Multinomial Logistic Regression, and MLP algorithms through adjustment settings to come up with understandable and meaningful results. The comparison of the results produced by all the models showed that the encouraging results obtained with MLP classifier is the most appropriate for classification and weather variability forecasting. MLP model with selected parameters presents an interesting (the highest) forecasting accuracy in the results. MLP classification algorithm can be used to generate rules to classify weather variability parameters such as maximum temperature, minimum temperature, relative humidity, dew point, rainfall, and sunshine to sunny, cloudy, and rainy target classes. Findings of this research revealed that these parameters have strong positive relationship with weather forecasting in meteorology sectors. Also, the detailed discussion on the discovered patterns was done with domain experts. It was proved that this research study can be a significant support to the meteorologists in decision making processes. Also, the improvement in the accuracy of weather variability forecasting can play a vital role in taking better precautionary measures. The outcome of this research can be used by meteorology agencies to help the meteorologist to make consistent forecasts, to support agricultural institutions and other business organizations.

Overall, the major contribution of this research is to find out the relationships among variables (parameters) that bring weather variability, data mining techniques offer great capacity in supporting meteorology agency towards making better decisions. The extracted knowledge generated rules (patterns) by the model can be a significant new knowledge-based input for improving the decision-making processes.

## VI. RECOMMENDATIONS

In due consideration of the results of the research, the following recommendations are forwarded for future researches in the weather variability forecasting domain:

- In this study, the forecasting model used seven climatic parameters such as rainfall, maximum temperature, minimum temperature, relative humidity, dew point, sunshine, and wind speed. However, some other important parameters like precipitation, evaporation, environmental/sea level pressure, wind direction, cloud coverage, and much more remains and needs to be considered for further study to check the impacts on the model's accuracy. In data collection, the researcher used only one station; however, the study can focus on expanding the scope for other weather-sensitive stations.
- Even though there are many data mining techniques, but this research used only the three classification algorithms such as Naïve Bayes, Multinomial Logistic Regression, and MLP. The other data mining algorithms like the K-NN algorithm, Vector Machine, Weighted Bayesian, are overlooked for testing, and forecasting in this research study. It might be important to design a forecasting model for weather variability forecasting using these different algorithms too.

## ACKNOWLEDGMENTS

It is great pleasure for us to express our heartfelt gratitude to the National Meteorology Agency and Awash Melkasa Meteorology station, Ethiopia for providing the relevant datasets required for this research. Also, authors express their gratitude to the professional scientists and the staff of the station for their consistent support and inputs to complete this study on time. Special thanks to our expert adviser Prof. DP Sharma for his consistent technical support and guidance despite of his hectic engagements in multidisciplinary assignments.

## REFERENCES

- [1] M Ramzan Talib, Toseef Ullah, M Umer Sarwar, M Kashif Hanif and Nafees Ayub, "Application of Data Mining Techniques in Weather Data Analysis," IJCSNS International Journal of Computer Science and Network Security, vol. 17, no. 6, pp. 22-28, June 2017.
- [2] D. Santhi Jeslet, S. Jeevanandham, "Climate Change Analysis Using Data Mining Techniques," International Journal of Advance Research In Science And Engineering, (IJARSE), vol. 4, no. 03, pp. 46-53, March 2015.
- [3] S. Kotsiantis, A. Kostoulas, S. Lykoudis, A. Argiriou, K. Menagias, "Using Data mining Techniques for estimating minimum, maximum and average daily temperature values," International journal of Mathematical, Physical and engineering science, vol. 1, no. 1307-7465, pp. 16-20, January 2008.
- [4] Y. Radhika and M. Shashi, "Atmospheric Temperature Prediction using Support Vector Machines," International Journal of Computer Theory and Engineering, vol. 1, no. 1, pp. 55-58, April 2009.
- [5] M. E. Dereje, "Meteorological Data analysis for creating predictive model that support weather forecasting using data mining techniques," p. 4, 15 July 2014.

- [6] Kristine Inchausti, Michele McLeod, Stacie Pierpoint, "Weather-forecasting," Annenberg Foundation, 23 July 2016. [Online]. Available: <https://www.learner.org/exhibits/weather/forecasting.html>. [Accessed 7 December 2017].
- [7] Kapil Khandelwal, Durga Prasad Sharma, "Hybrid Reasoning Model for Strengthening the problem solving capability of Expert Systems," (IJACSA) International Journal of Advanced Computer Science and Applications, vol. 4, no. 10, pp. 88-94, 2013.
- [8] Andrea Taylor, Thomas Kox, David Johnston, "Communicating High Impact Weather: Improving warnings and decision making processes," International Journal of Disaster Risk Reduction, vol. 30, no. 4, pp. 18-24, 2018.
- [9] P. V. B. N. Ganesh P. Gaikwad, "Different Rainfall Prediction Models And General Data Mining Rainfall Prediction Model," vol. 2, no. 7, July - 2013.
- [10] Nevonproject, "weather forecasting using data mining," NevonProjects, 12 March 2012. [Online]. Available: <http://nevonprojects.com/weather-forecasting-using-data-mining/>. [Accessed 11 December 2017].
- [11] Alazar Baharu, Durga Prasad Sharma, "Performance Metrics for Decision Support in Big Data vs. Traditional RDBMS Tools & Technologies," (IJACSA) International Journal of Advanced Computer Science and Applications, vol. 7, no. 11, pp. 222-228, 2016.
- [12] M. K. KELEŞ, "An Overview: The Impact of Data Mining Applications on Various Sectors," Technical Journal , Vol. 3, No. 11, pp. 128-132, 2017.
- [13] M Ramzan Talib, Toseef Ullah, M Umer Sarwar, M Kashif Hanif and Nafees Ayub, "Application of Data Mining Techniques in Weather Data Analysis," IJCSNS International Journal of Computer Science and Network Security, vol. 17, no. 6, pp. 22-28, June 2017.
- [14] D.P. Sharma and Kapil Khandelwal, "Knowledge-Based Systems, Problem Solving Competence and Learnability," Springer-Verlag Berlin Heidelberg 2011, vol. 250, no. 2011, pp. 543-547, 2011.
- [15] P. P. Sondwale, "Overview of Predictive and Descriptive Data Mining Techniques," vol. 5, no. 4, , April 2015.
- [16] Durga Prasad Sharma, "Integrating Multi Criteria Decision Making Model With Geographic Information System For Land Management," International Journal of Decision Science & Information Technology, Vol. 3, No. 1, pp. 32-42, 2011.
- [17] Nishchala C. Barde, Mrunalinee Patole, "Classification and Forecasting of Weather using ANN, k-NN and Naïve Bayes Algorithms," International Journal of Science and Research (IJSR), vol. 5, no. 2, pp. 17-42, February 2016.
- [18] N. V. D. K. C. T. Abhishek Saxena, "the review of weather prediction using artificial neural networks," International Journal of Engineering Research & Technology, vol. 2, no. 11, p. 222.340, (November - 2013).
- [19] S.Saraswathi and Dr. Mary Immaculate Sheela, "Comparative Study of Different Clustering and Decision Tree for Data Mining Algorithm," International Journal of Computer Science and Mobile Computing, vol. 3, no. 11, p. 422 – 428, November 2014.
- [20] V. Goebel, Knowledge Discovery in Databases (KDD) - Data Mining (DM), Department of Informatics, University of Oslo, 2014.
- [21] Ravindra Changala, D.Rajeswara Rao, T Janardhana Rao, P Kiran Kumar, Kareemunnisa, "Knowledge Discovery Process: The Next Step for Knowledge Search," vol. 3, no. 5, May 2015.
- [22] Thu Zar Phyu, Nyein Nyein Oo, "Performance Comparison of Feature Selection Methods," EDP Sciences, vol. 42, no. 06002, p. 4, 2016.
- [23] Jiawei Han, Micheline Kamber, Jian Pei, Data Mining Concepts and Techniques, Simon Fraser University: Elsevier, 2012.
- [24] Abhishek Saxena, Neeta Verma ,Dr K. C. Tripathi, "A Review Study of Weather Forecasting Using Artificial Neural Network Approach," International Journal of Engineering Research & Technology (IJERT), vol. 2, no. 11, pp. 2029-2035, November - 2013.
- [25] Meera Narvekar, Priyanca Fargose, "Daily Weather Forecasting using Artificial Neural Network," International Journal of Computer Applications, vol. 121, no. 22, pp. 9-13, July 2015.
- [26] Yükle, "Data Mining: Practical Machine Learning Tools and Techniques, Second Edition," Artificial Intelligency , 6 June 2017. [Online]. Available: <http://genderi.org/data-mining-practical-machine-learning-tools-and-techniques-se.html?page=111>. [Accessed 30 December 2017].
- [27] Shweta Kharya, Sunita Soni, "Weighted Naive Bayes Classifier: A Predictive Model for Breast Cancer Detection," International Journal of Computer Applications , vol. 133, no. 9, pp. 32-37, January 2016.
- [28] Jiawei Han, Micheline Kamber, Jian Pei, Data Mining Concepts and Techniques, Elsevier., 2012.
- [29] Peng Cao, Xiaoli Liua, Jian Zhang, Dazhe Zhao, Min Huang, Osmar Zaiane, "norm regularized multi-kernel based joint nonlinear feature selection and over-sampling for imbalanced data classification," Neurocomputing, vol. 234, no. x, pp. 38-57, 19 April 2017.

# A Recommender System for Mobile Applications of Google Play Store

Ahlam Fuad<sup>1</sup>, Sahar Bayoumi<sup>2</sup>, Hessah Al-Yahya<sup>3</sup>

Department of Information Technology  
College of Computer and Information Sciences  
King Saud University, Riyadh, Saudi Arabia<sup>1,2,3</sup>  
Institute of Graduate Studies and Research,  
Alexandria University, Alexandria, EGYPT<sup>2</sup>

**Abstract**—With the growth in the smartphone market, many applications can be downloaded by users. Users struggle with the availability of a massive number of mobile applications in the market while finding a suitable application to meet their needs. Indeed, there is a critical demand for personalized application recommendations. To address this problem, we propose a model that seamlessly combines content-based filtering with application profiles. We analyzed the applications available on the Google Play app store to extract the essential features for choosing an app and then used these features to build app profiles. Based on the number of installations, the number of reviews, app size, and category, we developed a content-based recommender system that can suggest some apps for users based on what they have searched for in the application's profile. We tested our model using a k-nearest neighbor algorithm and demonstrated that our system achieved good and reasonable results.

**Keywords**—Application profile; content-based filtering; Google play; mobile applications; recommender systems

## I. INTRODUCTION

Recent years have witnessed massive growth in mobile devices with an increasing number of users as mobile devices have become part of every component of modern life. The smartphone market has grown dramatically, and users can now take advantage of various features in applications, which can easily be obtained from centralized markets, such as Google Play. Google Play is Google's official store and portal for Android apps that was launched in 2008 and accumulated more than 1 million downloadable and ratable applications now [1]. In December 2018, the number of available apps in the Google Play App Store was nearly 2.6 million [2].

Due to the substantial and growing number of available mobile applications in application stores, it becomes necessary to provide a system that identifies user interest based on what the system believes the user likes through his/her profile. Using a user profile would support an efficient and personalized application filtering system.

The general idea of filtering is to get a sub-collection of applications based on a specified category. There are different approaches to performing information filtering, including classification and recommendation. Classification is a step taken to reduce the sparseness of the input space by classifying applications into predefined interest categories. The applications in stores are labeled according to a high-level

and store-specific classification method. This approach is limited by the fact that it depends entirely on the textual description available from the store [3].

Moreover, a recommender system is another way to filter information and is widely used in several domains. It is a decision-making tool that helps developers predict what a user will like or dislike from a list of applications. It provides personalized information by learning the user's interests from tracing through his/her interactions. It is also an excellent option for search fields as a recommender system that lets users discover more applications [3].

In this work, we explore a method of constructing recommender systems for apps in the Google Play app store based on the app profile. Issues related to modeling app preferences and choosing a set of recommended apps were investigated. Furthermore, a k-nearest neighbor classification approach (KNN) to classify apps based on the most influential attributes of apps within categories proposed. A prototype system is then built as a proof of concept, which tracks application profiles and then presents recommended applications to the user. Therefore, the research aim to answer the following question:

What are the most significant attributes of an application profile that could be used for developing a recommender system?

The remaining sections of this paper are organized as follows. Section II presents an overview and background of the topic. Section III discusses the related work of analyzing apps in app stores and app recommender systems. In Section IV, the methodology is introduced, and the results are explained and discussed in Section V. Finally, conclusions and future work are presented in Section VI.

## II. BACKGROUND

In this section, we discuss the background information and knowledge domains required for developing a recommender system.

### A. Pearson's Correlation Coefficient

Pearson's correlation coefficient is also referred to as Pearson product-moment correlation coefficient (PPMCC) or the bivariate correlation, is a measure of the linear correlation between two variables [4]. It evaluates how well the

relationship between two variables can be described. The statistic defined in the range  $[-1, +1]$  which indicates how strongly the two variables are associated, where -1 indicates total negative linear correlation and 1 indicates total positive linear correlation. A value of 0 indicates no correlation.

### B. K-Nearest Neighbor Classification (KNN)

The k-nearest neighbors' algorithm is a type of instance-based supervised learning approach. It is one of the simplest and most commonly used classification techniques and is easy to learn and implement and robust to noise. It is used mostly for classification and sometimes for predictive regression problems, in which a number of nearest neighbors of each data point are used based on the value of k, which represents how many nearest neighbors are to be considered to determine the class of a test sample data point. In other words, the KNN algorithm finds solutions by identifying similar objects. It is also called lazy learning because the function is only approximated locally and all computation is postponed until classification. This rule preserves the complete training set throughout the learning process and assigns to each query a class represented by the most frequent label of its KNN in the training set. One of the significant drawbacks of KNN is that becomes slow as the size of the data in use increases [5], [6].

### C. Recommender Systems

Recommender systems (RSs) are techniques and software tools that provide users with suggestions for information or items that may be of interest to the user. Those suggestions will improve the user's decision-making processes, such as choosing what music to listen, what things to buy, or what apps to install. Thus RSs are the most popular and powerful tools in e-commerce [7]. Coincidence is one of the major stimuli for RSs to help the user discover things he did not look for explicitly. The essential computational task of RSs is predicting the subjective evaluation which the user gives to an item. These predictions can be computed by using predictive models with common characteristics. For example, the ratings of the user's previously purchased items can be exploited. Recommendation systems can be classified into three major categories to generate a list of recommendations based on a particular prediction technique [8], [9].

1) *Content-based recommender systems*: Content-based recommendation approaches analyze the descriptions of items rated by a user previously to build a user profile of his interests based on the items' features. Later, this profile will help to suggest additional items with similar properties. Content-based recommendation systems use methods that are focused on the items' characteristics or descriptions. These methods build a profile for each user, which is called a content-based profile that conserves the features of the previously viewed items. Then, the RS will get the most suitable details for the user by comparing the information in the generated profile and the descriptions of items [7]. For example, we have our items: A, B, C, and D. Tom likes items B, C, and D; John wants A, B, and C; and Sozy likes C. Therefore, by comparing John's and Tom's liked items, it is apparent that they both like B and C, and then the recommender system conclude that B and C are similar. If

Sozy likes C, then item B should be recommended to him.

2) *Collaborative recommender systems*: Collaborative recommender approaches collect feedback information from all the users who rate the items. These approaches build a model based on the user's past behavior and the similar decisions of the other users. Thus, this model can be used to predict items the user may be interested. For example, Sozy again likes C and D. We need a recommender system to search for a person with similar preferences to Sozy, so we can notice that Tom also likes C and D. Therefore, he is the user who is identical to Sozy. Because he also likes B, B is recommended to Sozy.

Content-based approaches mostly perform better than collaborative filtering, especially when the data is extremely sparse. Merging both methods may improve the results {Suggesting Points-of-Interest via Content-Based, Collaborative, and Hybrid Fusion Methods in Mobile Devices}.

3) *Hybrid recommender systems*: Hybrid recommender systems have been developed by combining the abilities of both collaborative and content-based recommendations. These systems were introduced due to the limitations of the two techniques described above. Hybrid recommender approaches have been implemented using several methods: by applying the content-based and collaborative-based predictions separately and then combining both of them, by adding content-based capabilities to the collaborative-based approach (or vice versa), or by unifying the approaches into one model. Hybrid methods provide more accurate recommendations than simple approaches (collaborative methods and content-based methods).

## III. RELATED WORK

This section discusses the state of the art of within two directions: data analysis techniques and recommender systems. The related studies divided into three categories: analyzing applications in different app stores, applications based on similarity measures, and application based on recommender systems.

### A. Application Analysis Studies

The author of the study [10] aimed to analyze app store data. They extracted feature information from a set of data collected from Blackberry apps using data mining in order to analyze the technical, business, and customer issues of apps. The results of this work indicate a strong correlation between the rank of app downloads and the customer rating and no relationship between price and rating, nor between price and downloads. These results partially match those observed in [11], where the study aimed to analyze the Google app store in order to identify correlations among app features, and the authors found a strong relationship between the number of downloads and price as well as between participation and price. In a recent study [12], the authors investigated the factors which impact the rating of Google play store apps. They analyzed 10,840 apps, and they indicated that app ratings help to get more downloads. Furthermore they found

that the used keyword in the app title plays an important role in determining the higher and lower ratings.

Studies [13] and [14] aimed to analyze the characteristics of apps extracted from app stores. Interestingly the experiments of [13] proved that the app size, the number of promotional images displayed on the app's web store page, and the app SDK version are the most influential factors in defining high-rated apps. Studies [1] and [15] used the Causal Impact Release Analysis tool to facilitate app store analysis.

Studies [16] and [17] introduced a novel approach for app classification utilizing features extracted from both web knowledge and relevant real-world context. Then, they integrated these extracted features into a machine learning model (Maximum Entropy (MaxEnt)) for training an app classifier.

### B. App Similarity Studies

The study [18] introduced a classification system in order to classify mobile apps. They mined 5,993 apps from both the Apple and Google app stores and then classified them based on support vector machines (SVMs). As a result of this study, the automated app classification system achieved an excellent accuracy. Another study [19] proposed a novel technique for measuring the similarity among apps based on agglomerative hierarchical clustering techniques. They mined data for 17,877 apps from the Google and BlackBerry app stores. The empirical results of this study indicate an improvement over the existing categorization quality of both stores. In another study [11], the authors aimed to build clusters of similar apps using a probabilistic topic modeling technique and a k-means clustering method. The results showed that the Google Play categorization system does not respect application similarity.

The study [20] addressed the application classification issue and introduced a novel method for classifying apps using two methods. The first method used a neural language model applied to smartphone logs to embed apps into a low-dimensional space, while the second one used the k-nearest neighbors' classification method in the embedding space; the experimental results showed that the second proposed approach outperformed the current state of the art.

In a recent study [21], the authors introduced a classification method for local mobile app using deep neural network. They evaluated a dataset of Google Play to demonstrate the effectiveness of their method. Their results outperformed the baseline method by 5.5% related to F1 score. This study focused just on classifying local apps such as "Travel & Local" in the store.

A new framework for app categorization (FRAC+) has been proposed in [22], which is based on a data-driven topic model to suggest the appropriate categories for an app store, as well as to detect miscategorized apps. Experiments with the proposed system have shown that it is aligned with the new categories of the Google Play store.

### C. Application Recommender System Studies

There is considerable literature available on both recommendation systems and mobile recommendation systems with various descriptions of recommendation systems

in general [1], [23], [24]. The authors in [25] discussed the incorporation of recommender systems in the mobile application domain. They used a hybrid recommender system to deal with the added complexity of context and recommend appropriate mobile applications to users. Thus, this approach provides positive ratings. Therefore, based on this study, users can select from among several content-based or collaborative filtering components.

A new efficient framework called "SimApp" was proposed to detect similar applications using an online kernel learning algorithm [26]. They crawled real data from the Google app store and extracted a multi-modal heterogeneous data set. Their outcomes indicate the efficiency of the proposed framework. The similarity of items may help the application of content-based recommender systems. Another study introduced a framework based on the incorporation of version description features into app recommendation [27]. Another study [28] described the implementation of a hybrid recommender system that employed five different filtering techniques to help users when choosing a new application to download from a market. This system was also able to solve many common problems found in collaborative recommender systems that reduce the quality of the generated predictions. The study was based on using information collected from different users to support users with recommendations based on their history. The results showed good performance in terms of mean absolute error (MAE) and users' satisfaction.

The study [8] discussed assisting the users in choosing the appropriate application using recommendations. The author proposed a recommender system for mobile applications by integrating two methods: tracking user behavior to get his preferences to find new and similar apps to their used ones and utilizing the user's context in order to provide him with useful recommendations by using the Google Play Engine. While in the study [29], the authors proposed a recommender method for apps based on graph techniques. Interestingly, the proposed method can recommend apps without the need for specifying user preferences. Another paper proposed a recommender system for the mobile application market by understanding the mobile user's preferences and usage patterns for the types of applications they select and the online downloading process. The authors collected data from Google Play and then used statistical analysis and a pilot survey to find app features that influence user choices [30]. In [31], the authors proposed a novel structural user choice model (SUCM) to learn fine-grained user preferences by exploiting the hierarchical taxonomy of apps (tree hierarchy of apps). Also, they designed an efficient learning algorithm to estimate the model parameters. They used a diverse dataset of 52,483 users, 26,426 apps, and 3,286,156 review observations. The outcomes of this study show that SUCM consistently outperforms state-of-the-art Top-N recommendation methods by a significant margin. The study in [32] proposed a novel method using a unified model that combines content-based filtering with collaborative filtering, harnessing information from both ratings and reviews. This study applied topic modeling techniques to the review text and aligned the topics with rating dimensions to improve prediction accuracy. Another study [33] proposed a unified model VAMF for the

version-aware mobile app recommendation problem to address the data sparsity issue by incorporating review text from both the version level and the app level and modeling version based correlations of version-level temporal correlations and app-level aggregate correlation. They also proposed an efficient algorithm to solve the model and analyze its optimality and complexity. They used a Google Play dataset that contained the reviews for all of its versions and the descriptions of its latest version. The experiments conducted in this study on a large dataset showed that the proposed method outperforms comparable methods in prediction accuracy and that the proposed algorithm can be linearly scaled.

The study in [34] introduced a sequential approach for modeling the popularity of mobile apps by collecting data from 15,045 apps. They produced a popularity-based hidden Markov model (PHMM) for a variety of tasks, including app recommendation and review spam detection, and demonstrated its usefulness in ranking fraud detection. The experimental results validated both the effectiveness and efficiency of the proposed popularity modeling approach. Another study built on a hidden Markov model where the authors proposed a mechanism for modeling three main factors governing the app installation behavior of smartphone users: short-term context, co-installation pattern, and random choice. Then, a heterogeneous hidden Markov model (heterogeneous HMM) was used to incorporate these main factors. They used a combination of app installation data from the installation records of 9009 users with a portion of the Netflix data set from 54,314 users on 3561 movies. The experimental results indicated that the proposed system can outperform other methods consistently under different experimental settings [35].

The study [36] was generally focused on recommending independent items to users who were suggested by a hybrid cross-platform app recommendation (STAR) system. Another study [37] introduced recommender systems on mobile platforms based on user profiles generated from the installed apps. They improved on existing machine learning models to predict user profiles. The results of this study showed an increase in these models' predictive accuracy. Furthermore, study [38] introduced a recommender system (Vanilla) that considers social and contextual information processes. The system allows the comparison of different recommendation techniques. Besides this, Vanilla includes eleven contextual dimensions and a mechanism for analyzing the influence of social networks on app consumption. They found that the new proposed approach has a strong correlation with previous approaches and better efficiency than other techniques. A recent study proposed a context-aware approach for mobile app recommendation using tensor analysis (CAMAR) [39]. They conducted data analysis on Google Play Store and Apple's App Store in order to find the mobile apps characteristics. They utilized an effective tensor-based framework to integrate the features on users and apps and app category information to facilitate the app recommendation performance. Thus, they demonstrated the effectiveness of their proposed method.

A considerable amount of literature has been published recently regarding recommender systems for mobile apps. One study proposed a recommender system for mobile apps based on user reviews using topic modeling techniques and probability distributions to represent apps features [40]. Hence, this study aimed to construct a user profile based on his installed apps in order to identify his preferences. Therefore, they found that user reviews, extracted from datasets that were crawled from the Apple App Store, represented apps features efficiently. Another study [41] introduced a mobile sparse additive generative model (Mobi-SAGE) to recommend apps. They crawled an extensive collection of apps from the 360 App Store in China. The results of their study demonstrated that the proposed model outperformed other existing state-of-the-art methods.

According to the literature review, many models were developed using variety of features to support users selecting applications. Table I show a summary for researches discussed in the last section regarding the platform and used features for deployed recommender systems.

TABLE I. SUMMARY OF RECOMMENDER SYSTEM RESEARCHES

| Ref. | Platform   | Attributes   |
|------|--|--|
| [25] | play.tools framework                             | Ratings.   |
| [26] | Android  | Name, category, description, developer, update, permissions, and app logo/images.                    |
| [27] | Apple  | Version-categories, genre, and ratings.  |
| [28] | Android, Apple                                   | User history, Tags used to define the applications by the user, and user satisfaction.               |
| [29] | Apple, Android, Blackberry and Windows App store | Apps installed on the user's phone.  |
| [30] | Android  | Cost, app logo/ image, gender, and types of downloaded applications.                                 |
| [31] | Android  | Category tree.   |
| [32] | Amazon Dataset                                   | Ratings and reviews.   |
| [33] | Android  | # of users, # of versions, and # of ratings.   |
| [34] | Apple  | Trend based Applications, rating, review spam detection, and ranking fraud detection.                |
| [35] | Android  | User installation behavior, user preferences, and Modeling random choice.                            |
| [36] | Apple  | Application Rating between different platforms (iPhone-iPad platform and iPhone-iPad-iMac platform). |
| [37] | Android  | Cost, ratings, and user profile based on the installed applications.                                 |
| [38] | Android  | Categories and ratings.  |
| [39] | Android, Apple                                   | User's preference, app category, and features of multiple views.                                     |
| [40] | Apple  | User preferences and reviews.  |
| [41] | Android  | User interests, ratings, and privacy preferences.  |



As shown, most of previous researches used either collaborative or hybrid approaches for building recommender systems which increase system complexity.

Our research aims to study the inter-relation between app attributes to select the most significant features. Furthermore, develop a content-based recommender system using the selected attributes.

#### IV. METHODOLOGY

In this section, we illustrate our proposed system, which includes five steps, as shown in Fig. 1. In the first step, we acquired the dataset. In the second step, we prepared the dataset for analysis to complete the other steps. In the third step, we analyzed the data to identify the correlations among the features. In the fourth step, we designed a suitable recommender system model. The final step is to test the recommender system model and make some recommendations.

The rest of this section explains the steps of our approach in more detail.

##### A. Data Acquisition

The dataset used to achieve this study is consists of 10841 apps scraped from the Google Play store, which publicly available on the Kaggle website [42], where its most recent update provided two months ago. This dataset provides detailed information from Google about the apps on the Google Play store. Thus, it includes 13 attributes with three datatypes: String, Categorical and Numeric. Only the “reviews” belongs to the numeric data type with values range from 0 to approximate 78Million reviews. Table II shows the attributes based on data types.

The dataset includes 33 different categories (shown in Fig. 2) and 118 genres, which define the sub-category for each application.

An initial statistical summary about the numerical data shown in Table III to provide a deep understanding about each attribute for further processing.

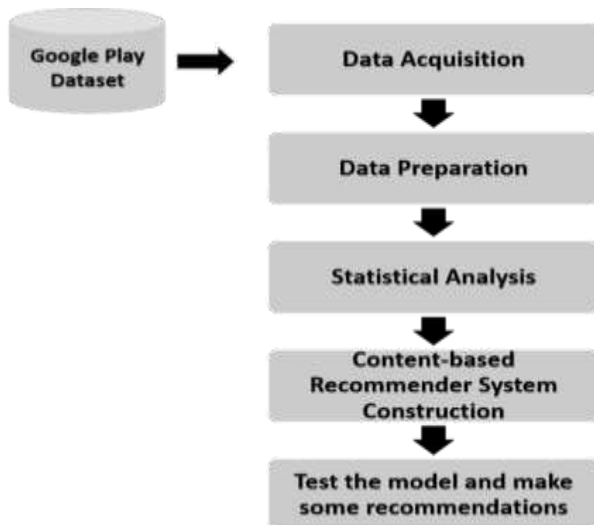


Fig. 1. Proposed System Methodology.

TABLE II. THE GOOGLE PLAY STORE ATTRIBUTES CLASSIFIED BY DATA TYPE

|             | Categorical attributes |              |
|-------------|------------------------|--------------|
| Installs    | Type                   | Genres       |
| Category    | Content Rating         | Android Ver. |
|             | Rating                 |              |
|             | String attributes      |              |
| App name    | Size                   | Price        |
| Last Update | App Ver.               |              |
|             | Numeric                |              |
|             | Reviews                |              |



Fig. 2. A Cloud Showing the Categories.

TABLE III. STATISTICAL SUMMARY FOR THE NUMERICAL ATTRIBUTES

|          | Count | Mean     | STD      | Min      | Max      |
|----------|-------|----------|----------|----------|----------|
| Rating   | 8196  | 4.17E+00 | 5.37E-01 | 1.00E+00 | 5.00E+00 |
| Reviews  | 9660  | 2.17E+05 | 1.83E+06 | 0.00E+00 | 7.82E+07 |
| Size     | 9660  | 3.18E+07 | 3.48E+07 | 8.70E+03 | 1.05E+08 |
| Installs | 9660  | 7.78E+06 | 5.38E+07 | 0.00E+00 | 1.00E+09 |
| Type     | 9660  | 7.82E-02 | 2.68E-01 | 0.00E+00 | 1.00E+00 |
| Price    | 9660  | 1.10E+00 | 1.69E+01 | 0.00E+00 | 4.00E+02 |

##### B. Data Preparation

Data preparation is a necessary step to analyze the data correctly, and to facilitate understanding of the relationships among the data and to gain useful insights. As shown from Table III; the Kaggle dataset contains some missing values for "rating" attribute. In addition, inconsistencies in the attributes "size", "installs", "price" and "reviews" by remove unwanted information.

The data preparation process of the dataset using Python are applied in three consecutive steps as follows:

- 1) Removed duplicated rows which reduce the dataset by 1181 records.
- 2) Convert string and categorical datatypes into numeric to allow further data analysis.
- 3) Insure consistency of numeric attributes through the following:

a) Convert the characters 'K' and 'M' within the “Size” attribute to a numeric values.

b) Propagate last valid observation forward using forward fill method to replace the "Varies with device" value to get a numeric value within the "Size" attribute.

c) Convert the characters '+', '\$', and 'M' from "Installs," "Price," and "Reviews," into numeric values.

d) The last attribute that needed to be converted was "Type," where we mapped the string values to numeric ones.

Regarding the missing values for the "rate" attributes, a null value kept for the associate records as they represent 15.15% of the total dataset. There are different reasons for the missing values within the "rate" attributes such as: new released application or not common for users. Therefore, we decided to include the records with the null value.

### C. Statistical Analysis

Pearson's correlation coefficient is used to measure the linear correlation between the numerical features of the apps in the Google Play store. Pearson's correlation coefficient is a statistical measure used to determine the strength of the relationship between paired data [4]. Fig. 3 shows the Pearson's correlation coefficient heatmap between the numeric features for the dataset. The correlation coefficients between attributes is the ground truth that help in choosing the most prominent features for further use in building the recommender system.

According to the heatmap; the attributes "reviews", "size", and "installs" are the most correlated attributes while other attributes are not.

### D. Recommender System Construction

When looking for app a common attribute to be specified is the category. Then, a list of all apps under the specified category are shown. Sorting the apps to guide you to the best is restricted by choose one attribute. Based on the correlation coefficient and the importance of category attribute for the user; we decided to include all the four attributes ("reviews", "size", "installs", and "category" ) in defining the app profile for a content-based recommender system. The app profile consists of 37 columns; the first column is the app id within the dataset, and the next 33 columns represent the 33 categories of apps and three columns for "reviews", "installations", and "size". Furthermore, each column/feature scaled by its maximum absolute value for efficient calculations.



Fig. 3. Pearson's Correlation Coefficient Heat Map.

### E. Classification using K-NN

A K-nearest neighbors (K-NN) algorithm; as an unsupervised machine learning; is used to measure the similarity between apps using their profiles. The nearest neighbor algorithm uses a "brute" algorithm and "cosine" metric.

## V. RESULTS AND DISCUSSION

Our proposed recommender system developed based on building a profile for each app using the most significant attributes. However, Pearson's correlation coefficient (as in Fig. 3) showed that "reviews", "size" and "installs" are the most significant correlated positively attributes. The highest correlated pair is "Installs" and "Reviews" with value 0.63. Thus, obviously highlight that users prefer to download apps that intensively reviewed. Fig. 4 shows a log scale for the relationship between the "installs" and the "reviews".

The second significant correlation between "Size" and "Installs" with value 0.19 shows a considered level of importance of the application size for users. The log scale relationship (Fig. 5) shows increase number of installs for small size applications while still large applications attract users.

However, most popular mobile apps, especially game apps, tend to be feature-rich, which implies that additional code and assets can pump up file sizes. Generally, the statistic indicates an increase in the number of mobile game app downloads from Google Play worldwide. In 2018, a total of 29.4 billion mobile games were downloaded globally across Google's app store [43]. For example, the famous PUBG Game, which is sized at 1.6 GB for Android, is considered to be the most downloaded mobile game in the last quarter, which is a free, high-resolution game with excellent graphics and details [44].

Furthermore, a less significant correlation between "Size" and "Reviews" with a value 0.16 showed the importance of the "Size" attribute along with the "reviews" which highlight the user's need to optimize their storage use. Fig. 6 shows that apps with small sizes are with more reviews, which therefore more installs.

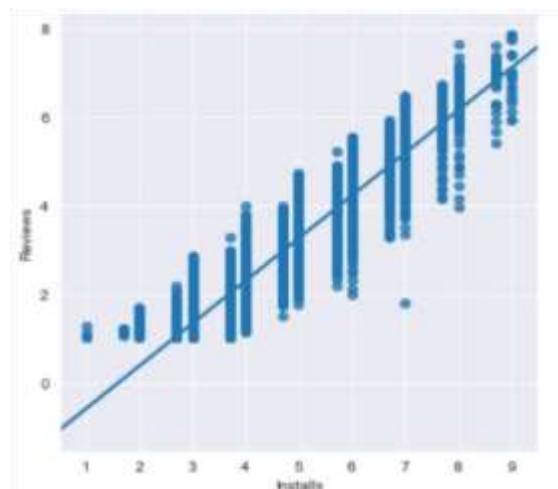


Fig. 4. Correlation between Installs and Reviews.

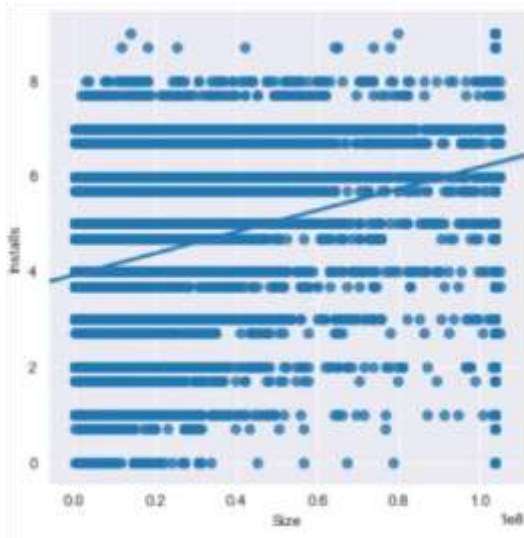


Fig. 5. Correlation between Size and Installs.

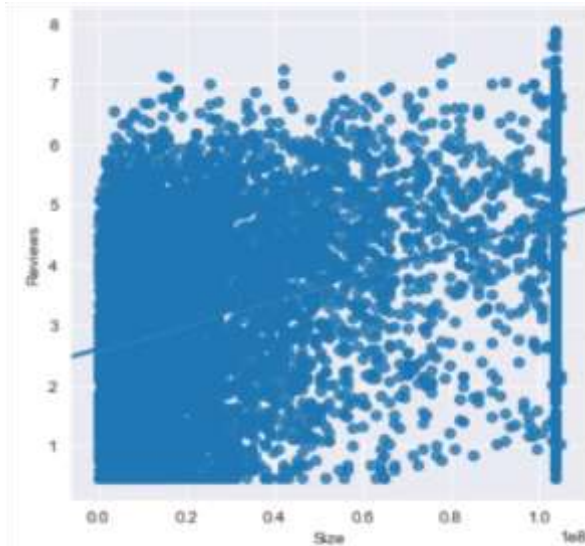


Fig. 6. Correlation between Size and Reviews.

To evaluate the developed recommender system; three case studies are used for different categories.

Case study 1: From the Social category; the “Facebook” application with id “8823”. By applying the “K-NN” to the matrix of profiles of all applications, four apps are recommended and sorted based on their K-NN metrics. Table IV shows the recommended apps and their distances from the Facebook app.

Case study 2: From the Game category, the “Candy Crush Saga” application with id “7484” is used for testing. By applying the “K-NN”, four recommended games are shown on Table V along with their K-NN distances.

Case study 3: From the education category, the “Wikipedia” application with id “8452” is used for testing. Table VI shows the recommended apps and their corresponding distance from the “Wikipedia” app.

TABLE IV. RECOMMENDED APPS FOR FACEBOOK CASE STUDY (ID=8823)

| ID   | App          | Reviews  | Size     | Installs | distance |
|------|--------------|----------|----------|----------|----------|
| 8824 | Instagram    | 6.66E+07 | 1.04E+08 | 1.00E+09 | 4.86E-04 |
| 8827 | SnapChat     | 1.70E+07 | 1.04E+08 | 5.00E+08 | 3.68E-03 |
| 3325 | Facebook Lit | 8.61E+06 | 1.04E+08 | 5.00E+08 | 4.81E-03 |
| 8830 | Google+      | 4.83E+06 | 1.04E+08 | 1.00E+09 | 5.86E-02 |

TABLE V. RECOMMENDED APPS FOR CANDY CRUSH SAGA CASE STUDY (ID=7484)

| ID   | App                     | Reviews  | Size     | Installs | Distance |
|------|-------------------------|----------|----------|----------|----------|
| 7485 | Dream League Soccer2018 | 9.88E+06 | 7.76E+07 | 1.00E+08 | 1.21E-02 |
| 6947 | Temple Run 2            | 8.12E+06 | 6.50E+07 | 5.00E+08 | 2.55E-02 |
| 8762 | My Talking Tom          | 1.49E+07 | 1.04E+08 | 5.00E+08 | 2.76E-02 |
| 7508 | Subway Surfers          | 2.77E+07 | 7.97E+07 | 1.00E+09 | 5.93E-02 |

TABLE VI. RECOMMENDED APPS FOR WIKIPEDIA CASE STUDY (ID=8452)

| ID   | App                        | Reviews  | Size     | Installs | distance |
|------|----------------------------|----------|----------|----------|----------|
| 9587 | English Hindi Dictionary   | 3.84E+05 | 1.04E+08 | 1.00E+07 | 1.65E-02 |
| 8455 | Dictionary-Merriam-Webster | 4.54E+05 | 1.04E+08 | 1.00E+07 | 1.26E-01 |
| 9496 | Dictionary                 | 2.64E+05 | 1.04E+08 | 1.00E+07 | 1.76E-01 |
| 8463 | Moon+ Reader               | 2.34E+05 | 1.04E+08 | 1.00E+07 | 1.87E-01 |

## VI. CONCLUSIONS

Mobile app recommendation based on only application installation records is a challenging task. In this paper, we proposed a model that seamlessly combines content-based filtering with application profiles. Thus, we used a real-world app dataset from Google Play to analyze app information and then utilized the most effective content to build a content-based recommender system. Based on our results, the most influential factors in choosing an app are the number of installs, number of reviews, app size, and category. Finally, we introduced some examples to prove that our system achieved good and reasonable results.

## VII. FUTURE WORK

Although our proposed recommender system was originally designed for app recommendation from the Google Play store, we believe it can also be applied to other stores as well as other domains, such as book recommendation, music recommendation, movie recommendation, and food recommendation. Therefore, we believe that some possible future studies using the same experimental set up are possible.

Also, we aim to build a benchmark dataset from various application stores which could be used by researchers for building AI systems and recommender systems.

#### ACKNOWLEDGMENT

This research project was supported by a grant from the "Research Center of the Female Scientific and Medical Colleges", Deanship of Scientific Research, King Saud University. The authors thank the Deanship of Scientific Research and RSSU at King Saud University for their technical support.

#### REFERENCES

- [1] W. Martin, F. Sarro, Y. Jia, Y. Zhang, and M. Harman, "A survey of app store analysis for software engineering," *IEEE Trans. Softw. Eng.*, vol. 43, no. 9, pp. 817–847, 2017, doi: 10.1109/TSE.2016.2630689.
- [2] "Google Play Store: number of apps 2018 | Statista." [Online]. Available: <https://www.statista.com/statistics/266210/number-of-available-applications-in-the-google-play-store/>. [Accessed: 30-Mar-2019].
- [3] G. Adomavicius and A. Tuzhilin, "Toward the next generation of recommender systems: a survey of the state-of-the-art and possible extensions," *IEEE Trans. Knowl. Data Eng.*, vol. 17, no. 6, pp. 734–749, Jun. 2005, doi: 10.1109/TKDE.2005.99.
- [4] Y. Mu, X. Liu, and L. Wang, "A Pearson's correlation coefficient based decision tree and its parallel implementation," *Inf. Sci. (Ny)*, vol. 435, pp. 40–58, Apr. 2018, doi: 10.1016/j.ins.2017.12.059.
- [5] M.-A. Amal and B.-A. Ahmed, "Survey of Nearest Neighbor Condensing Techniques," *Int. J. Adv. Comput. Sci. Appl.*, vol. 2, no. 11, 2011, doi: 10.14569/ijacsa.2011.021110.
- [6] S. Bafandeh, I. And, and M. Bolandraftar, "Application of K-Nearest Neighbor (KNN) Approach for Predicting Economic Events: Theoretical Background."
- [7] F. Ricci, L. Rokach, B. Shapira, and P. B. Kantor, *Recommender Systems Handbook*. 2011, doi: 10.1007/978-0-387-85820-3.
- [8] V. Viljanac, "RECOMMENDER SYSTEM FOR MOBILE APPLICATIONS," *Multimed. Tools Appl.*, vol. 77, no. 4, pp. 4133–4153, Feb. 2018, doi: 10.1007/s11042-017-4527-y.
- [9] R. G. De Souza, R. Chiky, and Z. K. Aoul, "Open source recommendation systems for mobile application," *CEUR Workshop Proc.*, vol. 676, pp. 55–58, 2010.
- [10] A. Finkelstein, M. Harman, Y. Jia, F. Sarro, and Y. Zhang, "Mining App Stores: Extracting Technical, Business and Customer Rating Information for Analysis and Prediction," *UCL Res. Notes*, vol. 13, p. 21, 2013.
- [11] S. Mokarizadeh and M. Matskin, "Mining and Analysis of Apps in Google Play," no. January 2013, pp. 527–535, 2013, doi: 10.5220/0004502005270535.
- [12] A. Mahmood, "Identifying the influence of various factor of apps on google play apps ratings," *J. Data, Inf. Manag.*, vol. 2, no. 1, pp. 15–23, 2020, doi: 10.1007/s42488-019-00015-w.
- [13] Y. Tian, M. Nagappan, D. Lo, and A. E. Hassan, "What are the characteristics of high-rated apps? A case study on free Android Applications BT - IEEE International Conference on Software Maintenance and Evolution," pp. 301–310, 2015.
- [14] M. Ali, M. E. Joorabchi, and A. Mesbah, "Same App, Different App Stores: A Comparative Study," *Proc. - 2017 IEEE/ACM 4th Int. Conf. Mob. Softw. Eng. Syst. MOBILESoft 2017*, pp. 79–90, 2017, doi: 10.1109/MOBILESoft.2017.3.
- [15] W. Martin, F. Sarro, and M. Harman, "Causal impact analysis for app releases in google play," pp. 435–446, 2016, doi: 10.1145/2950290.2950320.
- [16] H. Zhu, E. Chen, H. Xiong, H. Cao, and J. Tian, "Mobile app classification with enriched contextual information," *IEEE Trans. Mob. Comput.*, vol. 13, no. 7, pp. 1550–1563, 2014, doi: 10.1109/TMC.2013.113.
- [17] H. Zhu, H. Cao, E. Chen, H. Xiong, and J. Tian, "Exploiting Enriched Contextual Information for Mobile App Classification," pp. 1617–1621, 9781450311564.
- [18] G. Berardi, A. Esuli, T. Fagni, and F. Sebastiani, "Multi-store metadata-based supervised mobile app classification," *Proc. 30th Annu. ACM Symp. Appl. Comput. - SAC '15*, pp. 585–588, 2015, doi: 10.1145/2695664.2695997.
- [19] A. A. Al-Subaihin et al., "Clustering Mobile Apps Based on Mined Textual Features," *Proc. 10th ACM/IEEE Int. Symp. Empir. Softw. Eng. Meas. - ESEM '16*, pp. 1–10, 2016, doi: 10.1145/2961111.2962600.
- [20] V. Radosavljevic et al., "Smartphone App Categorization for Interest Targeting in Advertising Marketplace," 2017, pp. 93–94, doi: 10.1145/2872518.2889411.
- [21] K. Ochiai, F. Putri, and Y. Fukazawa, "Local app classification using deep neural network based on mobile app market data," *2019 IEEE Int. Conf. Pervasive Comput. Commun. PerCom 2019*, pp. 186–191, 2019, doi: 10.1109/PERCOM.2019.8767416.
- [22] D. Surian, S. Seneviratne, A. Seneviratne, and S. Chawla, "App Miscategorization Detection: A Case Study on Google Play," *IEEE Trans. Knowl. Data Eng.*, vol. 29, no. 8, pp. 1591–1604, 2017, doi: 10.1109/TKDE.2017.2686851.
- [23] A. Arampatzis and G. Kalamatianos, "Suggesting Points-of-Interest via Content-Based, Collaborative, and Hybrid Fusion Methods in Mobile Devices," *ACM Trans. Inf. Syst.*, vol. 36, no. 3, pp. 1–28, 2017, doi: 10.1145/3125620.
- [24] H. Cao and M. Lin, "Mining smartphone data for app usage prediction and recommendations: A survey," *Pervasive Mob. Comput.*, vol. 37, pp. 1–22, 2017, doi: 10.1016/j.pmcj.2017.01.007.
- [25] W. Woerndl, C. Schueller, and R. Wojtech, "A Hybrid Recommender System for Context-aware Recommendations of Mobile Applications," in *2007 IEEE 23rd International Conference on Data Engineering Workshop, 2007*, pp. 871–878, doi: 10.1109/ICDEW.2007.4401078.
- [26] N. Chen, S. C. H. Hoi, S. Li, and X. Xiao, "SimApp: A Framework for Detecting Similar Mobile Applications by Online Kernel Learning," *Proc. Eighth ACM Int. Conf. Web Search Data Min.*, pp. 305–314, 2015, doi: 10.1145/2684822.2685305.
- [27] J. Lin, K. Sugiyama, M.-Y. Kan, and T.-S. Chua, "New and Improved: Modeling Versions to Improve App Recommendation," *Proc. 37th Int. ACM SIGIR Conf. Res. Dev. Inf. Retr.*, pp. 647–656, 2014, doi: 10.1145/2600428.2609560.
- [28] E. Costa-Montenegro, A. B. Barragáns-Martínez, and M. Rey-López, "Which App? A recommender system of applications in markets: Implementation of the service for monitoring users' interaction," *Expert Syst. Appl.*, vol. 39, no. 10, pp. 9367–9375, Aug. 2012, doi: 10.1016/j.eswa.2012.02.131.
- [29] U. Bhandari, K. Sugiyama, A. Datta, and R. Jindal, "Serendipitous recommendation for mobile apps using item-item similarity graph," *Lect. Notes Comput. Sci. (including Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinformatics)*, vol. 8281 LNCS, pp. 440–451, 2013, doi: 10.1007/978-3-642-45068-6\_38.
- [30] A. Yasin, L. Liu, R. Fatima, and W. Jianmin, "Designing the Next Mobile App Recommender System for the Globe," *2017 14th International Symposium on Pervasive Systems, Algorithms and Networks & 2017 11th International Conference on Frontier of Computer Science and Technology & 2017 Third International Symposium of Creative Computing (ISPAN-FCST-ISCC)*, 2017, pp. 491–500, doi: 10.1109/ISPAN-FCST-ISCC.2017.44.
- [31] B. Liu, Y. Wu, N. Z. Gong, J. Wu, H. Xiong, and M. Ester, "Structural Analysis of User Choices for Mobile App Recommendation," *ACM Trans. Knowl. Discov. Data*, vol. 11, no. 2, pp. 1–23, Nov. 2016, doi: 10.1145/2983533.
- [32] G. Ling, M. R. Lyu, and I. King, "Ratings meet reviews, a combined approach to recommend," in *Proceedings of the 8th ACM Conference on Recommender systems - RecSys '14*, 2014, pp. 105–112, doi: 10.1145/2645710.2645728.
- [33] Y. Yao, W. X. Zhao, Y. Wang, H. Tong, F. Xu, and J. Lu, "Version-Aware Rating Prediction for Mobile App Recommendation," *ACM Trans. Inf. Syst.*, vol. 35, no. 4, pp. 1–33, Jun. 2017, doi: 10.1145/3015458.
- [34] H. Zhu, C. Liu, Y. Ge, H. Xiong, and E. Chen, "Popularity Modeling for Mobile Apps.," vol. 45, no. 7, pp. 1303–1314, 2015.

- [35] V. C. Cheng, L. Chen, W. K. Cheung, and C. kuen Fok, "A heterogeneous hidden Markov model for mobile app recommendation," *Knowl. Inf. Syst.*, vol. 57, no. 1, pp. 207–228, 2018, doi: 10.1007/s10115-017-1124-3.
- [36] T.-S. Chua et al., "Cross-Platform App Recommendation by Jointly Modeling Ratings and Texts," *ACM Trans. Inf. Syst.*, vol. 35, no. 4, pp. 1–27, 2017, doi: 10.1145/3017429.
- [37] R. M. Frey, R. Xu, C. Ammendola, O. Moling, G. Giglio, and A. Ilic, "Mobile recommendations based on interest prediction from consumer's installed apps—insights from a large-scale field study," *Inf. Syst.*, vol. 71, pp. 152–163, 2017, doi: 10.1016/j.is.2017.08.006.
- [38] D. F. Chamorro-Vela et al., "Recommendation of Mobile Applications based on social and contextual user information," *Procedia Comput. Sci.*, vol. 110, pp. 236–241, 2017, doi: 10.1016/j.procs.2017.06.090.
- [39] T. Liang et al., "CAMAR: a broad learning based context-aware recommender for mobile applications," *Knowl. Inf. Syst.*, vol. 62, no. 8, pp. 3291–3319, 2020, doi: 10.1007/s10115-020-01440-9.
- [40] K. P. Lin, Y. W. Chang, C. Y. Shen, and M. C. Lin, "Leveraging Online Word of Mouth for Personalized App Recommendation," *IEEE Trans. Comput. Soc. Syst.*, vol. 5, no. 4, pp. 1061–1070, 2018, doi: 10.1109/TCSS.2018.2878866.
- [41] H. Yin, W. Wang, L. Chen, X. Du, Q. V. Hung Nguyen, and Z. Huang, "Mobi-SAGE-RS: A sparse additive generative model-based mobile application recommender system," *Knowledge-Based Syst.*, vol. 157, pp. 68–80, 2018, doi: 10.1016/j.knosys.2018.05.028.
- [42] "Google Play Store Apps | Kaggle." [Online]. Available: <https://www.kaggle.com/lava18/google-play-store-apps/>. [Accessed: 16-Mar-2019].
- [43] "Global app stores mobile games downloads 2018 | Statistic." [Online]. Available: <https://www.statista.com/statistics/661553/global-app-stores-mobile-game-downloads/>. [Accessed: 30-Mar-2019].
- [44] "PUBG most downloaded mobile game last quarter, but revenue flags | GamesIndustry.biz." [Online]. Available: <https://www.gamesindustry.biz/articles/2018-05-04-pubg-most-downloaded-mobile-game-last-quarter-but-fails-to-make-an-impact-with-revenue>. [Accessed: 08-Sep-2020].

# Factored Phrase-Based Statistical Machine Pre-training with Extended Transformers

Vivien L. Beyala<sup>1</sup>, Perrin Li Litet<sup>3</sup>  
Computer Science  
URAIA-University of Dschang  
Dschang, Cameroon

Marcellin J. Nkenlifack<sup>2</sup>  
Computer Science  
URIFIA-University of Dschang  
Dschang, Cameroon

**Abstract**—This paper presents the development of a cascaded hybrid multi-lingual automatic translation system, by allowing a tight coupling between the two underlying research approach in machine translation, namely, the neuronal (deterministic approach) and statistical (probabilistic approach), while fully taking advantage of each method in order to improve translation performance. This architecture addresses two major problems frequently occurring when dealing with morphologically richer languages in MT, that is, the significant number unknown tokens generated due to the presence of out of vocabulary (OOV) words, and size of the output vocabulary. Additionally, we incorporated factors (additional word-level linguistic information) in order to alleviate data sparseness problem or potentially reduce language ambiguity, the factors we considered are lemmatization and Part-of-Speech tags (taking into consideration its various compounds). We combined a fully-factored transformer and a factored PB-SMT, where, the training data is pre-translated using the trained fully-factored transformer, and afterwards employed to build an PB-SMT system, parallelly using the pre-translated development set to tune parameters. Finally, in order to produce the desired results, we operated the FPB-SMT system to re-decode the pre-translated test set in a post-processing step. Experiments performed on translations from Japanese to English and English to Japanese reveals that our proposed cascaded hybrid framework outperforms the strong HMT state-of-the-art by over 8.61% BLEU and 7.25% BLEU, respectively, for validation set, and over 8.70% BLEU and 7.70% BLEU, respectively, for test set.

**Keywords**—Machine translation; transformer; statistical machine; morphologically rich; hybrid

## I. INTRODUCTION

Machine translation has known an improvement in the state-of-the-art performance by the intervention of Transformers [1] which is a new paradigm in Neural Machine Translation (NMT) [2] [3] powered by frameworks of sequence to sequence learning, thus rivaling since then the factored statistical machine translation paradigm [4] which has achieved the state-of-the-art in SMT frameworks [5] [6]. However, the fundamental design of NMT models which imposes them to make reliable the input representation of a word by observing several instances of that word in multiple examples, and make them to eventually face coverage issues during the computational complexity control by limiting the input and output vocabulary sizes, greatly affects their translation performance when processing rare or OOV (out of vocabulary) words (which are those neither included in the

vocabulary nor seen in the training data set, therefore mapped to an UNK token since being considered as unknown words) for languages that are morphologically rich and of low resources (such as Cameroon local languages and some national well known languages namely Arabic, Czech, German, Italian and Turkish). Though having fluent translations in most cases, NMT face challenges in modeling languages syntactic and semantic deeper aspects.

As such, for low-resource (or small corpus) and morphologically rich language conditions, the necessity to incorporate for the surface level words various linguistic annotations was found to resolve semantic ambiguities and data sparseness, thus leading to better translation of rare words or OOVs and greater generalization capacity as illustrated [4] when addressing this issue for the traditional SMT architecture [7] by proposing the factored translation model. This linguistic annotations or factors include features such as lemmas, stems, morphological classes, roots, data-driven clusters, data-driven clusters, part-of-speeches, constituency parsing and compounds. With the vision of alleviating data sparseness and reducing language ambiguity, such extra features may be of enormous benefits when added to both NMT and Phrase-based SMT frameworks.

However, the aim of improving translation performance has inspired much research works through the combination of NMT and SMT paradigms [8] [9] [10] [11] in order to fully take advantage of each system's strength, and therefore overcoming the deficiencies of meaningless translations (those with meanings totally different compared to source sentences) and limited vocabulary size usually faced by pure NMT models, although its strong language modeling capacities. By contrast, the hard word alignment technic of PBSMT models reflects the source sentences adequacy extremely well, thus helping to some extent to restore the meaning of source sentence whenever wrong translations are produced. The framework proposed by [12] is very close to our work in the global context and overall all architecture but as compared to theirs, ours integrates outperforming paradigms in both the NMT and PBSMT frameworks, that is, Factored Transformers and Factored SMT, respectively. Also, we used linguistic features taking into consideration compounds bot at the NMT (augmenting its embedding layer so as to learn various compositional input representations at different granularity levels) and SMT levels and finally, we proposed a novel UNK replacement algorithm. Our experimental findings reveals that our hybrid model provide consistently and significantly better

translation quality for morphologically rich and low resourced languages when coming across rare and unknown words than the state-of-the-art of hybrid translation models.

This paper is organized as follows: A literature review is performed in Section 2. We discuss the factorization process with the integration of compounds in Section 3. In Section 4, we describe the transformer operation with the incorporation of linguistic factors in detail. Section 5 detail our proposed neural hybrid MT framework. In Section 6, the results of two sets of experiments on Japanese to English and English to Japanese tasks are reported measured by their BLEU score. Finally, in Section 7 we summarize our findings and outline future plans.

## II. ANALOGOUS RESEARCHES

By using a combination of different modules, paradigms, resources and approaches, many researchers have explored Hybrid MT systems. In order to produce publishable quality translations, corrections of repetitive errors have to be implemented through the development of various automatic or semi-automatic post-processing techniques, human post-edition usually still have to be operated on the overall resulting MT output [13] [14]. Although human post-editing (PE) is needed over MT outputs, MT output post-edition more often remains cheaper and faster as compared to performing human evaluation from scratch. The authors in [15] [16], and [17] revealed that in some cases productivity can be increased as well as the quality of human translations exceeded by the quality of MT plus PE. More to that, a further optimization of the PE process needs to be done aiming at a time saving and cost-effective use of MT [13].

The authors in [18] and [19] brought out the idea of exploiting machine translation systems combined linearly using different paradigms has been successfully operated over SMT and rule-based MT (RBMT). As such, the systematic errors produced by the RBMT system were corrected by this automatic PE (APE) system based on PB-SMT, hence leading to the reduction of post-editing effort. For translation into a morphologically rich language, a rule (20 hand-written rules)-based approach for English-Czech MT outputs APE at the morphological level was applied by [20] and [21], based on the most frequent errors encountered in translation. Words morphosyntactic categories such as case, number, person, and gender as well as dependency labels are efficiently corrected by this approach. Intuitively, one useful way to improve the APE performance is by source-language information integration in APE. The author in [22] proposed a pipeline in order to overcome data sparsity issues, where through task-specific dense features the best pruned phrase table and language model are selected. More to that, they found that consistent improvements in all language pairs can be obtained by including source language information into statistical APE. The author in [23] considered the potential links of individual alignments occurrences and used an arbitrary number of alignments generated by different models (including both a refine model and minimum Bayes risk based models) by constructing over the 1-best alignments from multiple alignments [24] [25] weighted alignment matrices, rather than performing the combination of exactly two bidirectional

alignments as proposed [26] and [27]. The works presented by [28] were motivated based on the fact that word alignment quality is constraint by word alignment-based reordering of source words, with the principal objective of producing monotone source and target chunk alignments through the reordering of source chunks. We argue that the problem of long-range reordering can be reduced to only short-range, intra-chunk reordering by obtaining monotone chunk associations from monotone word alignments while some source language syntax is preserved. The assumption is founded on the reflection that translation is performed by human translators much preferably at chunk level rather than at the word level.

Also, translation outputs produced by an SMT were either re-ranked in a post-processing step using NMT [29] [30] [31] [32] [33], or used to produce an NMT system [10]. Another scenario involves re-ranking the translation outputs produced by an NMT in a post-processing step by using an SMT [12], or guiding translation in NMT by integrating an SMT into an NMT, as they revealed significant translation quality improvement over the Chinese-English translation tasks during experiments [9] [34]. In the works of [34], an NMT architecture is trained in an end-to-end manner where at each NMT decoding step, based on decoding information additional recommendations scored by an auxiliary classifier are offered by the SMT in order to generate words, and the SMT recommendations are combined with NMT generations exploiting a gating function while jointly taking part in the training process.

The several aforementioned attempts to improve MT system's performance did not still properly handled the issues faced by morphologically rich and low resourced languages, and long-term dependency modelling. We argue that, in order to limit the vocabulary size words could equally be split into sub-word units as proposed [35]. Also, lexical probabilities could be integrated into the NMT as successfully investigated [36]. Another latitude to achieve more monotone translation could be to exploit pre-reordering as experimented [37], and finally but not the least, the NMT translation of rare words could be improved in a post-processing step as suggested [38].

## III. WORD COMPOUNDING AND FACTORIZATION

In order to reduce the rate Out-Of-Vocabulary (OOV) occurrences and the amount of bilingual data when processing morphologically rich languages, factored models are majorly used. Factorization consists of splitting and retrieving from a given word linguistic information/factors such as dependency information, syntactic information, part-of-speech tags and lemma, using Tree-Tagger [39] and integrating it as a vector into a translation system. Machine translation from one morphological rich language to another has been a tedious task especially when not having enough required morphological information on the source side, since to have an exact target language word-form, word compounding is pronounced useful and highly productive [27] [40] [41] since it leads to sparse data problems and increases the vocabulary size. As such, integrating word compounding in the pre-processing phase has proven to be useful to add extra

morphological information to the linguistic/morphological factors of the source and target languages.

Compounding is operated at the level of POS Tag, where minimized part-of-speech tag are produced by refining POS-tags from the Tree-Tagger using a dependency parser to add morphological information including gender, number, case, verbs, person for nouns, definiteness, pronouns, determiners and adjectives, provided that both tools agreed on the POS-tag. And in case of disagreement the Tree-Tagger POS-tags were chosen. Morphologically rich language compound are formed by joining words, inserting filler letters (example: -s, -en, -er, -ien) or from the end of all but the last word remove letters (example: -en, -n) of the compound [42].

#### A. Compound Splitting

The morphologically rich data language model is POS tagged and employed to compute the adverbs, adjectives, negative particles, verbs and nouns frequencies. Then making use of the adjusted version of the corpus-based method proposed by [27], each adjective and noun splits into known words from the corpus also proper names are not split since it would give rise to errors if translated in parts, while permitting filler additions and truncations. Also, due to the fact that compound parts often contain the base form, lemmas are equally used to calculate word frequencies in addition to surface form. As hint, more splits are gotten when using the arithmetic mean of the frequencies of its parts rather than the geometric mean, where the highest arithmetic value is validated. Each compound parts length was limited to 4 characters and the number of parts for adjectives particularly was restricted to  $\leq 2$  with minimum words length to be split  $\geq 7$ .

All compound parts but the last were marked with the symbol # so as to be handled as separate words.

Special POS-tag are assigned to split words parts based on the compounds last word's POS, with both the full word and the last part receiving the same POS. Finally, words containing hyphens are split based on this same algorithm, and different POS-tags are assigned to their parts, with hyphens left at the end of all but last part. Factorization with compound splitting is integrated in a pre-processing step for training and translation of both the Transformer framework and the Phrase-based statistical machine framework.

#### B. Compound Merging

For translation into the morphologically rich language, the split compounds are merged based on POS through a post-processing step at the outputs of both the Transformer framework and the Phrase-based statistical machine framework. As such, if a compound-POS is possessed by a word and a matching POS possessed by the following word, they are merged. Alternatively, a hyphen is added to the word in case the next POS does not match, thus allowing for coordinated compounds.

We used the merging algorithm proposed by [41] based on [40], with this algorithm the advantage is that unseen compounds can be merged and coordinated compounds handled.

### IV. INCORPORATING LINGUISTIC FACTORS INTO THE TRANSFORMER

Our principal innovation over the standard encoder decoder based Transformer architecture is that we express the encoder input and decoder output as a combination of features such as [43] [44] [45]. Our generalized model supports an arbitrary number of input features.

It is on a number of well-known linguistic features that we focused in this paper, having as empirical question of knowing to which extend does providing linguistic features to both encoder and decoder improves the translation quality more specially in morphologically richer languages when using the transformer paradigm.

In order to better integrate linguistic factors in our NMT framework, we extended the Transformer architecture propose by [1] which employs multiple stacked layers of an encoder-decoder structure. Two sub-layers constitute the encoder layer, which are a self-attention sub-layer succeeded by a position-wise feed-forward sub-layer. Similarly to the encoder, the decoder has an additional sub-layer which serves at preventing information about future output positions to be incorporated by a given output position during training through masking in its self-attention. For all positions in a sequence, the transformer model computes attention scores using as query each position's input representation. The input representations weighted average are computed then using the previously obtained attention scores. More generally, the attention is identified as query and key/value vector pairs mapping to an output. As such, our work is an extension of [1] by the integration of additional linguistic factors. Considering that we have  $L$  layers of annotations for linguistic factors, and  $N$  training parallel sentences from the training data  $\{\{x^{(n,l)}, y^{(n,l)}\}_{l=0}^L\}_{n=1}^N$  where the  $n$ -th sentence pair word sequence is denoted in layer zero as  $x^{(n,0)}$  and its length denoted as  $|x^n|$ , the annotations of its  $L$  layers are denoted by  $\{x^{(n,l)}\}_{l=1}^L$ , with the target sentence denoted as  $y^{(n,l)}$ . In other words, for each feature we look up separate embedding vectors, and concatenate them. The total embedding size is matched by the concatenated vectors length, and the internal structure of the transformer's encoder and decoder is maintained. According to this setting we extended our standard encoder-decoder based Transformer architecture, operating as follows:

Given the input sequence  $x = (x_1, \dots, x_n)$  of  $n$  elements where  $x_i \in \mathbb{R}^{d_x}$  on which each attention head operates, and from which a new representation  $z = (z_1, \dots, z_n)$  of same length is computed where  $z_i \in \mathbb{R}^{d_z}$ . The weighted sum of a linearly transformed input elements will be computed from each output representations as [46]:

$$z_i = \sum_{j=1}^n \alpha_{ij} (x_j W^V) \quad (1)$$

Equally, a softmax function is used to compute each weight coefficient,  $\alpha_{ij}$  as:

$$\alpha_{ij} = \frac{\exp e_{ij}}{\sum_{k=1}^n \exp e_{ik}} \quad (2)$$



And compatibility function which compares two input elements is used computed from  $e_{ij}$ :

$$e_{ij} = \frac{(x_i W^Q)(x_j W^K)^T}{\sqrt{d_z}} \quad (3)$$

To enable efficient computation, a scaled dot product was chosen for compatibility function. Where we have as unique parameter matrices per layer  $W^Q, W^K, W^V \in \mathbb{R}^{d_x \times d_z}$ .

Input representations in multi-headed self-attention are linearly mapped to lower-dimensional spaces firstly, and one multi-headed self-attention layer's output is formed by the concatenation of several attention mechanisms output vectors (provided that each attention mechanism is identified as a head). Thus in the first self-attention layer the vector for position  $i$  for a single attention head  $\vec{h}$  is computed as:

$$\vec{h}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \quad (4)$$

And multi-head attention computed as:

$$\text{Multi}\vec{h}(Q, K, V) = \text{Concat}(\vec{h}_1, \dots, \vec{h}_h)W^O \quad (5)$$

given the function computing the resulting vector as:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (6)$$

Adding sinusoids of various wave-lengths enables the self-attention paradigm to encode positional information.<sup>1</sup>

#### A. Beam Search Integration with Factors

We extended our beam search procedure in order to find the best sequences by dealing with multiple word features (outputs), for simplicity reasons we have one beam responsible for generating lemmas and another beam responsible for generating the concatenation of the different factors. With the help of a toolkit such as MACAON [47] or even the more specialized KyTea [48], we performed the grammatical and morphological analysis. While taking into consideration the context, the lemma and factors for each word is output using the MACAON/KyTea POS-tagger [49]. In the various outputs, the generation of the lemmas and factors are made in a synchronous stream thus leading to sequences with different length sizes, ending each after the generation of the  $\langle eos \rangle$  (end-of-sequence) symbol, and creating by such, multiple representations of the  $\langle eos \rangle$  symbol in an output word. Due to the fact that lemmas carry most of the meaning and are closer to the final objective, we constricted the length size of the factors sequence to be equal to that of the lemma sequence. This implies that when the lemma sequence generation has ended we stop the generation of factors while ignoring their  $\langle eos \rangle$  symbol, therefore avoiding both longer and shorter factors sequences.

In order to generate the next word in the sequence, the feedback (previous word) is employed taking into consideration its various features (outputs), in this case, in order to obtain full benefit of both feedback outputs we

performed the tanh (non-linear) transformation of both embedding concatenation, thus having more information and learning better by their combination. Given as:

$$Prev_w(y_{t-1}) = \tanh([y_{t-1}^L; y_{t-1}^F] \cdot W_{Prev_w}) \quad (7)$$

Where, the previous output  $y_{t-1}$  feedback is computed by  $Prev_w, W_{Prev_w}$  are trained weights, with  $y_{t-1}^L$  and  $y_{t-1}^F$  the embedding of the lemma and factors generated at previous time step, respectively.

Finally, for each partial hypothesis we did the cross product of the output spaces of the best generated lemma and factors hypotheses, thus associating each factor hypothesis to each lemma hypothesis. Also, having  $k$  as beam size, the  $k$  – best combinations was kept for each sample. Equally, in order to get the word candidate when having the lemma and some factors, the MACAON toolkit was used. In situations where name entities are processed therefore having no factors found, the lemma was outputted by the system.

## V. HYBRID MACHINE TRANSLATION SUCCESSION

Although the translations produced by NMT are more fluent than those of SMT, it still does not fully and explicitly exploit the source information as compared to SMT. Thus, sometimes generating translations that are quite different from the source sentence original meaning [50] and some other times may mistakenly ignore some words during source sentence translations causing other words to be repeatedly translated [51].

If we consider as “intermediate language (another language)” the translation produced by the output of the NMT, to some extent we may amend the duplicated and meaningless translations, by building a translation model and operating a word alignment using an SMT.

Therefore we propose a factored multi-engine hybrid MT system consisting of an NMT and SMT framework, illustrated in Fig. 1.

Firstly, a preprocessing phase in this pipeline is performed by the transformer, which consist of training the transformer system using the initial factored training data, translating the training data, development set and test set into factored pre-translations; secondly, a target-target SMT system is built using the factored pre-translated training data, with parameters tuned using the pre-translated development set; and finally, the desired output is produced by decoding the pre-translated test set using the tuned SMT system.

When using the transformer to perform the pre-translations, if there is an occurrence of OOV in the source sentence, an ‘UNK’ token is generated by the transformer when translating the training data, development set and test set. We therefore propose a simpler and efficient technique to replace in the translation sentence, the “UNK” token by the corresponding source word. This method is known as the “labeled UNK replacement algorithm”, which alleviates the weaknesses faced by the UNK replacement algorithm inspired from [52] proposed by [12]. The technic is presented in Algorithm 1.

<sup>1</sup>We exploited relative positional encoding as emphasized [46] [60] so as to improve performance with respect to machine translation and relation classification, respectively.

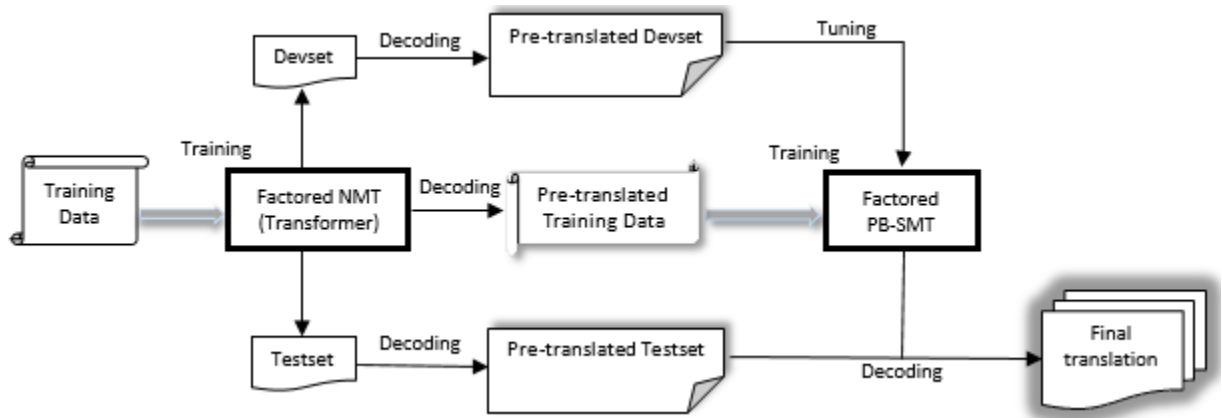


Fig. 1. The Factored Hybrid Transformer-PBSMT Framework.

As such, our algorithm will simply traverse the translation and replace the UNK token they encounter with their corresponding source word (the key at that position), if in the vocabulary there is no existence of the source word. Reference [12] proposed a naïve algorithm to do the UNK replacement, facing the weakness of eventually having between the source sentence and the target sentence different word order, thus creating wrong replacements.

---

**Algorithm 1** Labelled UNK replacement by source words

---

**Require:** The translation  $e_1^m$  with UNK tokens from the Transformer.

With  $e$  as an array of key and value pairs where each key is a source word, and the value the corresponding translation or UNK token.

e.g. for the French sentence: “un chat noir” to translate in English, we may have the corresponding  $e$  below.

|        |              |            |           |
|--------|--------------|------------|-----------|
| {un:a} | {noir:black} | {chat:UNK} | {eos:eos} |
|--------|--------------|------------|-----------|

Considering that we have the following sentence, the replacement will thus be done as:

```
1: procedure LABELLED UNK REPLACEMENT
2:   for  $i = 1$  to  $m$  do
3:     if  $e_i.value == UNK$  then
4:        $e_i.value = e_i.key$ 
5:     end if
6:   end for
7: end procedure
```

---

Finally, to post-process these unknown words, instead of using a back-off dictionary [53], we engage by considering more context a factored phrase based SMT system to perform the desired translation. In the factored PB-SMT, for decoding and training we applied the Moses toolkit [7], for sequence models we used SRILM [46] to train a 5-gram language model, and for word alignments creation we employed Giza++ [54], using for feature weights tuning the MERT (Minimum Error Rate Training) [26].

## VI. EXPERIMENTS

In order to verify our proposed framework, we selected translations between Japanese and English languages, noting that Japanese is drastically different in terms of word order

and has a far richer grammatical structure as compared to English language.

For fair comparison we re-implemented the hybrid frameworks proposed by [10], training our models using a machine with 8 NVIDIA P100 GPUs.

### A. Datasets and Setup

We used as training data Part-1 of the JP-EN Scientific Paper Abstract Corpus (ASPEC-JE) for JP-EN translation task which contains 1M sentence pairs, with the 1,790 sentence pairs contained in the development/validation set, and the 1,812 sentence pairs contained in the test set [55], provided that for the validation and test sets each sentence at the source side has only one reference.

For factorization at both the Transformer and the PBSMT level, we used Lemma and POS tags with compounds (as explained in Section 2 above) as input and output features, which can be produced either by using the MACAON toolkit [47] or the more specialized KyTea [48] especially for the Japanese data.

Due to the fact that unknown words cannot be generated when using Byte-Pair Encoding (BPE) [56] since they are all encoded as BPE units, we thus keep words as translation units. Besides this, incorrect words are sometimes produced by BPE units generation during the final word level processing, thus does not lead to any noticeable improvement in terms of %BLEU [57]. We used the PB-SMT system described in section 4 above. Also, we used as NMT system the transformer [1] default settings with some variants, setting mini batches of size 80, and having as 60 the maximum length of a sentence, with a size of 600 for word embeddings. Parallely, we have as input and output vocabulary size set to 45K. We reshuffled the training corpus between epochs, and trained the models with the AMSGrad optimizer [58], while at every 5,000 mini batches on the validation set, we validated the model through BLEU (BiLingual Evaluation Understudy) scores, and at every 30,000 performed model safeguard.

We only utilize the baseline transformer system pre-translated training data and devset as input to the SMT engine for its training and tuning. For tuning, the optimized configuration file settings for our translation model is found using Batch MIRA (equally known as k-best MIRA) [59] [60],

which is a version of MIRA (a margin-based classification algorithm) working within a batch tuning framework when we have sparse features OR using Minimum error rate training (MERT), but the use of more than about 20-30 features cannot be supported. After which the pre-translated test set is re-decoded utilizing the tuned SMT system.

### B. Evaluation and Results

Through bootstrap re-sampling significance test we calculated the statistical significance [61], and also, case-insensitive BLEU scores were used to report all results.

Table I shows the BLEU score based translation results for  $JP \leftrightarrow EN$  with non-reordered data, considering as baseline systems a standard PB-SMT [62] for statistical based translations and a NMT proposed by [3] for neuronal based translations. Thus, we observe that:

- The hybrid translation system where the SMT system is used to pre-translate data which serves as input to the NMT, performs significantly gets worse than both the baseline NMT systems and the FNMT system, when operating on  $JP \rightarrow EN$  and  $EN \rightarrow JP$  languages. The baseline SMT systems has been outperformed in %BLEU points by all the SMT $\Rightarrow$ NMT systems on  $JP \rightarrow EN$  and  $EN \rightarrow JP$ , except for the  $JP \rightarrow EN$  validation set which reports a decrease in result of **- 0.18 BLEU** points.
- The hybrid NMT $\Rightarrow$ SMT model results indicates that the translations produced by the baseline NMT system are re-decoded by the NMT $\Rightarrow$ SMT pipeline, leading to a significant improvement of **+ 1.25 BLEU** points and **+ 1.13 BLEU** points on the  $JP \rightarrow EN$  validation and test sets translation performance, respectively, and also, a significant improvement of **+ 1.41 BLEU** points and **+ 1.96 BLEU** points on the  $EN \rightarrow JP$  validation and test sets translation performance, respectively, compared to the baseline NMT system. As compared to the factored NMT system, the hybrid Factored NMT $\Rightarrow$ SMT model results indicates a slight but noticeable improvement of **+ 0.41 BLEU** points and **+ 0.42 BLEU** points on the  $JP \rightarrow EN$  validation and test sets translation performance, respectively, and also, a significant improvement of **+ 1.25 BLEU** points and **+ 1.65 BLEU** points on the  $EN \rightarrow JP$  validation and test sets translation performance, respectively.
- Finally, we observe that the hybrid model where translations produced by the factored transformer at both its input and output (fully-factored transformer), and which are further re-decoded by the factored SMT, outperforms the translations on the  $JP \rightarrow EN$  validation set generated by the fully-factored transformer, and the transformer, by **+ 0.86 BLEU** points and **+ 2.74 BLEU** points, respectively, and also, translations on the  $JP \rightarrow EN$  test set generated by the fully-factored transformer, and the transformer, by **+ 1.04 BLEU** points and **+ 2.86 BLEU** points, respectively. Similarly, both the translations on the  $EN \rightarrow JP$  validation set generated by the fully-

factored transformer, and the transformer, by **+ 1.06 BLEU** points and **+ 2.66 BLEU** points, respectively, and those on the  $EN \rightarrow JP$  test set generated by the fully-factored transformer, and the transformer, by **+ 1.27 BLEU** points and **+ 3.25 BLEU** points, respectively, are as such outperformed by our proposed hybrid system.

### C. Discussion

From the above results with reference to the state of the art, we analyze that:

As compared exceptionally to [10] framework consisting of an SMT $\Rightarrow$ NMT pipeline which has a higher computational complexity due to the integration of the source information into both the SMT and NMT (concatenating at this level the pre-translated and source sentences as input), and other state of the art hybrid frameworks particularly [12] consisting of an NMT $\Rightarrow$ SMT pipeline, our hybrid MT pipeline is more simpler, viable and efficient, by employing source-side information only during the transformer training and exceptionally during OOVs processing, thus favoring its faster computation. Analytical studies for rare/OOV word impact on the translation quality were operated over the Scientific Paper Abstract Corpus (ASPEC-JE) for Japanese-to-English, sorted by the words average inverse frequency and validation sentences were split into groups with comparable numbers of rare words independently evaluated. All target words which occur in the training data for each number of sentence occurrence less than N times were replaced by the UNK token, for all analyzed systems. Given  $N \in \{0K, 0.5K, 1K, 1.5K, 2K, 2.5K, 3K\}$ . Thus, a higher occurrence of rare words is obtained for large N, hence in the reference only the most frequent words are exploited. Meanwhile a lesser occurrence of rare words is obtained for lower N, using hence more words. We observed that our best performing model (FF-Transformer $\Rightarrow$ FSMT) considerably outperforms the state of the art both stand-alone and hybrid MT systems on sentences with many OOV words, as a greater occurrence of OOV words implies an increased amount of data size. This boost in performance can be justified by the fact that attention mechanisms which makes up the Transformer operates better on larger data sizes.

We point out that, attention mechanisms are used by neural networks to encode each position while relating two distant words of both the inputs and outputs with respect to itself, by which the training can be accelerated through parallelization. An attention mechanism is a technique created for paying attention to specific words, which have proven to be useful to address the bottleneck issues that arise when handling long sentences with complicated dependencies between words, as it is harder for the context vector to capture all the information contained in the sentence due to the sequential order of word processing. More precisely, the Attention technique focuses on part of a subset of the information it is given, provided that for each input word one hidden state vector is produced. These vectors can then be concatenated, averaged or (even better!) weighted in order to give higher importance to words from the input sentence, most relevant to decode the next word of the output sentence.

TABLE I. RESULTS OF VARIOUS HYBRID (NMT-SMT) MACHINE TRANSLATION EXPERIMENTS PERFORMED ON JP→EN AND EN→JP WHERE, “♣” INDICATES THE BEST TRANSLATION PERFORMANCE

| SYSTEM                            | JP-EN      |        | EN-JP      |        |
|-----------------------------------|------------|--------|------------|--------|
|                                   | Validation | Test   | Validation | Test   |
| SMT [62]                          | 18.46      | 17.79  | 27.71      | 26.54  |
| NMT [3]                           | 24.66♣     | 24.94♣ | 35.72♣     | 35.48♣ |
| FACTORED SMT                      | 18.87      | 17.91  | 27.84      | 26.80  |
| FACTORED NMT                      | 25.70♣     | 25.83♣ | 36.57♣     | 36.31♣ |
| SMT⇒NMT [10]                      | 18.28      | 17.92  | 27.82      | 27.98  |
| T⇒SMT                             | 25.91      | 26.07  | 37.13      | 37.44  |
| FACTORED NMT⇒SMT [12]             | 26.11♣     | 26.25♣ | 37.82♣     | 37.96♣ |
| TRANSFORMER [1]                   | 31.98      | 32.09  | 42.16      | 42.41  |
| FULLY-FACTORED TRANSFORMER        | 33.86      | 33.91  | 44.01      | 44.39  |
| FULLY-FACTORED TRANSFORMER ⇒ FSMT | 34.72♣     | 34.95♣ | 45.07♣     | 45.66♣ |

Also, due to the larger vocabulary of the test set by the integration of factors during the PB-SMT post-processing translation, we experienced in our proposed framework a significant decrease in rate of OOVs as compared to the NMT system, of 1.06% and 5.37%, respectively.

We emphasize that, the results on the ASPEC Japanese-to-English corpus should be interpreted with caution. It is the expectation that the attention based HMT when used on longer sentences will show their true potential. In order to investigate on the effect of translating long sentences, sentences of similar lengths having unknown words to the models included were grouped together and the BLEU score was computed per group. The results are delineated in Fig. 2, analyzed over the full validation set.

We observe on Fig. 2 that the buckets of longer sentences are more effectively handled by our Transformer based HMT

(purple curve) due to its integrated Attention mechanism at both the encoder and decoder levels as compared to the winning entry recurrent based HMT (green curve) in which the Attention mechanism is integrated only at the level of the decoder, hence as sentences become longer the quality does not degrade. While at shorter sentence lengths, it is observed that our outperforming model performs worse, indicating that although the attention mechanism speeds up training, it is likely not very important and may potentially be redundant. More to that, higher perplexities are produced when operating Attention mechanisms over short sentences, as the model becomes less certain about its predictions than without it.

And we believe that, translations performance will be improved if phrases corrected and reordered are considered. We shall dive deeper by considering this fact in future work.

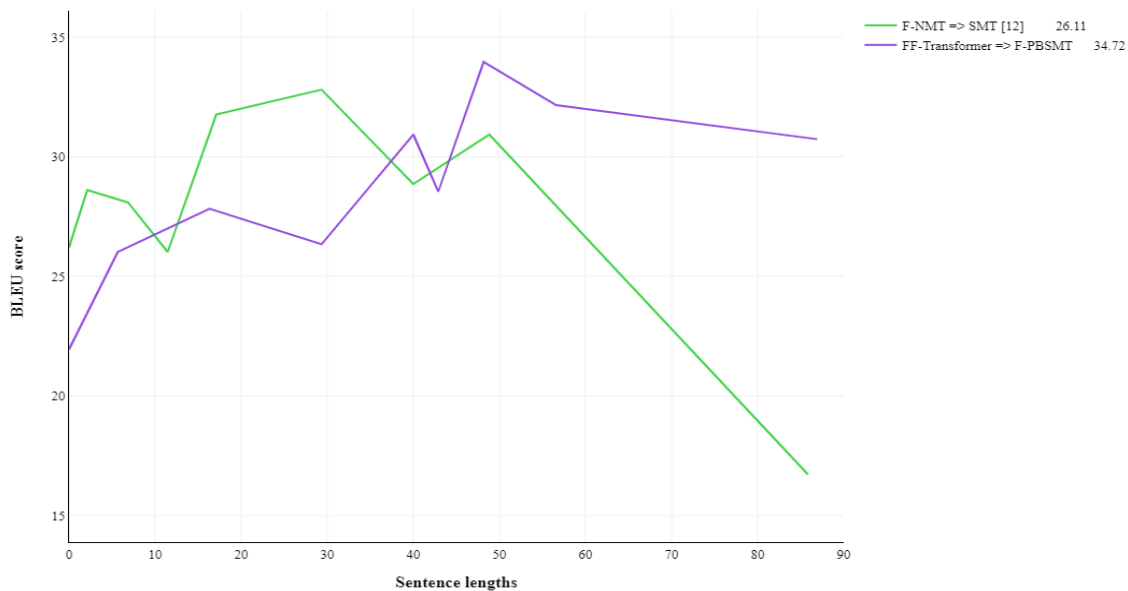


Fig. 2. Length Analysis – Impact of Attention Mechanism on Translation Qualities as Sentences become Longer Performed on ASPEC-JE Data.

## VII. CONCLUSION

We have proposed a novel HMT framework cascaded as a Fully-Factored Transformer  $\Rightarrow$  Factored SMT pipeline consisting of integrated linguistic factors at both the source language and target language of the transformer model, and linguistic factors at source language (pre-translated language) of the SMT model. The considered linguistic factors were lemmatization, part-of-speech tagging (taking into consideration its various compounds). Our experimental results on *JP*  $\leftrightarrow$  *EN* language pairs clearly revealed that our proposed HMT framework with integrated linguistic factors outperforms the state-of-the-art HMT frameworks, in terms of both perplexity and BLEU points. More to that, we observed an OOV rate reduction, due to the generation of new word forms derived from the integrated additional linguistic resources.

As future work, we aim to explore whether the integration of a grammatical error detection and correction (GEC) process [34] will further help in reducing the rate of OOVs. Also, use compositional learned word representations from smaller orthographic symbols inside the words such as character n-grams, which can easily fit in the model vocabulary.

## VIII. CONFLICTS OF INTEREST STATEMENT

The authors whose names are listed above certify that they have NO affiliations with or involvement in any organization or entity with any financial interest (such as honoraria; educational grants; participation in speakers' bureaus; membership, employment, consultancies, stock ownership, or other equity interest; and expert testimony or patent-licensing arrangements), or non-financial interest (such as personal or professional relationships, affiliations, knowledge or beliefs) in the subject matter or materials discussed in this manuscript.

## REFERENCES

- [1] Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, Kaiser Ł, and Polosukhin I. Attention is all you need. In Proceedings of NIPS, 2017.
- [2] Devlin J, Zbib R, Huang Z, Lamar T, Schwartz R, Makhoul J. Fast and robust neural network joint models for statistical machine translation. In: Proceedings of the 52st Annual Meeting of the Association for Computational Linguistics (ACL 2014), vol 1, pages 1370–1380, 2014.
- [3] Bahdanau D, Cho K, Bengio Y. Neural machine translation by jointly learning to align and translate. In Proceedings of the International Conference on Learning Representations (ICLR), 2015.
- [4] Koehn Philipp and Hoang Hieu. Factored Translation Models. In: Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL), Association for Computational Linguistics, pages 868–876, 2007.
- [5] Kalchbrenner Nal and Blunsom Phil. Recurrent continuous translation models. EMNLP. Pages 1700–1709, 2013.
- [6] Stymne S, Holmqvist M, and Ahrenberg L. Effects of morphological analysis in translation between German and English. In: Proceedings of the Third ACL Workshop on Statistical Machine Translation, 2008.
- [7] Koehn P, Hoang H, Birch A, Callison-Burch C, Federico M, Bertoldi N, Cowan B, Shen W, Moran C, Zens R, Dyer C, Bojar O, Constantin A, Herbst E. Moses: Open source toolkit for statistical machine translation. In: Proceedings of the 45th Annual Meeting of the ACL, demonstration session, pages 177–180, 2007.
- [8] Cho Kyunghyun, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder-decoder for statistical machine translation. In: Proceedings of the Conference on Empirical Methods on Natural Language Processing, pages 1724–1734, 2014.
- [9] He W, He Z, Wu H, Wang H. Improved neural machine translation with SMT features. In: Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, pages 151–157, 2016.
- [10] Niehues J, Cho E, Ha TL and Waibel A. Pre-translation for neural machine translation. In: Proceedings of the COLING, pages 1828–1836, 2016.
- [11] Vaswani Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In: Proceedings of NIPS, 2017.
- [12] Du Jinhua and Way Andy. Neural pre-translation for hybrid machine translation. In: Proceedings of MT Summit XVI, MT Summit XVI, vol 1, pages 27–40, 2017.
- [13] Report T. Maschine translation post-editing guidelines published (Technical report). TAUS, 2010. Retrieved from <http://www.cnlg.org.cn/tauscngl-machine-translation-post-editing-guidelines-published>.
- [14] Roturier J. Deploying novel MT technology to raise the bar for quality: a review of key advantages and challenges. In Proceedings of the twelfth machine translation summit, 2009.
- [15] Fiederer R., and OBrien S. Quality and machine translation: a realistic objective. Journal of Specialised Translation, 11, pages 52–74, 2009.
- [16] Koehn P. A process study of computer-aided translation. Machine Translation, 23(4), pages 241–263, 2009.
- [17] Palma D. A., and Kelly N. Project management for crowd sourced translation: How user-translated content projects work in real life. In Translation and localization project management: The art of the possible, pages 379–408, 2009.
- [18] Dugast L., Senellart J., and Koehn, P. Statistical post-editing on systran's rule-based translation system. In Proceedings of the second workshop on statistical machine translation. Prague, Czech Republic: Association for Computational Linguistics, pages 220–223, 2007.
- [19] Simard M., Goutte C., and Isabelle P. Statistical phrase-based post-editing. In Proceedings of naac1-hlt, pages 508–515, 2007.
- [20] Rosa R., Mareček D., and Dušek, O. Depfix: A system for automatic correction of czech mt outputs. In Proceedings of the seventh workshop on statistical machine translation, pages 362–368, 2012.
- [21] Mareček D., Rosa R., Galuščáková P., and Bojar O. Two-step translation with grammatical post-processing. In Proceedings of the sixth workshop on statistical machine translation, pages 426–432, 2011.
- [22] Chatterjee R., Turchi M., and Negri M. The fbk participation in the wmt15 automatic postediting shared task. In Proceedings of the tenth workshop on statistical machine translation, Lisbon, Portugal, pages 210–215, 2015. Retrieved from <http://aclweb.org/anthology/W15-3025>. (Association for Computational Linguistics. URL).
- [23] Tu Z., Liu Y., He Y., Genabith J., Liu Q., and Lin S. Combining multiple alignments to improve machine translation. In The 24th international conference of computational linguistics, pages 1249–1260, 2012.
- [24] Liu Y., Xia T., Xiao X., and Liu, Q. Weighted alignment matrices for statistical machine translation. In Proceedings of the 2009 conference on empirical methods in natural language processing, Singapore: Association for Computational Linguistics, pages 1017–1026, 2009.
- [25] Tu Z., Liu Y., Liu Q., and Lin S. Extracting hierarchical rules from a weighted alignment matrix. In Proceedings of 5th international joint conference on natural language processing, pages 1294–1303, 2011.
- [26] Och FJ. Minimum error rate training in statistical machine translation. In: Proceedings of the 41st Annual Meeting of ACL, pages 160–167, 2003.
- [27] Koehn Philipp and Knight K. Empirical methods for compound splitting. In: Proceedings of the tenth conference of EACL, pages 187–193, 2003.
- [28] Pal Santanu. A Hybrid Machine Translation Framework for an Improved Translation Workflow. Saarland University, Saarbrücken, Germany, 2018.

- [29] Ding S., Duh K., Khayrallah H., Koehn P., and Post M. The jhu machine translation systems for wmt 2016. In Proceedings of the conference on statistical machine translation, Berlin, Germany, 2016.
- [30] Farajian M., Chatterjee R., Conforti C., Jalalvand S., Balaraman V., Gangi M., and Federico M. In Fbk's neural machine translation systems for iwslt 2016. In Proceedings of the 13<sup>th</sup> workshop on spoken language translation, Seattle, USA, pages 8–15, 2016.
- [31] Lee H. G., Lee J., Kim, J. S., and Lee, C. K. NAVER machine translation system for wat 2015. In Proceedings of the 2nd Workshop on Asian Translation (WAT2015), Kyoto, Japan, pages 69–73, 2015.
- [32] Neubig G., Morishita M., and Nakamura S. Neural re-ranking improves subjective quality of machine translation: NAIST at WAT2015. In Proceedings of the 2nd Workshop on Asian Translation (WAT2015), Kyoto, Japan, pages 35–41, 2015.
- [33] Zhao Y., Huang S., Chen H., and Chen J. Investigation on statistical machine translation with neural language models. In Proceedings of Chinese Computational Linguistics and Natural Language Processing Based on Naturally Annotated Big Data, pages 175–186, 2014.
- [34] Wang X, Lu Z, Tu Z, Li H, Xiong D, Zhang M. Neural machine translation advised by statistical machine translation. In: Proceedings of the AAAI Conference on Artificial Intelligence, San Francisco, California, USA, 2017.
- [35] Sennrich Rico, Haddow Barry, and Birch Alexandra. Edinburgh Neural Machine Translation Systems for WMT16. In Proceedings of the 1st conference on machine translation, Berlin, Germany, pages 368–373, 2016a.
- [36] Philip Arthur, Graham Neubig, and Satoshi Nakamura. Incorporating discrete translation lexicons into neural machine translation. In Conference on Empirical Methods in Natural Language Processing (EMNLP), Austin, Texas, USA, November 2016.
- [37] Kay Rottmann and Stephan Vogel. Word Reordering in Statistical Machine Translation with a POS-Based Distortion Model. In Proceedings of the 11th International Conference on Theoretical and Methodological Issues in Machine Translation (TMI 2007), Skövde, Sweden, 2007.
- [38] Thang Luong, Hieu Pham, and Christopher D. Manning. Effective approaches to attention-based neural machine translation. In Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, Lisbon, Portugal. Association for Computational Linguistics, pages 1412–1421, 2015.
- [39] Schmid H. Probabilistic part-of-speech tagging using decision trees. In: Proceedings of the International Conference on New Methods in Language Processing, pages 44–49, 1994.
- [40] Popović M, Stein D and Ney H. Statistical machine translation of German compound words. In: Proceedings of FinTAL - 5th International Conference on Natural Language Processing, pages 616–624, 2006.
- [41] Stolcke A. SRILM - an extensible language modeling toolkit. In: Proceedings of the International Conference on Spoken Language Processing (ICSLP), pages 901–904, 2002.
- [42] Langer S. Zur Morphologie und Semantik von Nominalkomposita. Tagungsband der 4 Konferenz zur Verarbeitung natürlicher Sprache (KON- VENS) pages 83–97, 1998.
- [43] Alexandrescu A, Kirchoff K. Factored Neural Language Models. In: Proceedings of the Human Language Technology Conference of the NAACL, Companion Volume: Short Papers, Association for Computational Linguistics, pages 1–4, 2006.
- [44] Wang Y, Wang L, Wong DF, Chao LS, Zeng X, Lu Y. Factored Statistical Machine Translation for Grammatical Error Correction. In: Proceedings of the Eighteenth Conference on Computational Natural Language Learning: Shared Task, pages 83–90, 2014.
- [45] Youzheng Wu, Xugang Lu, Hitoshi Yamamoto, Shigeki Matsuda, Chiori Hori, Hideki Kashioka. Factored Language Model based on Recurrent Neural Network. In Proceedings of COLING 2012: Technical Papers, pages 2835–2850, 2012.
- [46] Shaw P, Uszkoreit J, Vaswani A. Self-Attention with Relative Position Representations. arXiv preprint, 2018.
- [47] Nasr A, Béchet F, Rey JF, Favre B and Roux JL. Macaon, an nlp tool suite for processing word lattices. In Proceedings of the ACL-HLT, pages 86–91, 2011.
- [48] Graham Neubig and Shinsuke Mori. Word-based partial annotation for efficient corpus construction. In Proceedings of the 7th International Conference on Language Resources and Evaluation, 2010.
- [49] Neubi G, Nakata Y and Mori S. Pointwise prediction for robust, adaptable japanese morphological analysis. In: Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies, pages 529–533, 2011.
- [50] Sutskever Ilya, Vinyals Oriol and Le Quoc V. Sequence to sequence learning with neural networks. Neural Information Processing Systems (NIPS) Montréal pages 3104–3112, 2014.
- [51] Toral A, Sánchez-Cartagena V. M. A multifaceted evaluation of neural versus phrase-based machine translation for 9 language directions. In: Proceedings of the Conference of the European Chapter, Association for Computational Linguistics, 2017.
- [52] Jean S, Cho K, Memisevic R, Bengio Y. On using very large target vocabulary for neural machine translation. In: Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing, pages 1–10, 2015.
- [53] Luong T, Sutskever I, Le Q, Vinyals O, Zaremba W. Addressing the rare word problem in neural machine translation. In: Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing, pages 11–19, 2015b.
- [54] Och FJ, Ney H. A systematic comparison of various statistical alignment models. *Computational Linguistics* 29(1): pages 19–51, 2003.
- [55] Nakazawa T, Yaguchi M, Uchimoto K, Utiyama M, Sumita E, Kurohashi S and Isahara H. ASPEC: Asian Scientific Paper Excerpt Corpus. In: Proceedings of the Tenth International Conference on Language Resources and Evaluation, Portorož, Slovenia, 2016.
- [56] Sennrich R, Haddow B, and Birch A. Neural machine translation of rare words with subword units. 2015.
- [57] Papineni K, Roukos S, Ward T and Zhu WJ. BLEU: a method for automatic evaluation of machine translation (PDF). ACL-2002: 40th Annual meeting of the Association for Computational Linguistics 9416:311– 318, CiteSeerX 10.1.1.19, 2002.
- [58] Sashank J, Reddi Satyen Kale and Sanjiv Kumar. On the convergence of Adam and beyond. Published as a conference paper at ICLR, 2018.
- [59] Cherry Collin and Foster George. Batch Tuning Strategies for Statistical Machine Translation. National Research Council Canada, 2012.
- [60] Bilan I, Roth B. Position-aware Self-attention with Relative Positional Encodings for Slot Filling. arXiv preprint, 2018.
- [61] Koehn Philipp. Statistical significance tests for machine translation evaluation. In: Proceedings of the Conference on Empirical Methods on Natural Language Processing, pages 388–395, 2004.
- [62] Koehn Philipp, Franz Josef Och, and Daniel Marcu. Statistical Phrase-Based Translation. HLT/NAACL, 2003.

# Reward-Based DSM Program for Residential Electrical Loads in Smart Grid

Muthuselvi G<sup>1</sup>

Research Scholar

School of Electrical Engineering

Vellore Institute of Technology, Vellore, Tamilnadu, India

Saravanan B<sup>2</sup>

Associate Professor

School of Electrical Engineering

Vellore Institute of Technology, Vellore, Tamilnadu, India

**Abstract**—There is a positive attitude towards the use of different strategies for engaging in demand response (DR) programs in energy markets through the innovation and trend of smart grid technologies. In this paper, a reward-based approach is proposed that enhances the involvement of customers in the DR program by assuring the customer's comfort level. Most of the previous works considered limited controllable loads like thermal loads for demand side management (DSM). In this approach thermal and all active electrical loads are considered for the analysis. Comfort indicator is used for the analysis which is an important parameter for measuring comfort of each resident. This technique significantly reduces the utility reward cost and maximizes the user satisfaction level compared with existing approach. The numerical example considered in this work illustrates the fruitfulness of the proposed approach. This problem is formulated as mixed-integer linear programming (MILP) and solved by using CPLEX solver in General Algebraic Modelling Software (GAMS).

**Keywords**—DSM; LSE; RLA; smart grid; reward

## I. INTRODUCTION

The solution to demand-supply problems in the power supply system is an efficient DR program. DSM is the customer based DR program in a smart electric grid by changing the regular use of electricity. Demand Side Management (DSM) programs enable load-serving entities (LSE) to manage the electric loads on the user side. Customer interaction and responsiveness are the two critical factors of the DR program. This program requires a collaborative relationship between LSE and consumers to achieve the customer load changes that benefit consumers, LSE and society as a whole. The categories of DSM programs are time-based and incentive/penalty based program. The time-based program depends on electricity prices that vary over time. The incentive program is based on fixed rewards (or) time-varying incentives. These programs play a significant role in reducing demand in peak periods. With the introduction of the new reward-based DR program consumers are encouraged to reduce their loads during peak hours.

The DR program classification is as shown in Fig. 1. The dynamic pricing schemes like Real Time Pricing (RTP), Time Block Pricing (TBP), Critical Peak Pricing (CPP), and Time of Use (ToU) are commonly used in the DR program to lower peak demand by encouraging the users. In [1], the author proposed a Home Energy Management System (HEMS) with price-based DR programs for reducing the consumer's

electricity cost by transferring the ON peak load usage to the OFF peak periods. During peak hours of the day, utilities have control over Direct Load Control (DLC) to shed the load of the consumer. The utility commonly control loads of the customer remotely in Interruptible / Curtailable Service (CS) to achieve the required load reduction level. Several researchers [2-9] suggested the implementation of optimal DR scheduling via incentives. Demand Bidding (DB) is a process that encourages consumers to involve actively in electricity usage trading. It offers incentives for accepting to reduce their electricity usage during peak load periods [10], [11]. Incentive-based DR programs is the most powerful tool for handling peak load and attract more consumers in the DR program. There are several methods and strategies explored to reduce peak demand, utility reward cost and consumer electricity costs. Existing methods have less concentration in the improvement of the participation factor of the consumer in the DR program.

DR program reduces the electricity cost of consumers by motivating them to use less power consumption in high-priced periods [12] and more power consumption in low-priced periods. The author demonstrated the DR program of residential areas by encouraging customers to voluntarily minimize their daily energy use by scheduling available resources in [13].

Optimization approaches to modify the customer's load curve in response to the changes in the electricity cost through incentive is analyzed in [14]. In this approach complexity faced on customer satisfaction level is analyzed at minimal percentile [15]. Previous studies focused mainly on the minimization cost at the consumer level and not dealt with the revenue cost of utility received from the grid operator [16], [17], and [19]. Smart Home Energy Management Systems (HEMS) uses real-time information under various schemes to manage loads of residential communities. Price based HEMS are discussed [23-25] under restricted controllable appliances. In [21], the author formulated an optimization approach as a Multi-Constrained Mixed-Integer Problem (MCMIP) that schedules the controllable appliances based on consumer preferences. A researcher proposed an Adaptive Differential Evolution Algorithm (ADEA) to find the optimal schedule of appliances in the sectors like residential, industrial, and commercial [22]. In this work, peak load minimization and reduction in consumer's electricity bills are considered as the objectives.

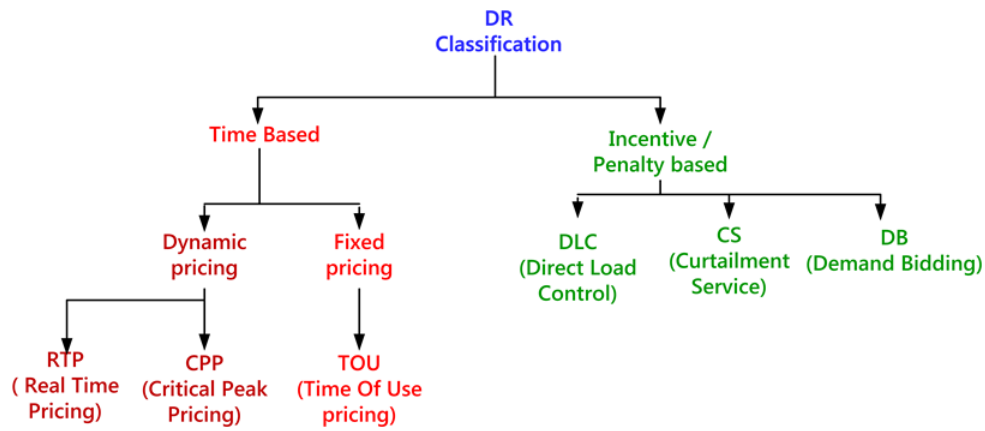


Fig. 1. Types of DR Programs.

The impacts, as mentioned earlier in existing methods, are rectified in this proposed approach. In this approach, controllable (electric vehicle, clothes dryer, dish washer, pool pump, energy storage), thermal and all active electrical loads are considered for the analysis. LSE sends the preferred demand reduction request (PDRR) to all the Residential Load Aggregator (RLA). Then RLA generates the optimal scheduling of appliances based on preferred demand reduction limits (DRL) and the willingness to compromise the load demand of all the residents. Customers are granted a reward depending on the participation in the DR program. New reward rates and demand reduction requests are informed to each consumer through home energy communication network port.

The remaining section of this article is structured as follows. Section 2 discusses the modelling structure of various appliances. Section 3 illustrates the reward structure used in the analysis. Section 4 with the problem formulation for the objectives and constraints. Section 5 deals with the case study considered in this approach. Finally, Sections 5 and 6 discuss the results and conclusion, respectively.

## II. OVERVIEW OF APPLIANCES MODELLING

This section describes the modelling of residential loads. Modelling parameters for appliances are considered from [18] and [20]. Controllable/Shiftable and uncontrollable/non-shiftable are the two classifications of appliances in the residential sector. Uncontrollable appliances have fixed power and time of operation. Controllable devices run within the desired working time based on their power consumption and time i.e. they may allow the operating schedule to be changed.

### A. Controllable Loads

1) *EWH Model*: When the current temperature exceeds the setpoint value, the condition of EWH is OFF. If the current temperature is below the minimum temperature required by EWH, then the status is ON. If the temperature is between the minimum required temperature and the setpoint temperature, then the EWH state follows the status of the previous time are shown in (1). The power consumption is calculated using equation (2).

$$S_{EWH,h,t} = \begin{cases} 0, T_{EWH,h,t} > T_{EWH,s} \\ 1, T_{EWH,h,t} < T_{EWH,r} \\ S_{EWH,h,t-1}, T_{EWH,r} \leq T_{EWH,h,t} \leq T_{EWH,s} \end{cases} \quad (1)$$

$$p_{EWH,h,t} = P_{EWH} \cdot S_{EWH,h,t} \quad (2)$$

2) *AC Model*: It is one of the significant thermal controlled residential appliances. When the room temperature is conquer ( $T_{AC,sp} + T_{AC,DB}$ ), AC is ON. When the room temperature is below ( $T_{AC,sp} - T_{AC,DB}$ ) the status of AC is OFF. Otherwise, it follows the previous status. If the state of AC is ON, the rated power is consumed. In the OFF state, no power will be consumed. The consumed power of AC for resident 'h' at time 't' is illustrated as follows:

$$S_{AC,h,t} = \begin{cases} 0, T_{AC,h,t} < T_{AC,sp} - T_{AC,DB} \\ 1, T_{AC,h,t} > T_{AC,sp} + T_{AC,DB} \\ S_{AC,h,t-1}, T_{AC,sp} - T_{AC,DB} \leq T_{AC,h,t} \leq T_{AC,sp} + T_{AC,DB} \end{cases} \quad (3)$$

$$p_{AC,h,t} = P_{AC} \cdot S_{AC,h,t} \quad (4)$$

3) *Clothes dryer model*: Equation (5) details the ON/OFF status of the CD. The state is ON when the total accumulated time is lower than the needed time to complete the job. The state is OFF when the accumulated time higher than or equal to the time necessary to initiate the job. The expression for power consumption of CD is in equation (6).

$$S_{CD,h,t} = \begin{cases} 0, T_{CD,acc} \geq T_{CD,r} \\ 1, T_{CD,acc} < T_{CD,r} \end{cases} \quad (5)$$

$$p_{CD,h,t} = P_{CD} \cdot S_{CD,h,t} \quad (6)$$

4) *Electric vehicle model*: The State-of-Charging (SoC) of the battery at a time 't' is the ratio of remaining  $EV_{rem}(t)$  or residual capacity at that time to the maximum battery capacity, as shown in equation (7). ON/OFF status of EV is illustrated in equation (8). It is ON, when the SoC is less than the maximum capacity and is OFF, when the SoC is greater than or equal to the maximum state of charge. The power consumption of EV is calculated using equation (9).



$$SOC(t) = \frac{EV_{rem}(t)}{EV_{max}} \quad (7)$$

$$S_{EV,h,t} = \begin{cases} 0, & SOC_{h,t} \geq SOC_{max} \\ 1, & SOC_{h,t} < SOC_{max} \end{cases} \quad (8)$$

$$p_{EV,h,t} = P_{EV} \cdot S_{EV,h,t} \quad (9)$$

5) *Dishwasher model*: The representation of the state of Dish Washer (DW) is, as shown in equation (10). The status of DW is ON when the cumulative ON time is lower than the needed ON time to complete dishwashing. It is OFF when the overall ON time greater than or equal to the needed ON time to finish that cycle. Equation (11) represents the power consumption of DW.

$$S_{DW,h,t} = \begin{cases} 0, & T_{DW,acc} \geq T_{DW,r} \\ 1, & T_{DW,acc} < T_{DW,r} \end{cases} \quad (10)$$

$$p_{DW,h,t} = P_{DW} \cdot S_{DW,h,t} \quad (11)$$

6) *Cloth washer model*: Equation (12) illustrates the status of Cloth Washer (CW). The condition of CW is ON when the cumulative ON time is lower than the needed time to complete the washing job. The status is OFF when the overall ON time greater than or equal to the needed ON time. The power consumption of CW is calculated using equation (13).

$$S_{CW,h,t} = \begin{cases} 0, & T_{CW,acc} \geq T_{CW,r} \\ 1, & T_{CW,acc} < T_{CW,r} \end{cases} \quad (12)$$

$$p_{CW,h,t} = P_{CW} \cdot S_{CW,h,t} \quad (13)$$

7) *Pool pump model*: The status of Pool Pump (PP) is represented in equation (14). The state is ON when the time taken for the operation is less than the desired total operating time and is OFF when its operating time exceeds the expected total running time. Power consumption of PP is obtained using equation (15).

$$S_{PP,h,t} = \begin{cases} 0, & T_{PP,acc} \geq T_{PP,r} \\ 1, & T_{PP,acc} < T_{PP,r} \end{cases} \quad (14)$$

$$p_{PP,h,t} = P_{PP} \cdot S_{PP,h,t} \quad (15)$$

### B. Uncontrollable Loads

Uncontrollable loads are loads that fix their mode operation in time and power consumption. The loads, including TV, computer, lighting loads, fan, and refrigerator, are examples considered in this category for the analysis. Each appliance's power ratings are regarded as of [19].

## III. REWARD-BASED DR FRAMEWORK

### A. Reward Rate Structure

In this proposed study, the reward-based DR model is formulated. In this framework, LSE sends the preferred demand reduction request to all RLA. Then RLA generates the optimal scheduling of appliances based on preferred demand reduction limits (DRL) and the willingness to compromise the load demand of all the residents are generated by RLA. Customers are granted a reward depending on the participation in the DR program. The reward rate  $Rw_2$  is given to the houses that are willing to compromise their demand during the PDRR event.  $Rw_3$  reward rate is awarded to the houses those who are not willing to compromise their demand during the PDRR event. But the  $Rw_3$  reward rate will be used only during the emergency (or) rare events. The choice between  $Rw_1$  and  $Rw_2$  plays a significant role in the optimization problem.

Table I represents the different reward-based rate structures. Here is a simple example that explains the reward rate for the residents based on their willingness to compromise and preferred demand reduction limits. In this, rewards  $Rw_1$ ,  $Rw_2$ , and  $Rw_3$  are considered as 20, 40, 60 cents / kW.

Here this include three houses A, B, and C. The preferred Demand Reduction Limits (DRL) of each resident is as shown in Table I. The total demand for all the residents is 35.6 kW. Assume LSE expects 30% of load reduction from RLAs. So RLA makes the optimal strategies that satisfy the requisition, which is given by LSE. In this example, 11.7 kW is expecting to reduce during that particular period. Depending on the DRL desired by each house, the LSE specification that met with 8.3 kW (total power minus DRL). Houses A and C are agreed to compromise their demands for the remaining kW. So the  $Rw_2$  reward rate is given to houses A and C because of their comfort index, which is higher than 1. House B is getting  $Rw_1$  reward rate because the willingness to compromise is '0'. An emergency is a rare occurrence case. Reward rate  $Rw_3$  is provided, if the compromise is equal to '0' and CI value is higher than 1. In this case, house 'B' will handle the emergency case and get the  $Rw_3$  reward rate.

### B. Comfort Index

The optimal scheduling of appliances of all the residents and the participation of customers depends upon the factor of CI. The design of the CI considers both thermal and other active appliances participating in the DR program. If the value of the CI is higher than '1', the residents are in an uncomfortable zone. When the CI is less than or equal to '1', the residents are in a comfortable area.

TABLE I. EXAMPLE OF REWARD RATE STRUCTURE

| House | Power(kW) | Demand reduction limit (kW) | Compromise (1=yes,0=No) | Reward (\$) |            |           |
|-------|-----------|-----------------------------|-------------------------|-------------|------------|-----------|
|       |           |                             |                         | Normal      | Occasional | Emergency |
| A     | 14.3      | 10.6                        | 1                       | $Rw_1$      | $Rw_2$     | -         |
| B     | 12.8      | 9.5                         | 0                       | $Rw_1$      | -          | $Rw_3$    |
| C     | 8.5       | 7.2                         | 1                       | $Rw_1$      | $Rw_2$     | -         |

The normalized value of CI is calculated using the equation (16).

$$CI_n = n1 CI_{H,h,t} + n2 CI_{R,h,t} \quad (16)$$

Where,  $CI_{H,h,t} = CI_{AC,h,t} + CI_{WH,h,t}$ , CI of thermal appliances.

$CI_{R,h,t}$  = CI of remaining active appliances

Where  $n1$  and  $n2$  are weight factors and  $(n1 + n2) = 1$ .

The CI of AC is calculated as in equation (17),

$$CI_{AC,h,t} = \left| \frac{2T_{AC,h,t} - T_{AC,Lo,h} - T_{AC,Hi,h}}{T_{AC,Hi,h} - T_{AC,Lo,h}} \right| \quad (17)$$

The CI of EWH is illustrated as in equation (18),

$$CI_{WH,h,t} = \left| \frac{2T_{WH,h,t} - T_{WH,Lo,h} - T_{WH,Hi,h}}{T_{WH,Hi,h} - T_{WH,Lo,h}} \right| \quad (18)$$

The CI of all the remaining active controllable appliances is as shown below,

$$CI_{R,h,t} = \left| \frac{TNA_h - TAA_h}{TNA_h - TRA_h} \right| \quad (19)$$

Where,  $TNA_h$  = Number of active appliances in resident 'h'.

$TAA_h$  = Number of allowed appliances to ON in resident 'h'.  $TRA_h$  = Total number of requested appliances to ON in resident 'h'.

#### IV. PROBLEM FORMULATION

The relation between reward rate and CI is as follows:

$$RWR_{h,t} = \begin{cases} Rw1, & \text{if } CI \leq 1 \\ R w2, & \text{if } CI > 1 \text{ and } compromise = 1 \\ R w3, & \text{if } CI > 1 \text{ and } compromise = 0 \end{cases} \quad (20)$$

From equation (20), if the CI value is less than or equal to '1' (i.e. Customer Satisfaction), reward 'Rw1' will be given. If the CI is more than '1' (i.e. Customer Dissatisfaction) and is willing to reduce demand, consumers are rewarded with 'Rw2'. If the consumer is not ready to compromise, the demand and the reward rate is 'Rw3' is provided to encourage them for active participation. It is a rate which is given during emergency but it is a rare case. These reward structures are used to attract the non-participant and not compromised customers to involve in the DR program.

$$RWR_{h,t} = R w_1 b_{h,t} + R w_2 (1 - b_{h,t}) comp_h + R w_3 (1 - b_{h,t}) (1 - comp_h) \quad (21)$$

Where,  $b_{h,t}$  is a binary variable and  $comp_h$  is a compromise for the house 'h'.

By considering the following objective function and constraints, the optimization problem is formulated. The primary objective of this proposed approach is to reduce the utility reward costs and CI. While reducing CI, the satisfaction level of the user's comfort will increase. The objective function is as shown in equation (22).

$$C = \min \sum_{h=1}^H RWR_{h,t} + k \cdot \sum_{h=1}^H CI_h \quad (22)$$

With the following constraints:

$$\sum_{h=1}^H CDR_{h,t} \geq PDRR_t \quad (23)$$

$$-B b_{h,t} < (DRL - P_{a,h,t}) \leq B(1 - b_{h,t}) \quad (24)$$

$$P_{a,h,t} = \sum_{n=1}^N P_{h,n,t} \cdot S_{h,n,t} \quad (25)$$

$$RW_{h,t} = (P_{T,h,t} - P_{a,h,t}) \cdot RWR_{h,t} \quad (26)$$

Where,  $k$  is the CI weight factor,  $CDR_{h,t}$  is the Consumer Demand Reduction of the resident 'h' (kW) at time 't' and  $B$  is a positive constant,  $S_{h,n,t}$  is the status of the appliance (1 or 0). Equation (23) represents the total demand reduction of all the residents should be higher than or equal to PDRR. It is not possible to achieve a reduction in demand beyond the resident's total power consumption and is expressed in equation (24). The overall power consumption is corresponding to the sum of power consumption of active appliances in all the residents, as shown in equation (25). The reward rate for each resident is calculated using equation (26). Hence, this problem is formulated as mixed-integer linear programming (MILP).

#### V. OPTIMAL SOLUTION APPROACH

##### A. Conventional Approach

In [2], the concept of CI is used for measuring consumer comfort level by considering the thermal appliances like AC and EWH alone. BONMIN solver in GAMS is used to verify the performance of the approach. Two DRR schemes such as 30% and 60% are performed and validation is done only for AC.

##### B. Proposed Approach

The above problem statement is the mixed of continuous variable  $P_{a,h,t}$ , discrete variable (reward rate) and binary variable  $b_{h,t}$ . Therefore it can be modeled as MILP.

By substituting equation (21) and (26) in (22), the new objective function is obtained as follows in equation (27).

$$C_{New} = \min_{P_{a,h,t}} \left\{ \begin{aligned} & P_{T,h,t} [(Rw_1 - Rw_3) b_{h,t} + Rw_3] - \\ & B_{new,i,t} (Rw_1 - Rw_3) - P_{a,h,t} Rw_3 + k CI + \\ & (P_{T,h,t} - P_{a,h,t}) [(Rw_2 - Rw_3) (1 - b_{h,t}) comp_h] \end{aligned} \right\} \quad (27)$$

Subject to the following constraint (28) along with the constraints (23), (24), (25) and (26).

$$0 \leq B_{new,i,t} \leq P_{T,h,t} b_{h,t} \quad (28)$$

Where,  $B_{new,i,t} = P_{a,h,t} b_{h,t}$ . In this approach, as  $b_{h,t}$  is considered as binary variable  $B_{new,i,t}$  is zero if  $b_{h,t} = 0$  otherwise  $B_{new,i,t} = P_{a,h,t}$  ( $b_{h,t} = 1$ ).

Then the utility reward cost and average comfort for each house are calculated. This MILP for reducing the utility reward costs and CI, which is guaranteed the finding global optimum operation work has been solved by using CPLEX solver in General Algebraic Modelling Software(GAMS).

### C. Case Studies

The proposed DSM strategy for the residential consumer is performed on 10, and 500 residents and input data are taken from [2]. The test system considers thermal, controllable and other active appliances that are participated during the DR program. The load data of the 10 residents are recognized for this analysis as given in Table II. The total demand of 10 residents is 157.8 kW. Table III shows the total demand reduction of 10 residents for various CDR during peak hours. DR program of each resident is done based on user-preferred load and DRL. DRL is the threshold value for demand reduction and decided by each user. Compromise '1' represents the resident is willing to compromise their demand during the peak hours. Compromise '0' represents the resident is not ready to reduce their demand.

Table IV shows the average percentage of comfort and reward for various PDRR. In this case study Rw1, Rw2, and Rw3 are considered as 20, 40, 60 cents / kW (5 min). The comfort percentage is the measure of the number of times the power consumption of the residents is within the user preferred power range. The time duration taken for implementing each PDRR is 5 minutes. The reward rate and comfort percentage of each PDRR for 10 residents are explained below. In this analysis, six PDRR are discussed like 10%, 20%, 30%, 40%, 50% and 60%.

- PDRR#1 with 16kW/20min (10%): RLA receives a request of 16 kW for the time length of 20 minutes from LSE. Ten residents' total demand is 157.8 kW. Approximately 10% of the load from the total demand is decreased. All residents are within their reduction limit of preferred demand. Therefore, for all residents, the percentage of comfort is 100% and with the Rw1 reward. According to their lower power utilization and a wide range of demand reduction limits, Resident 3 gets more reward. Resident 2 does not earn any reward because of their lower comfortable power range.
- PDRR#2 with 32kW/20min (20%): RLA receives a request of 32 kW from LSE for 20 minutes. Approximately 20% of the load from the total demand is reduced. All residents are within their preferred demand reduction limits. Therefore the percentage of comfort is 100% for all the residents and with the

reward of Rw1. Resident 7 procured more rewards due to their less power utilization and a wide range of demand reduction limits.

- PDRR#3 with 47kW/20min (30%): RLA receives a request of 47 kW for 20 minutes from LSE. Approximately 30 percent of the load is decreased from total demand. All residents are within their required demand reduction limit. The rate of comfort for all residents is therefore 100%, and the rate of reward is Rw1. Resident 7 received more rewards because of their lower power consumption and a wide range of reduction in demand.
- PDRR#4 with 63kW/20min (40%): RLA receives a request of 63 kW for the time length of 20 minutes from LSE. The total demand reduction of 10 residents is approximately 40% of total demand. All the residents are within their preferred demand reduction limits. Therefore the percentage of comfort is 100% for all the residents and with the reward of Rw1. Resident 5 has earned more reward because of their reduced power use and a wide range of reduced demand.
- PDRR#5 with 75kW/20min (50%): RLA receives a request of 79 kW for the time length of 20 minutes from LSE. Approximately 50% of the load is reduced from total demand. All the residents are within their preferred demand reduction limits. Therefore the percentage of comfort is 100% for all the residents excludes and with the reward of Rw1 and Rw2. Resident 7 procured more reward due to their less power utilization and a wide range of demand reduction limits.
- PDRR#6 with 95kW/20min (60%): RLA receives a request for the time length of 20 minutes from LSE of 95 kW. The demand reduction of 60% is achieved from total demand. Some of the residents are within their preferred demand reduction limits. The level of comfort for all residents is therefore 100%, except residents 1, 5, and 9 with a reward rate of Rw1 and Rw2. Resident 1 procured more rewards due to their less power utilization and a wide range of demand reduction limits. The reward cost of utilities will also be increased while the PDRR is increasing. During an emergency, the affected houses can receive reward rate Rw3.

TABLE II. TOTAL LOAD DATA OF THE 10 RESIDENTS

| Resident | Controllable devices (kW) |     |     |     |     |     |     | Uncontrollable Devices (kW) | Total power (kW) |
|----------|---------------------------|-----|-----|-----|-----|-----|-----|-----------------------------|------------------|
|          | AC                        | EWH | CD  | DW  | EV  | CW  | PP  |                             |                  |
| 1        | 1.4                       | 4.0 | 3.4 | 2.9 | 4.0 | 1.3 | 1.2 | 1.1                         | 19.3             |
| 2        | 1.2                       | 3.9 | 3.7 | 2.7 | 0   | 0.9 | 1.3 | 1.3                         | 15.0             |
| 3        | 1.5                       | 3.5 | 3.8 | 3.0 | 3.8 | 1.1 | 1.4 | 1.3                         | 19.4             |
| 4        | 1.6                       | 3.8 | 0   | 2.6 | 0   | 1.2 | 1.5 | 1.1                         | 11.8             |
| 5        | 1.3                       | 3.1 | 3.1 | 0   | 3.6 | 1.0 | 1.1 | 1.4                         | 14.6             |
| 6        | 1.2                       | 3.4 | 3.5 | 2.8 | 0   | 1.3 | 1.2 | 1.2                         | 14.6             |
| 7        | 1.1                       | 3.9 | 3.7 | 0   | 3.8 | 1.2 | 1.3 | 1.5                         | 16.5             |
| 8        | 1.5                       | 3.8 | 0   | 2.9 | 4.0 | 0.9 | 1.4 | 1.7                         | 16.2             |
| 9        | 1.5                       | 4.0 | 3.3 | 2.6 | 0   | 1.1 | 1.5 | 1.1                         | 15.1             |
| 10       | 1.3                       | 3.2 | 3.2 | 0   | 3.6 | 1.2 | 1.6 | 1.2                         | 15.3             |

TABLE III. TOTAL CDR OF 10 RESIDENTS

| Resident | Total power (kW) | DRL (kW) | Compromise (Yes=1, No=0) | CDR (kW) |      |      |      |      |     |
|----------|------------------|----------|--------------------------|----------|------|------|------|------|-----|
|          |                  |          |                          | 10%      | 20%  | 30%  | 40%  | 50%  | 60% |
| 1        | 19.3             | 11.4     | 1                        | 18.3     | 16.2 | 14.1 | 11.4 | 7.9  | 5.6 |
| 2        | 15               | 7.9      | 0                        | 14.5     | 12.4 | 10.7 | 7.9  | 7.9  | 7.9 |
| 3        | 19.4             | 12.8     | 1                        | 16.1     | 15.7 | 14.5 | 12.8 | 8.1  | 7.1 |
| 4        | 11.8             | 7        | 0                        | 11.1     | 9.5  | 8.4  | 7    | 7    | 4.5 |
| 5        | 14.6             | 7.1      | 1                        | 13.5     | 12.8 | 9.7  | 7.1  | 8.6  | 5.1 |
| 6        | 14.6             | 10.8     | 0                        | 12.4     | 11.5 | 11   | 10.8 | 10.8 | 6.5 |
| 7        | 16.5             | 10.6     | 1                        | 14.8     | 12.8 | 10.9 | 10.6 | 4.5  | 7.1 |
| 8        | 16.2             | 9.2      | 0                        | 14.1     | 10.9 | 11.2 | 9.2  | 9.2  | 9.2 |
| 9        | 15.1             | 9.7      | 1                        | 13.7     | 12.1 | 9.9  | 9.7  | 6.4  | 5.4 |
| 10       | 15.3             | 6        | 0                        | 13.5     | 12.3 | 10.1 | 8.5  | 8.5  | 4.7 |

TABLE IV. PERCENTAGE OF COMFORT AND REWARD FOR VARIOUS PDRR

| Resident | 10%              |             | 20%              |             | 30%              |             | 40%              |             | 50%              |             | 60%              |             |
|----------|------------------|-------------|------------------|-------------|------------------|-------------|------------------|-------------|------------------|-------------|------------------|-------------|
|          | Avg. Comfort (%) | Reward (\$) | Avg. Comfort (%) | Reward (\$) | Avg. Comfort (%) | Reward (\$) | Avg. Comfort (%) | Reward (\$) | Avg. Comfort (%) | Reward (\$) | Avg. Comfort (%) | Reward (\$) |
| 1        | 100              | 0.2         | 100              | 0.92        | 100              | 2.07        | 100              | 3.16        | 100              | 9.12        | 75               | 13.97       |
| 2        | 100              | 0           | 100              | 0.78        | 100              | 1.73        | 100              | 2.84        | 100              | 5.68        | 100              | 7.24        |
| 3        | 100              | 0.67        | 100              | 1.11        | 100              | 1.94        | 100              | 2.64        | 100              | 9.04        | 100              | 12.54       |
| 4        | 100              | 0.14        | 100              | 0.69        | 100              | 1.38        | 100              | 1.92        | 100              | 3.84        | 100              | 7.45        |
| 5        | 100              | 0.22        | 100              | 0.54        | 100              | 1.96        | 100              | 4.50        | 100              | 4.80        | 75               | 9.69        |
| 6        | 100              | 0.44        | 100              | 0.93        | 100              | 1.44        | 100              | 1.52        | 100              | 3.04        | 100              | 8.26        |
| 7        | 100              | 0.34        | 100              | 1.71        | 100              | 2.24        | 100              | 2.36        | 75               | 12.24       | 100              | 9.58        |
| 8        | 100              | 0.42        | 100              | 0.99        | 100              | 2.0         | 100              | 2.80        | 100              | 5.60        | 100              | 7.14        |
| 9        | 100              | 0.28        | 100              | 0.90        | 100              | 3.12        | 100              | 2.16        | 100              | 6.96        | 75               | 9.89        |
| 10       | 100              | 0.36        | 100              | 0.89        | 100              | 2.08        | 100              | 2.72        | 100              | 5.44        | 100              | 10.81       |

VI. RESULT AND DISCUSSION

A. 10 Resident System

In the current framework, thermal controlled devices EWH and AC are considered to determine the level of comfort. In this method, all the device status in residents is taken into account. The simulation result shows the Percentage Average Comfort (PAC) and Utility Reward Cost (URC). Fig. 2(a, b) demonstrates the effect of URC when considering different amounts of reduction in demand and time for a 10 resident system for existing and proposed method. Compared to the existing method [2], the simulation shows reduced URC for various PDRR during the time length. Fig. 3(a, b) shows PAC for existing and proposed method. PAC is same in both the methods up to 40% of PDRR for different time length. Some residents may be less comfortable while increasing the PDRR above 40%. The affected residents may get the reward rate  $R_w2$ . This method provides a higher level of comfort above 40% of PDRR. For the existing and proposed method, the PAC and average URC for 10 resident systems are as shown in Table V.

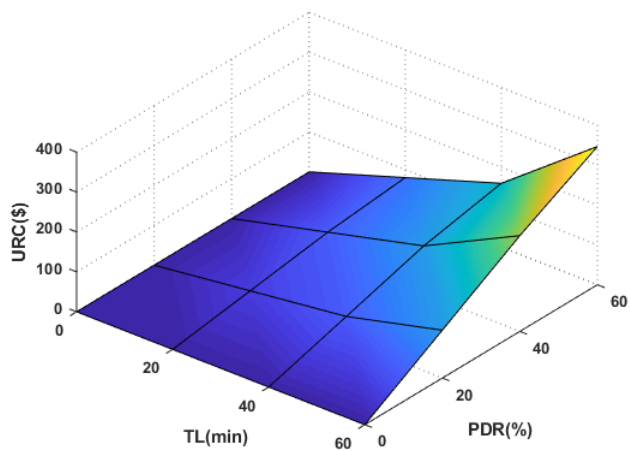
It is observed that the performance for 60% DRL, the average reward cost URC for existing, and these proposed methods are \$346.44 and \$207.82. PAC for existing and the proposed method are 52.5% and 88%.

B. 500 Resident System

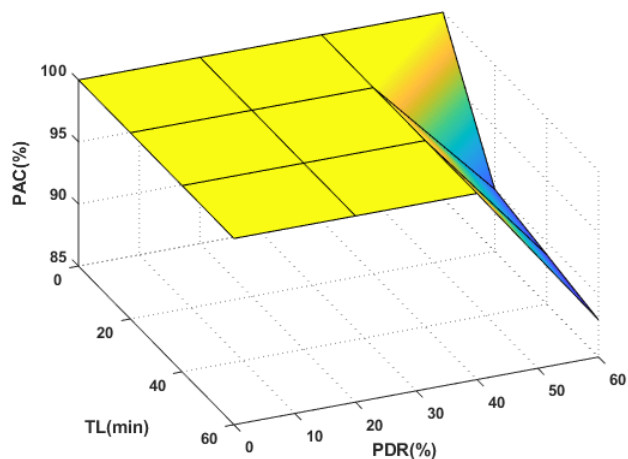
The data of large systems are taken from [2]. In this analysis 500 resident system is considered. Fig. 4(a, b) compares the behaviour of PAC of proposed method along with existing method for different time length and PDRR. While increasing the PDRR for an increasing time length there is a significant reduction of PAC. But when compared to the existing method the PAC is 60% higher. This leads to the reduction of overall utility reward cost for performing the entire PDRR program. Fig. 5(a, b) shows the effect of URC for different percentage of reduction in demand and time length for 500 resident system. For the current and proposed method, the PAC and URC for 500 resident systems are obtained as in Table VI. It is observed that for 80% PDRR, the reward cost URC for existing, and these proposed methods are \$13500 and \$10700. PAC for existing and the proposed method are 18.5% and 78.5%. The simulation result shows that this approach minimizes the utility reward cost significantly and increases the percentage of average comfort compared to the existing method.

TABLE V. RESULT COMPARISON FOR 10 RESIDENTS

| Approach    | Percentage of average Comfort (PAC) (%) | Utility reward cost (URC) (\$) |
|-------------|---|--------------------------------|
| Existing[2] | 52.5                                    | 346.44                         |
| Proposed    | 88                                      | 207.82                         |

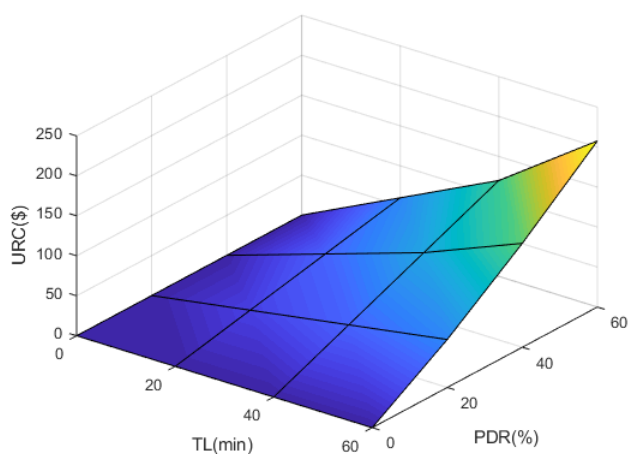


(a) Existing Method.



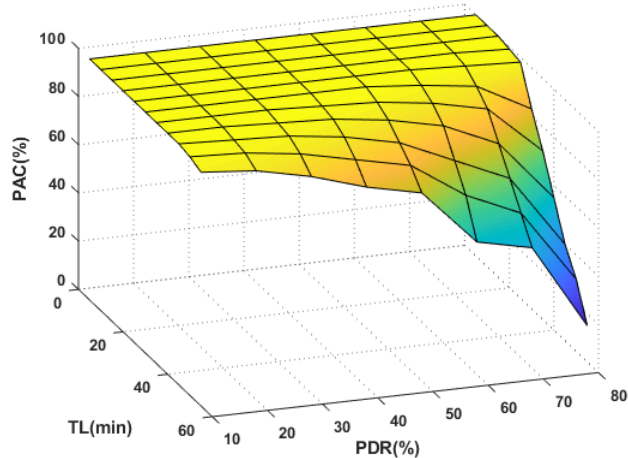
(b) Proposed Reward Method.

Fig. 3. Result of the Percentage of Average Comfortableness for 10 Residents System.

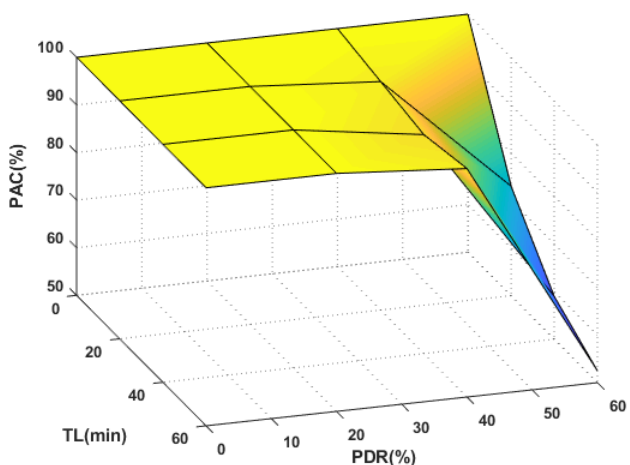


(b) Proposed Reward Method.

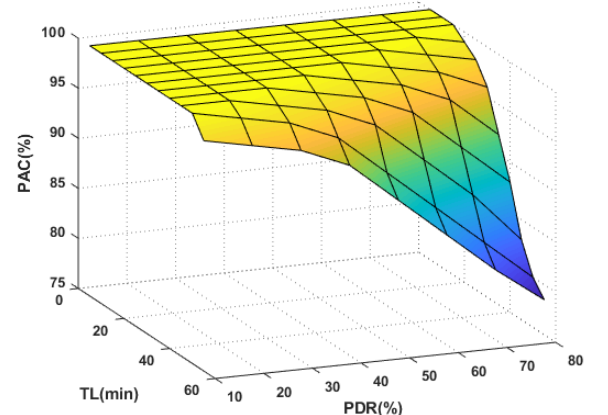
Fig. 2. Result of utility Reward Cost for 10 Residents System.



(a) Existing Method.



(a) Existing Method.



(b) Proposed Reward Method.

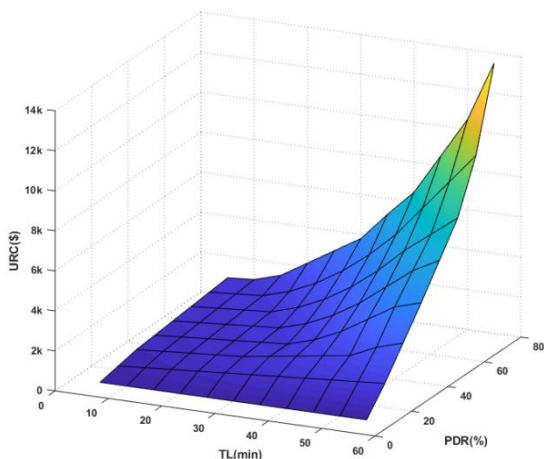
Fig. 4. Result of the Percentage of Average Comfortableness for 500 Resident Systems.

ACKNOWLEDGMENT

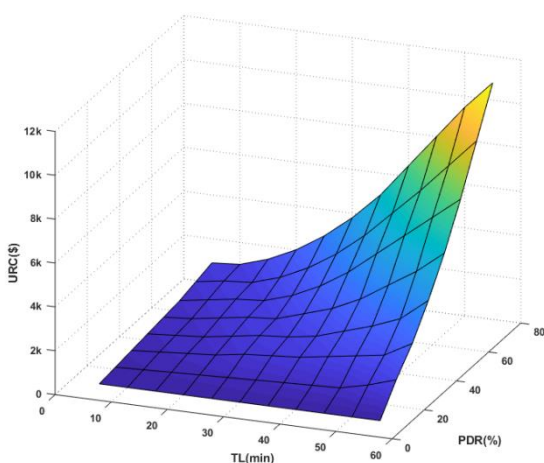
The authors would like to express their thanks to Vellore Institute of Technology for their continuous support.

REFERENCES

- [1] Yi Liu, Liye Xiao, Guodong Yao, and S.Q. Bu “Pricing-Based demand response for a smart home with various types of household appliances considering customer satisfaction” IEEE Access, vol.7,pp. 86463-86472, 2019.
- [2] Q. Hu, F. Li, X. Fang, and L. Bai, “A framework of residential demand aggregation with financial incentives,” IEEE Trans. Smart Grid, vol. 9, No. 1, pp. 497-505, Jan2018.
- [3] E. Shahryari, H. Shayeghi, B. Mohammadi-ivatloo, M. Moradzadeh, “An improved incentive-based demand response program in day-ahead and intraday electricity markets” Energy 155 .pp.205-214,2018.
- [4] Valles M, Bello A, Renesas J, Frias P,,” Probabilistic characterization of electricity consumer responsiveness to economic incentives”. Appl Energy .Feb.2018.
- [5] Ni, Z.; Das, A.A new incentive-based optimization scheme for residential community with financial trade-offs IEEE Access, Oct. 2018, 6, 57802–57813.
- [6] M. R. Sarker, M. A. Ortega-Vazquez, and D. S. Kirschen, “Optimal coordination and scheduling of demand response via monetary incentives,” IEEE Transactions on Smart Grid, vol. 6, no. 3, pp. 1341–1352, May 2015.
- [7] Paudyal. P, Ni. Z, “Smart home energy optimization with incentives compensation from inconvenience for shifting electric appliances”, Electrical Power and Energy Systems 109 (2019) 652–660.
- [8] Mengmeng Yu, Seung Ho Hong “Incentive-based demand response considering hierarchical electricity market: A stagelberg game approach” Applied Energy 203 (2017) 267–279.
- [9] C. Vivekananthan, Y. Mishra, G. Ledwich, and F. Li, “Demand response for residential appliances via customer reward scheme,” IEEE Trans.Smart Grid, vol. 5, no. 2, pp. 809–820, Mar. 2014.
- [10] Priti Paudyal, Prateek Munankarmi, Zhen Ni, Timothy M. Hansen, “A Hierarchical Control Framework with a Novel Bidding Scheme for Residential Community Energy Optimization” IEEE transaction on smart grid,2019.
- [11] Sayyad N, Kazem Z, Behnam Mohammadi-Ivatloo, “Robust bidding and offering strategies of electricity retailer under multi-tariff pricing”, Energy Economics 68 (2017) 359–372.
- [12] S. Moon and J.-W. Lee, “Multi-residential demand response scheduling with multi-class appliances in smart grid,” IEEE Trans. Smart Grid, vol. 9, no. 4, pp. 2518-2528, Jul. 2018.
- [13] Joo, I.Y., Choi, D.H, “Optimal household appliance scheduling considering consumer’s electricity bill target”, IEEE Trans. Consum. Electron., 2017, 63,(1), pp. 19–27.
- [14] Nan, S, Zhou, M, Li, G, “Optimal residential community demand response scheduling in smart grid”, Applied Energy 210 (2018) 1280–1289.
- [15] Sayyad Nojavan, Ramin Nourollahi, Hamed Pashaei-Didani, Kazem Zare, “Uncertainty-based electricity procurement by retailer using robust optimization approach in the presence of demand response exchange”, Electrical Power and Energy Systems 105 (2019) 237–248.
- [16] Shafie-Khah, M., Siano, P, “A stochastic home energy management system considering satisfaction cost and response fatigue”, IEEE Trans. Ind. Inf., 2018, 14, (2), pp. 629–638.
- [17] Z. Zhao, W. C. Lee, Y. Shin, and K.-B. Song, “An optimal power scheduling method for demand response in home energy management system,” IEEE Trans. Smart Grid, vol. 4, no. 3, pp. 1391–1400, Sep. 2013.
- [18] M. Pipattanasomporn, M. Kuzlu, and S. Rahman, “An algorithm for intelligent home energy management and demand response analysis,” IEEE Trans. Smart Grid, vol. 3, no.4, pp. 2166–2173, Dec. 2012.
- [19] Beaudin, M., Zareipour, H, “Home energy management systems: a review of modelling and complexity”, Renew. Sust. Energy Rev., 2015, 45, pp. 318–335.



(a) Existing Method.



(b) Proposed Reward Method.

Fig. 5. Result of Utility Reward Cost for 500 Residents System.

TABLE VI. RESULT COMPARISON FOR 500 RESIDENTS

| Approach    | Percentage of average Comfort (PAC) (%) | Utility reward cost (URC) (\$) |
|-------------|---|--------------------------------|
| Existing[2] | 19.5                                    | 13500                          |
| Proposed    | 78.5                                    | 10700                          |

VII. CONCLUSIONS

An efficient demand response program based on rewards is introduced in this paper. It takes into account all active electrical equipment involved in DR. In this analysis, the CI is the essential factor that defines the resident’s level of comfort. This valid reward-based scheduling method minimizes utility reward cost and increases the PAC. It is identified that the proposed approach maintain the average comfort of consumer while increasing the residents from 10 to 500. Result of case studies inferred that the reward-based demand response program provides a better cost solution to utility and consumers compared to state of art work. In future, the proposed approach should be improved for meeting the realistic constraints that can be evaluated using large scale system with real time data.

[20] F. Ruelens, B. J. Claessens, S. Vandael, B. D. Schutter, R. Babuska, and R. Belmans, "Residential demand response of thermostatically controlled loads using batch reinforcement learning," IEEE Trans. Smart Grid, vol. 8, no. 5, pp. 2149–2159, Sep. 2017.

[21] Shenglin Li, Junjie Yang, Wenzhan Song, and An Chen, "A Real-Time Electricity Scheduling for Residential Home Energy Management" IEEE Internet Of Things Journal, vol. 6, no. 2, Apr. 2019, pp.2602-2611.

[22] Muthuselvi G, Saravanan B, "Energy Consumption Scheduling Using Adaptive Differential Evolution Algorithm in Demand Response Programs", International Journal of intelligent engineering and systems (2019), vol. 12, no.5.

[23] K.Stenner, E.R.Frederiks, E.V.Hobman, and S.Cook, "Willingness to participate in direct load control: The role of consumer distrust", Applied Energy, vol. 189, pp.76–88, Mar. 2017.

[24] Seung-Jun Kim, Georgios B. Giannakis, "An online convex optimization approach to real-time energy pricing for demand response", IEEE Trans. Smart Grid 8, vol.6, pp.2784–2793, 2017.

[25] Yu Wang, Haiyang Lin, Yiling Liu, Qie Sun, Ronald Wennersten, "Management of household electricity consumption under price-based demand response scheme", Journal of Cleaner Production 204, 926-938,2018.

APPENDIX (NOMENCLATURE)

|                  |   |                |  |
|------------------|---|----------------|--|
| H -              | Number of residents.  | $S_{DW,h,t}$ - | ON-OFF status of Dish Washer (DW) in resident 'h' at time t.           |
| N -              | Number of appliances in a resident.   | $P_{Dw}$ -     | Power rating of DW (kW).   |
| $S_{EWH,h,t}$ -  | Current state (ON/ OFF) of Electric Water Heater (EWH) in resident 'h' at time t. | $p_{Dw,h,t}$ - | Power usage of DW in resident 'h' at time t (kW).                      |
| $S_{EWH,h,t-1}$  | Previous state of EWH in resident 'h'   | $T_{Dw,i,t}$ - | Current temperature of DW (°F).  |
| $P_{EWH}$ -      | Rated Power of EWH (kW)   | $T_{Dw,s}$ -   | Setpoint temperature of DW (°F).                                       |
| $p_{EWH,h,t}$ -  | Power usage of EWH at resident 'h' at time t (kW)                                 | $S_{CW,h,t}$ - | ON-OFF status of Cloth Washer (CW) heater in resident 'h' at time 't.' |
| $T_{EWH,h,t}$ -  | Current temperature of EWH (°F)   | $P_{CW}$ -     | Power rating of CW (kW).   |
| $T_{EWH,r}$ -    | Minimum required temperature of EWH (°F)  | $p_{CW,h,t}$ - | Power usage of CW in resident 'h' at time 't' (kW).                    |
| $S_{AC,h,t}$ -   | Current state (ON / OFF) of Air Conditioner (AC) at resident 'h' at time t.       | $T_{CW,h,t}$ - | Current temperature of CW (°F).  |
| $S_{AC,h,t-1}$ - | Previous State of Air Conditioner (AC) at resident 'h'                            | $T_{CW,s}$ -   | Setpoint temperature of CW (°F).                                       |
| $P_{AC}$ -       | Power rating of AC (kW)   | $S_{PP,h,t}$ - | ON-OFF status of Pool Pump (PP) at resident 'h' at time 't'.           |
| $p_{AC,h,t}$ -   | Power usage of AC in resident 'h' at time 't' (kW)                                | $P_{PP}$ -     | Power rating of PP (kW).   |
| $T_{AC,h,t}$ -   | Current temperature of AC in resident 'h' at time 't' (°F).                       | $p_{PP,h,t}$ - | Power usage of PP in resident 'h' at time t (kW).                      |
| $T_{AC,DB}$ -    | Dead band temperature of AC (°F).   | $T_{PP,h,t}$ - | Current temperature of PP (°F) of resident 'h' at time t.              |
| $T_{AC,sp}$ -    | Setpoint temperature of AC (°F).  | $T_{PP,s}$ -   | Setpoint temperature of PP (°F).                                       |
| $S_{CD,h,t}$ -   | Current state (ON / OFF) of Clothes Dryer (CD) in resident 'h' at time t.         | $T_{Lo,h}$ -   | Minimum temperature of the room in resident 'h' (°F).                  |
| $T_{CD,acc}$ -   | Accumulated ON time temperature of CD (°F).                                       | $T_{Hi,h}$ -   | Maximum temperature of the room in resident 'h' (°F).                  |
| $T_{CD,r}$ -     | Required ON time temperature of CD (°F).  | $T_{RM,h}$ -   | Room temperature of resident 'h' (°F).                                 |
| $p_{CD,h,t}$ -   | Power consumption of CD at resident 'h' at time t (kW).                           | $P_{T,h,t}$ -  | Overall power usage of resident 'h' at the time 't'.                   |
| $P_{CD}$ -       | Power rating of CD (kW).  | $P_{h,n,t}$ -  | Power usage of appliance 'n' of resident 'h' at the time 't'.          |
| $P_{EV}$ -       | Power rating of EV (kW).  | $P_{a,h,t}$ -  | Actual Power consumption of resident 'h' at time 't'.                  |
| $S_{EV,h,t}$ -   | ON-OFF status of Electric Vehicle (EV) in resident 'h' at time t.                 | $P_{DRR}$ -    | Preferred Demand Reduction Request                                     |
| $p_{EV,h,t}$ -   | Power charge of EV in resident 'h' at time t (kW).                                |                |  |
| $RWR_{h,t}$ -    | Reward rate of resident 'h' at time t.  |                |  |
| $SOC_{h,t}$ -    | Battery charging state of resident 'h' during time period t (%).                  |                |  |
| $SOC_{max}$ -    | EV-Maximum charging rate of the battery (%).                                      |                |  |

# SQ-Framework for Improving Sustainability and Quality into Software Product and Process

Kamal Uddin Sarker<sup>1</sup>

Aziz Bin Deraman<sup>2</sup>

Faculty of Ocean Engineering  
Technology and Informatics  
University Malaysia Terengganu  
Terengganu, Malaysia

Raza Hasan<sup>3</sup>

Department of Information  
Technology, Malaysian University of  
Science and Technology  
Selangor, Selangor  
Malaysia

Ali Abbas<sup>4</sup>

Department of Computing  
Middle East College  
Knowledge Oasis Muscat  
Oman

**Abstract**—Sustainability is one of the most important quality factors and it integrates some other quality factors in the product too. Moreover, it makes an effective workflow and improves user satisfaction. A manager can meet the target by controlling a project but sustainability is more versatile. Quality factors are the measuring criteria of a product while sustainability drives to make the quality product, efficient project, and successful organization so it is a package of strategy, tasks, processes, technologies, and stakeholders. It is observed that there is a lacking of sustainability practice in software engineering like other engineering communities. There are many software developing models that exist with limited scope in quality control for sustainability. Given the aforementioned viewpoint, this research proposes a new software project management framework, “SQ-Framework”. Its hybrid structure consists of the features of methods, quality models, and sustainability. The execution guidance of “SQ-Framework” is provided according to “Karlskrona manifesto”. A manager can use the framework to improve the management process of a project, a developer can integrate quality factors with sustainability into the product, an executive could be motivated to integrate quality and sustainability strategy in the organization, and the users get inspiration to practice sustainability.

**Keywords**—Software project management; sustainable project; sustainable product; sustainable and quality model; system development methodology

## I. INTRODUCTION

A project is a temporary endeavor that generates a product or service within a fixed budget and time. Project management is a systematic approach that has a set of tasks, processes, guidelines, technologies, and stakeholders to meet the goal of the project on time, in budget with better user satisfaction. User satisfaction is measuring based on the quality criteria of the projects and sustainability is one of them. The evolving concept “quality” is varied according to the application domain and quite tricky to define. A product is known as a qualified product when it has distinct features than a similar type of other products. It aggregates multiple measuring attributes of this domain. It is a misleading and risky word and multi-attributes measuring system required to qualify a software system [1]. If without domain knowledge no one can define a terminology; without proper definition, you can not measure the product's quality, and without knowing measures can not do quality control of a system [2]. The software functionalities

are expressed in quantity ( e.g. source line of code or functional points) that specify the size of software and non-functional attributes are qualitative. Software quality measurement includes both in measures (completeness of the system and user satisfaction factors) and project management integrated quality attributes to a software product. Quality management is a process that aims to make sure the existence of quality factors in a soft product. Quality measures of the quality models and standardization companies guidelines help to gain quality products. The initial measurement of a system depends on the expected quality standard of the customer or user. Moreover, it also meets regulatory quality requirements from the developing teams. Quality management is not only limited to testing activities before released to the market but also maintains critical evaluation processes in different phases [3]. It also guides to develop a quality management culture in the enterprise [3]. Quality management consists of quality planning: method of measuring quality goals, quality control: defect identification and correction, and quality assurance: a set of actions in process management to ensure quality.

Sustainability becomes important in global concern due to the huge consumption of energy for the industrial revolution of the 20th century and in 1930s economists developed sustainable models for non-renewable resource management. Day by day its' value is adding in our life due to technological advancement and smarter lifestyle. In the 21st century, its importance is extending in personal, social, and corporate life. Mitchell Grant simply defined “meet present needs without compromising future needs” [4]. Sustainability ensures the wise-utilization of environmental, social, and economic resources that could offer the same for an upcoming generation [5]. Social sustainability promotes wellbeing by developing processes and structures in society. Sustainable business culture ensures human rights, fairness, diversity, and wellbeing practice in the enterprise. Reducing wastage of time, effort, and money of an organization is called economic sustainability; moreover, it suggests implementing technological applications to make the highest productivity. Environment-friendly farming and foresting practice by using sustainable energy and technology is the practice of environmental sustainability. Innovation in applications and software engineering for financial, social, and environmental sustainability is a common practice. But today researchers focus on integrating green soft features in the system and sustainability practice in the project.



Sustainability is a special quality factor in an application to reduce energy consumption by stopping the unnecessary processing cycles, implementing power-saving mode, applying efficient data structure and algorithm, and using green technology. Information technology (IT) engineers can bring Innovation in software with sustainability features and play the golden role in reducing carbon emissions. Business transformation with green information and communication (ICT) enhances economical sustainability know as ICT for sustainability. The ethical practice of ICT can improve social harmony and cultural exchanges for peace and happiness too. This article shows more importance of sustainability because it integrates other quality factors like re-usability, efficiency, and cost-effectiveness. The proposed framework illustrated a way of integrating sustainability with other quality factors without conflict; as well as the way of practice and keep up the quality of product and project. it brings innovation in cost-cutting, enhancing competitive advantages, and adding the value of software project [6].

“Introduction” consists of the working area of this research and the “literature review” presents the importance of research on sustainability and quality control practice in software industries. It also includes a comparison study on practicing quality models of the software product. Section-3 illustrates a framework for quality control and sustainability practice in software farm that is a hybrid framework of system development methodology and quality model. The aim of this framework is enhancing sustainability practice into the product, and project besides keeping quality control; so it is named "SQ-Framework" to The way of implementation for the proposed model is expanded and elaborated in the followed section immediately. Section five consists of a comparative analysis of the framework with system development approaches and quality models. Section 6 carries the guidelines for using the framework according to the “manifesto”. The last section makes the conclusion of the study with recommendations on future studies.

## II. STATE OF THE LITERATURE

### A. Role of Technology in Project Management

Advancement of technology changes the processes and methods of project management. Project management institute (PMI) has been publishing regular bulletins on project management opportunities, challenges, and technologies. The latest analysis showed higher importance in technical, leadership, business, and digital skill for measuring talent in project management [7]. It noticed that a single business, project, or big idea is not enough in this wrap-speed word to keep an organization at the top constantly. Moreover, your brilliant strategy or amazing product idea could not be a success for supply chain disruption or new technology, and projects can fail fast. Disrupting technology like AI and machine learning can be run by only smart people. The rapid technological change increases challenges into the process of conversion idea to reality. Moreover, it suggests three tenets for resolving the aforementioned challenges: adaptability or agility in process, regular training, and automatic design skills. Project management wants to turn their idea to reality but there is no super-secret formula to make project success; so change

management approach, design thinking approach, hybrid management approach, and agile approach become more popular [8].

### B. Reasons for Project Failure

Researchers, academicians, professionals, standardization organizations, and certification vendors of this domain are working to improve IT project success rate but till now about 45% project is challenging, 36% success, and 19% fail (2018) [9]. Recently, the success rate is increasing by benefits realization management (BRM), which is a powerful project handling approach and it aligns projects, programs, and portfolios to an organization’s overarching strategy [10]. A project does not fail for a single task, person, or process; but there is a cause-effect relationship among task, person, and process. A task is defined and executed by a person or team and the task is executed according to the predefined process so the outcome of the task could not be measured by individual parameters. A project may fail if there is a lack in i) the project and organizational strategy, ii) established accountable result measuring, iii) unambiguous checkpoint or consistent process, iv) consistent methodology for planning and executing, v) stakeholders’ involvement in requirement elicitation and change, vi) utilization maximum effort and vii) effectively use of tools and technology [11]. A quality model or framework helps to ensure quality in all aspects with factors: completeness, accuracy, efficiency, security and reliability, sustainability, usability-accessibility, portability, maintainability, etc.

### C. Quality Models in Project Management

Factor Criteria Model (FCM) is considered the first quality model in software engineering [11] that is developed in 1977 by Air Navy [12]. FCM consists of 11 quality factors that are mapped into three major phases (operation, revision, transition) of the software development and each factor is mapped with multiple criteria of 23. FCM is also known as McCall’s software quality model. After one year (1978), Boehm’s model [13] is developed with a hierarchy structure with 7 top-level quality factors and 15 bottom-level quality factors, one higher factor is linked to 2 or more lower-level factors. It’s clustering consist of three major areas: portability, maintainability, and utility. International Organization for Standardization (ISO) update its generalization model ISO 9000 by ISO 9126 in 1991 [15], but full adaptation was completed in 1992 for software quality measurement [14]. In 1995, Dromey’s quality model is proposed that distinguished a software to multiple product-properties and recommend for adding quality attributes to each product-property list. It is standing on three principles i) setting high-level quality attributes ii) identification product properties that affect quality, iii) linking on product properties and quality attributes [16]. In 2001, ISO restructured the quality view with updated version ISO 9126 - 1:2001 [17]. ISO sets guidelines for measuring software characteristics and international standard measurement of software quality into four sub-domains: “ISO/IEC 9126-1 (ISO/IEC, 2001a)”, “ISO/IEC 9126-2 (ISO/IEC, 2003a)”, “ISO/IEC 9126-3 (ISO/IEC, 2003b)”, and “ISO/IEC 9126-4 (ISO/IEC, 2001b)” respectively for "define and update the model", "define attributes of external measures", "define attributes of internal measures", and "define

the quality on uses". It refined the six main quality measures to 3-5 sub-lists without overlapping. FURPS is an acronym that stands for functionality, usability, reliability, performance, and supportability and each of these has a set of quality attributes presented by a quality model by Robert Gready [18]. Mobile devices were not like today in the developing periods of these models and sustainability was not a serious concern of this domain; moreover, these models show important on the product not process (see details in Appendix-1).

IBM rational software extended this by the name FUPRS+ (2000) with integrating requirements on design, implementation, and interfaces [19]. Software Assurance Technology Centre (SATC), developed a model SATC's quality model in 1996 to support NASA that assisting manager for cost-minimizing and identifying testing quality within four goals: requirements quality, product quality, implementation Effectively, testing effectively [20]. Dromey's model is extended with a hierarchical structure for explicitly specify object-oriented design is known as Bansiya's QMOOD Model and it focuses on the identification of qualified design components, patterns, characteristics, and matrices [21]. Kazman et al integrate the quality factors into the software life style architecture [22] and our study motivated that logically quality factors should work together. Capability Maturity Model Integration (CMMI) gives priority to the organization level too besides project-oriented to ensure the quality of a system. It consists of five maturity levels for integrating and reviewing quality aspects of a system. It is also aligned to ISO 9001 standards; moreover, it promotes the Software Engineering Institute (SEI) of the USA. While ISO 9001 performs quality actions in software development and maintaining stages, CMMI's framework focuses on the continuous improvement of a software process with explicit information [23]. ISO is the world standard organization and CMMI is developed by SEI at Carnegie Mellon University in Pittsburgh and the main difference is "ISO is an audit standard" and "CMMI is a process model" [24]. ISO, CMMI, and IEEE are pioneer organizations in software quality standardization besides national standardization institutions of technological advancement countries. ISO certifies the software firms according to ISO9001 standards, and IEEE has a computer society for certification, ISO and API consulting for software developing companies certification, and CMMI has certification for every CMMI maturity level. This research encourages enhancing business and projects' ability and quality by following standard models, methods, and guidelines without recommending to get the certification.

#### D. Methodologies for Software Project Management

A methodology defines and mentions the work process and management flow every single part of the process. Sometimes, it assumes that more specification kills project execution time and increases project time but imagine "If any error appears after few days that need to recover from the foundation part"; what dangerous the situation is? And the consequence is losing the trust of customers. A short-cut is an instant success but not good for long-term goals. Quality assurance ensures a perfect balance among technology, process, and people to produce a quality product or service, and methodology creates a combination of these three ingredients. The wrong choice of

process methodology is a major risk that can appear during software development [25]. Method, process, model, and framework are upgrading continuously. Scrum is a popular methodology in agile families but the scrum team often overlooks the quality assurance activities due to the tight schedule or early delivery thought agile is not anti-methodology or against the quality practice [26]. The author [26] also demands quality assurance to work with the scrum team to clarify the goal, responsibilities, way handling issues, setting up monitoring and controlling strategies, and finally stay on track to achieve the goals. So, it is an additional activity that going to resolving by the proposed framework.

PRINCE2 (Project in control environment) is a project management methodology developed by the UK government and widely recognized in government and private sectors. It sets roles for the manager, customer, and supplier that uses the "PRINCE board" to accept inputs from users, suppliers, and experts. Series of sequential activities are recorded in PRINCE2 quality practice register: quality identifier, product identifier, product titles, method, roles and responsibilities, planned date of quality check, planned date of sign-off of quality check, the actual date of quality check, the actual date of sign-off of quality check, result, and quality records [27]. It is generic and adaptable to any project that is embraced with ISO 9000 standards, but it is not exact for software project management [28]. The waterfall model is the first process model and commonly used. It has sequential phases and the previous phase provides feedback to the subsequent phase. Quality assurance is predefined at the project initiation stage and practiced in all phases. QATestLab is applying the Black Box testing approach in the waterfall model for quality assurance [29]. It shows the value for quality plans, standards documentation, reviews, verification, demonstration, and quality assurance implementation to every phase. The spiral model has the most influence on risk management and it is right for handing projects with complex functional dependencies [30].

#### E. ITC Product and Sustainability

The information factories do not spew out carbon smoke but they are not bereft of environmental impact for explosion demand of energy. Already, a demand of 200 terawatt-hours (TWh) electricity (2018) is recorded for data centers every year that is more than the energy consumption of some countries like Iran [31]. But the demand will reach 1000 TWh according to the best case practice by 2030 [32]. The same research showed that 2500 TWh would be the least demand by 2030 for ICT productions, networking, data centers, and consumer devices. So, sustainability becomes a series matter of fact in Information and Communication Technology (ICT) similarly engineering and environment studies [33]. In addition, according to the "Ericsson Energy and Carbon Report 2015": mobile subscribers will be 9.5 billion where 55% of mobile data will come from video streaming data by the end of 2020 [34]. Less than 2% of greenhouse gas (GHG) emissions happen for ICT and it could be an enabler to reduce 98% of GHG that emission is not related to ICT [35]. It is also noticeable that 8% energy of the European union's (2015) is consumed by ICT services and subscribers' devices [36]. Definitely, the number of data centers, ICT infrastructures, networking area,

subscribers, the volume of information, have weight data (video), sensor data for Internet of things (IoT), and so on inventions will increase the demand for electricity in ICT sectors. A software engineer can carry out sustainability in a system by choosing comparatively less power consuming technology, applying the effective algorithm and data structure, integrating power-saving features, and reducing unnecessary workflow. On the other side, the new invention of technology “blockchain” consumes 12.76 Twh per year, becomes a new concern of sustainability [37]. However, ecological practice in electronic wastage management by proper destroying, reusing, and recycling is another dimension of sustainability.

#### F. Sustainability and Quality in Software Product

A software project has three major dimensions for improving quality with sustainability practice (Fig. 1).

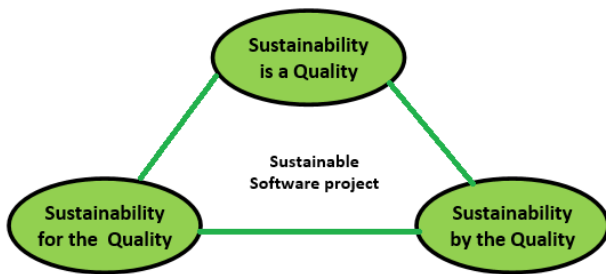


Fig. 1. Sustainability-Quality Relationship.

1) *Sustainability is a quality (sustainability practice in the product)*: According to David Bicknell, attention needs all phases from requirement analysis to deployment and operation phase of a system developing process cycle to create an energy-efficient application [38]. An application could be a sustainable system when it is featured with the auto switch to power saving mode when not need to display anything; auto shutdown actions for dormant applications; turn off wifi, Bluetooth, and global positioning features in a stable place; compression data transmission format; reduce screen brightness and screen timeout; sleep mode activation; reduce data transition steps; increasing resource sharing features; reduce log files; reducing unnecessary historical data; optimal power-consuming programming language, and omitting frivolous features, etc. it becomes more important for mobile systems such as robots, networking nodes, sensors, laptop, and mobiles. Technique and technology both are equally important to integrate green features. The project team should select the comparatively low power consuming data structure (green data structure) [39] and optimal power-consuming algorithm such as Low power Konsumption Algorithm (LOKA) is implemented to end-device of the ZigBee protocol [40].

2) *Sustainability by the quality (sustainability practice in the project)*: Most of the researchers noticed three dimensions for sustainability practice called the “triple bottom line approach” [41]. It is not common that most of them considered three factors of green project management, environmental sustainability practice, and financial sustainability practice. Project strategies and perspective could make suitable for

sustainability practice and it differs from project to project [42]. Some projects focus on social factors [44], while someone considered for financial [43], and others practice by environmental sustainability [43]; so sustainability practice is influenced by the subject of projects. In a software project, sustainability practice could be done by minimizing wastage of time, cost, and effort; maximize re-usability of existing assets: hardware, software, design, and code; and efficient resource allocation. Sustainability practice is not an individual action for a certain task or process of a project so linked with all phases, tasks, and stakeholders and it could be implemented by various methods [45]. This is a strategic level decision of an organization that is executing the project and projects’ sustainability practice should be aligned to the organization. A manager could achieve the target by controlling a project but sustainability is more versatile than just control [46]. Sustainability not only an act but also practice in life and all stakeholders’ participation is important [47].

3) *Sustainability for quality (sustainability by the product)*: An application in a wellbeing tool for humanity and a software project should consider the security concerns inside the system. A product is developed today but it could be backdated after two years for technological advancement that is not financially sustainable. If there is any lacking of protection the system could be abused by culprits that might make a problem for someone. So good software could contribute to sustainable social development by saving expenditure and time, ensuring privacy and security, and making social awareness. An application has an effect on the users’ mental and physical health, an embedded system has radiation effects, and kids’ game might be a cause for addiction (isolated from social gathering [48]). A smart device could broadcast radiation of non-ionizing radio-frequency that is hazardous for health [49].

### III. S-Q FRAMEWORK

1) *Stakeholders*: This model distributes the stakeholders in two major groups called internal (hired employee) who mainly responsible to execute the project and external stakeholders (not working in the developing company) who will support (client, user, sponsor, government agencies) to execute properly. Two types of internal stakeholders are proposed in the framework due to enjoy outsourcing facility or economical employees from the corner of the world (virtual). The employees who are working in the office are remarked as physical mode internal stakeholders. The external stakeholders are not specified due to the enjoy flexibility of mode according to the applications and customers.

2) *Quality factors (QFs)*: There are plenty of quality factors and these are not specified in the model to keep adjustment facility according to the requirements and nature of an application. But the Appendix-1 shows the collection QFs from well know quality models with sustainability, re-usability. While a mobile application is highly biased by power-saving features, a computer user ignores the same features on the desktop. This model proposed to develop a guideline to

integrate the elements from the universal set of quality factors (QFs) (universal set is the superset of all models and sustainability factors). All quality factors are not equally essential to implement for all software except accuracy, security, and efficiency. The importance of a factor depends on the nature of applications such as sustainability is one of the most important for mobile applications. Usability and accessibility are essential for online web applications for mass people.

3) *Green environment*: Sustainability practice is not a single task or unique action in an organization or even in an application. It is part and parcel of personal and corporate practice in a business organization or industry. A software engineer could practice in office management and add to soft products. The green background of the framework (Fig. 3) shows a sign of sustainability and it bears the importance of practice in a project management workplace. An organization should encourage to practice sustainable technology (hardware and application) and a carbon-reducing office environment. Re-usability of historical soft documents (design, code), using tools (information system), and implement sustainability features in the software projects. Moreover, the sustainability motto of the organization can motivate external stakeholders to accept green soft products.

4) *Phase replace by task*: The project is divided into a series of tasks (task1, task2, ....., task<sub>n</sub>) and these are closely adjacent because there is not significant transition between two tasks. It also stands for the flexibility of tasks so that the manager can adjust among tasks, teams, and resources of adjacencies. The size of the task is increasing means it is carrying the historical information and highly dependent on the earlier actions, and if anything is needed it can go back to get or update. It also gives the flexibility in the task for a minor period (margin) of adjustment. Tasks are not mentioned (like traditional phases) for keeping logical adjustment facility according to the application, resource allocation, and risk management. Though task<sub>n+1</sub> is bigger than task<sub>n</sub> in the diagram, the workload of task<sub>n+1</sub> could be less than task<sub>n</sub> because it carries the weight and dependency of previous stages.

5) *Sub-team*: A team could be developed from internal and external stakeholders who can work in the physical and virtual environment. Team with the required resource is tight to each other carries information that there is no transition gap between two adjacent teams. A team can get feedback from its neighbor and a team resource could be a common element of adjacent two teams. This scope has enhanced interaction and communication between two groups and mutual understanding will improve.

6) *Management*: A two-layer management body exists in the SQ-Framework to separate strategic activities and regular activities. A project manager is a key person to execute a project who has to report to the executive level. He/she has to work along to the mission and vision of the organization and business goals of the company too. Managers cooperate and execute organization level strategy of standardization to a project that helps to effectively handle his/her regular work such as manage, organize, coordinate, and keep control of it. A manager works in shell management who directly deals with internal and external stakeholders.

7) *Quality control spin*: The executive body is responsible to check assigned tasks of each team and the manager is the vital person of a project. Instance decisions, communication with stakeholders, risk management, change control, task review and approval, and quality assurance are the main activities. Moreover, it supports the core management for standard documentation development, practicing, assessing, and update regularly. Policy procedure and guidelines development for building a standardization record-keeping system and practice depends on core management. Other stakeholders will provide required input and feedback with accurate information. Information system (tool) and control language (e.g. descriptive logic) are proposed for record management and standardization respectively. Predicate logic, descriptive logic, or control language can reduce the ambiguity of an information system and it is familiar to the technical person of a software project.

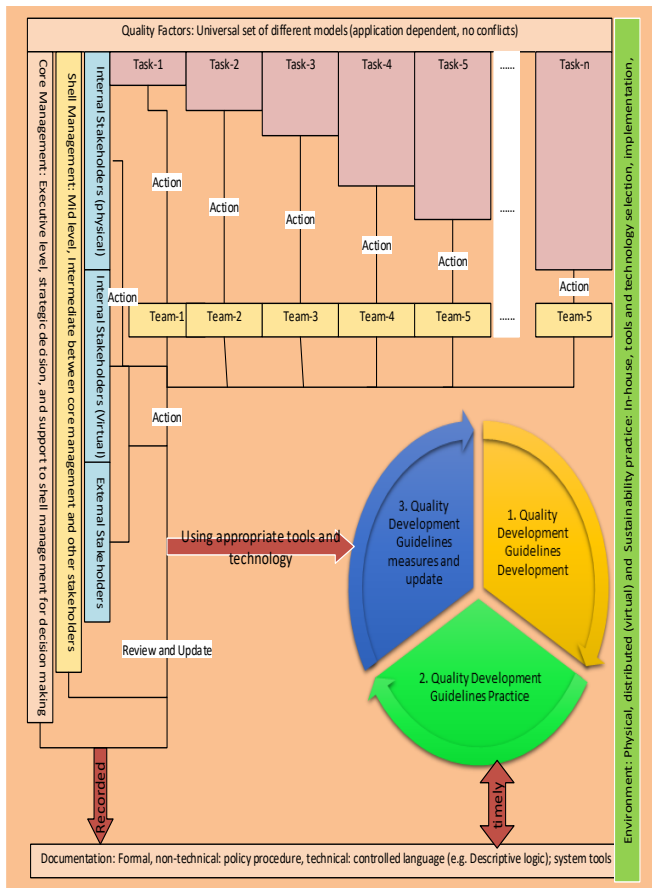


Fig. 2. S-Q Framework.

#### IV. FRAMEWORK EXECUTION MODEL

A 4-phase execution process model (Fig. 3) shows the way of SQ-Framework implementation by four sequential steps (left to right). These phases are connected with a logical sequence of actions but within a phase, the project team can execute parallel operations. The quality of a product is measuring in a particular stage but quality should be ensured in every stage. So quality practice will start from the organizations' strategy and the outcome will 100% fruitful only when the product is utilizing properly.

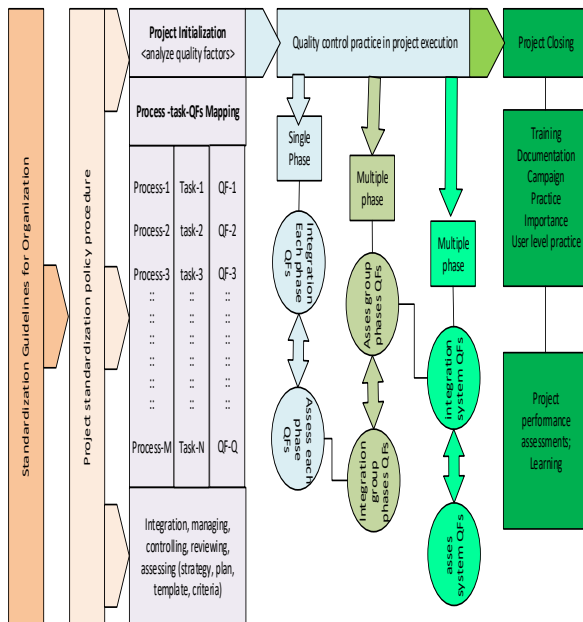


Fig. 3. Framework Execution Process.

**Phase-1 (Definition):** The executive body of an organization develops the strategy for the organization and is executed by subordinates. A project manager launches and executes the project according to the strategy of the organization. The manager can willingly accept and practice the sustainability and or quality guidelines of standardization but it must be practiced when the organization has influence. This phase asks to adapt and practice the standardization guidelines but not recommend for certification because any can practice without certification.

**Phase-2 (initiation):** This is the project initiation stage that analyzes the project scope, feasibility study, planning (hardware-software resources, time, cost, human resources, communication, acceptance, risk. etc.) for a traditional practice.

But, this phase is expected to do more with the strategy of integrating quality factors (QFs) and sustainability factors; specification of the way of implementation, measuring and assessing. Moreover, it suggests making the relationship with the process, task, and stakeholders of the project.

**Phase-3 (execution):** This is the quality implementation action into the product to each task and it would be finalized after a successful quality test. A group of qualified independent tasks could be integrated for the system and it has complete a successful integration test. This process will continue until qualified for the system test.

**Phase-4 (closing):** Phase-1 to 3 already enjoyed sustainability practice to the organization, for the project and integrated into the system; now this is the crucial phase for involving customers into the system and getting long-term advantages. It is also time to assess the performance of the system, learning from the project, and archiving the information for enjoying re-usability.

#### V. FRAMEWORK EXECUTION GUIDELINES

##### A. Karlskrona Manifesto

The Karlskrona Manifesto for Sustainability Design (KMSD) was the outcome of the Third International Workshop on Requirements Engineering for Sustainable Systems (RE4SuSy) at RE'14 in Karlskrona, Sweden [50]. A group of researchers from various disciplines brings the manifesto to leading an interdisciplinary sustainability research platform. The summary of Karlskrona Manifesto [51] is defined in Table I according to the nine principles classification of Oyedeji et al [52]. This specification will guide to create a relationship among stakeholders, task, process, and Karlskrona principle(s).

##### B. Guidelines of utilization SQ-Framework according to Karlskrona Principles (KP) and SQ-Framework Execution Model

SQ-Framework does not mention the phases of software project handling but any project has a set of logical sequences of steps with zero or more iteration that depends on the importance of phases. Documentation, management, quality control, communications are common activities and very much important for project management. Sustainability practice becomes one of them. But for better reflection, a set of the task is mentioned in Table II and KP of Table I is mentioned where it has more influence on each phase of the SQ framework execution model (Fig. 3). In Table II, QFs stands for quality factors and it carries a set of elements that are required according to application form the universal set of QFs (Appendix-1).

TABLE I. SUMMARY OF KARLSKRONA PRINCIPLES (KP)

| Karlskrona Principle (KP)  | Description   |
|--|---|
| Sustainability meets the needs of the future without compromising current demands (KP1). | KP1 is a brief definition of sustainability that could be achieved by introducing Innovation science and technology, management process, and the lifestyle of human beings. It assumes that against the trade-off mindset of the greedy approach but better in long-term business aspects. Aim of the sustainability practice is not only getting advantages now but also the future. |
| Sustainability for long-term and continuous practices (KP2).                             | Advantages of sustainability can enjoy after long-term practice and it could be measured for a long period with multiple scales, indicators, and aspects. This practice should be continuous and mutual participation.  |
| Sustainability is a systematic approach (KP3).   | Sustainability is a systemic set of actions and processes that can execute collectively. Its' need a common background and common platform to design and implement. Its property makes a relationship as an organization to society, society to the nation, nation to global.   |
| Sustainability has multiple dimensions (KP4).  | Sustainability practice could not be possible by a single property or dimension because naturally, it has a relation with finance, environment, energy, and social values. According to a circumstance we have to analyze sustainability outcomes from these aspects.   |
| Sustainability with multiple disciplines (KP5).  | Sustainability practice includes people from multiple societies and backgrounds so working in this area is becoming challenging for human interaction. It also addresses the challenges of multiple disciplines and perspectives.   |
| Sustainability is independent of the purpose of the system (KP6).                        | This principle represents the importance of sustainability and the emphasis on the value of the practice. Even initially is not yet mentioned for a particular task or sustainability is not primary focus but it has to be considered.   |
| Sustainability is a wider context (KP7).   | Sustainability is a part and parcel of a system and its surrounding environment too. It encourages to enhance the scope of dimensions and area of practicing or integrating elements that could be part.  |
| Sustainability is a precondition to system design (KP8).                                 | This is an amendatory condition that shows the significance study of sustainability in system design and development. Sustainability is an essential enabler for system design from different perspective levels and abstractions that could increase participation.  |
| Sustainability action in multiple levels (KP9)   | Look alternative better choice in every level of action and each process of a system to enjoy the most leverage on a system. It encourages comparison analysis on multiple options to accept a better.  |

TABLE II. MODEL UTILIZATION GUIDELINES ACCORDING TO “KARLSKRONA PRINCIPLES (KP)” AND “SQ- FRAMEWORK EXECUTION MODEL”

| Project tasks<br><variable numbers according to Fig.3. The manager has to choose based on the nature and functions of the software project> | Phases of SQ- Framework execution model from Fig. 3   |   |   |  |
|---|---|---|---|--|
|   | Phase-1 <identification, analysis, define and select for the project based on the organization's guidelines and individual attempt> | Phase-2 <map task, process, Sustainability, and quality factors (QFs) and select a better option from alternatives> | Phase-3 <implement according to the phase-1 and phase-2 but do not miss the scope of enhancement for quality> | Phase-4 <practice, learn, update for long term-goal: ensure participation, continuation, practice> |
| Requirement elicitation   | KP1 & QFs <discover, classification, specification, prioritization>   | KP2, KP3, KP6, KP9 & QFs <fixing for each requirement, what and how to implement>                                   | KP2 & KP6 <implementation QFs based on Phase-1 & 2 for all functional & non functional requirements >         | KP2 & KP9 <Verification, correction, learning for next time >                                      |
| Feasibility study   | KP1 & QFs < Analyze and setting for social, economic, technological, and select KP9>  | KP2-5 & DFs <Fixing and keeping alternatives for each factor by KP6 & KP7 >   | KP4-6 <implementation QFs based on Phase-1 & 2 for all aspects >  | KP2-5 <Verification, correction, integration, preparation for next project>                        |
| Planning  | KP3 & QFs < Allocation for time, risk, cost, resource, etc.>  | KP3, KP8, KP9 & QFs <allocate for each task and process to maximum utilization>                                     | KP3-5, & KP9 <execute QFs based on Phase-1 & 2 for all task according to plan >                               | KP3 <Verification, correction, realization and tools selection >                                   |
| Design  | KP8 & QFs < Identify by KP3-KP7 for the whole system too>   | KP4, KP5, KP6, KP9 <set relationship of the task, process, and QFs>   | KP6 & KP8 <implementation QFs based on Phase-1 & 2 for all functional & non functional according to design >  | KP8 <Verification, correction, upgrading KP1 focusing on product and process>                      |
| Development & unit testing  | KP7 & QFs <tools, process, template, time, person, method, specification >  | KP3, KP6, & QFs <method and process of unit test and purposes fixing>   | KP7 <implementation QFs based on Phase-1 & 2 for all task and verify>   | KP7 <Verification, correction, upgrading phase-1 >   |
| System Testing  | KP7, KP8 & QFs < template, time, person, method, specification >  | KP4, KP5, KP9 & QFs <method and process of unit test and purposes fixing>   | KP4, KP5, KP9 <integrate QFs based on Phase-1 & 2 for all requirements and verify according>                  | KP7 <Verification, Completeness, correction, upgrading phase-1 and KP1 >                           |
| Deployment & Maintenance  | KP4, KP5, and QFs <method, training, participation>   | KP4, KP5, KP9 & QFs <regulation, procedure, training, participation>  | KP4, KP5, KP9 <confirm utilization QFs based on Phase-1 & 2 for all by user >                                 | KP1-9 < Realization stakeholders involvement and upgrading KP-1 >                                  |
| Documentation   | KP1-3, and QFs <specify technology, methods, process>   | KP3-5, KP9 & QFs <fixing relationship with tool, format, and structure>   | KP3 <ensure QFs based on Phase-1 & 2 for documentation standard>  | KP3 <Alalysis and upgrade record-keeping, accept tools and techniques>                             |
| Management  | KP3, KP5, & QFs <identification and specification task, person, process, procurement>   | KP1-9 & QFs <fixing policy procedure to manage task, person, process, procurement >                                 | KP3 <track, monitor, control for QFs based on Phase-1 & 2 for declaration>                                    | KP1-9 <analyzised method, process, task and procedure, realize and upgrade>                        |

### C. SQ-Framework Implementation Algorithm Model

#### START

##### Phase-1 (Definition)

- 1.1 Study the system requirements (functional (user & expert ), non-functional (expert), sustainability (expert))
- 1.2 Set the guidelines: Policy-procedure according to organizations' objectives and strategy
- 1.3 Apply KP1: Make clear with explicit specification by definition
- 1.4 Analyze and finalize system requirements (functional, non-functional, sustainability)

For the product, process, and human being:

Apply KP2 to select technology, KP3 for planning, KP4 for the feasibility study, KP5 for managing (task, process, human), KP6 for smart light, desk, data analysis, etc., KP7 for an individual (energy-saving, efficient working process), KP8 for organization policy and design aspect, and KP9 for the selection of better process/tool/technology.

End of 1.4

#### End phase-1

##### Phase-2 (Initiation)

- 1.1 Develop relationship for a task, process, and quality factors( QFs)
  - 2.1.1 Apply KP6 for each task
  - 2.1.2 Apply KP6 for each process
  - 2.1.3 Apply kp6 for each QFs (resolve conflict e.g. faster service & outstanding visualization)

Apply KP7 for generating relationship of 1.1, 1.2, & 1.3

- 1.2 Allocate resource for 1 according to phase-1 (time, effort, money, asset, stakeholders)
- 1.3 Develop quality measuring, assessing criteria, method of assessing, and template for each task, process

#### End of phase-2

##### Phase-3 (Execution)

- 3.1 Apply execution for
  - 3.1.1 Individual task and the respective process by following appropriate  $KP_i (i=1,2,3,\dots,9)$
  - 3.1.2 Perform assessment (QFs) for each
  - 3.1.3 Apply for a group (3.1.1-2) of related tasks and the respective process by following appropriate  $KP_i (i=1,2,3,\dots,9)$
- 3.2 Apply for n number of 3.1 until the finish

#### End of phase-3

##### Phase-4 (Closing)

- 4.1 Justification by  $KP_i (i=1,2,3,\dots,9)$ 

For the process, task, product (Aspect: social, economic, environment, QFs)

Analyze KP1 based on assessment and If a required update
- 4.2 Justification for maintainability, reusability, and training by documentation standard by  $KP_i (i=1,2,3,\dots,9)$  if the required update in KP1
- 4.3 Justification for user participation for product  $KP_i (i=1,2,3,\dots,9)$  and if required update KP1
- 4.4 Analyze the project and accept learning for next practice by  $KP_i (i=1,2,3,\dots,9)$

End of Phase-4

END

## VI. FRAMEWORK ANALYSIS

The SQ-Framework (Fig. 2) is a comprehensive model for sustainable and qualified software development. The framework is not rigid for a specific software domain and it consists of brief guidelines with "Karlskrona manifesto" (Table II). The application of SQ-Framework is presented with a model "SQ- Framework execution model" (Fig. 3) and an algorithm (5.3). Key sustainability attributes are specified by the analysis of the impact of social, economic, and environment to integrate into the framework. There are plenty of software quality models as well as software development life cycles but there is a gap though both are working for software quality control. The framework demolishes the gap and brings in a single platform called "SQ-Framework" that performs project management, quality control, and sustainability practice because individually these are not effective. This model considered cross-platform compatibility in sustainability design for all stakeholders and showed importance on "quality product", "quality process", and "quality management". It aims to make sure the integration of quality factors but not what is the meaning of each factor but all models ignored the importance of sustainability according to the demand of current ages (Appendix-1). But, a reality the quality and sustainability are very closed to each other when a manager works for sustainability automatic quality will improve and vice versa. But this research considered sustainability practice by the system, in the project, and in the product and modeled by Fig. 1.

Features of current applications or systems are not the same as five years back and not will be the same for five years later. Technological advancement changes the working environment and business demands. Common uses of robotics, industrial revolution 4.0 and automation, smart environment, and web video data are increasing rapidly and applications are also biased by those. So, sustainability becomes one of the significant QFs and it also accelerates to integrate other QFs. Industry people could use this framework for systematic software engineering practice. They can modify this according to their demand. It will motivate the executive body for adopting sustainability strategy in the business organization, inspire the client to practice and implement software for sustainable practice, a software development team can develop policy procedures and guidelines to integrate quality and sustainability factors into a system, and manager can improve management process.

## VII. CONCLUSION AND FUTURE WORK

Sustainability and quality highlight each other and jointly both shine the product and developing company. Moreover, sustainability practice in the software industries is new and till now there is a scope for clarification for current sustainability perception. SQ-Framework showed where the scope for implementing sustainability and quality is, how to practice, and who will practice in the software industries. The concrete framework is helpful for engineering practice in soft products and measuring the quality. The milestone of this research is the development of SQ-Framework that reduces the gap of quality models and methodologies. The contribution consists of "identification importance of the sustainability"; "gap of

practicing quality and sustainability”; “the distance between quality model and system development methodologies”; “developing linking between quality and sustainability”; and describes with execution model, Karlskrona manifesto, and algorithm.

SQ-model expects explicit standard documentation, formal practice according to the recommendation of standardization organizations or owns developing standard quality strategy. That is not possible for Ad hoc basis or freelancing development as well as difficult for comparatively small companies who are developing small projects. But if any organization wants to be a standard company then SQ-Framework will guide, though the first time it has to do hard work. Moreover, it would be easier for an experienced manager to execute, modify, and practice.

This paper would guide the researchers to develop new methodologies according to the upcoming technology trend with the dimensionality reduction of complexity, standard documentation ontology, domain-oriented sustainability-quality models: For robotics, embedded system, communication systems, smart infrastructure with (internet of things) IoT, etc.

#### REFERENCES

- [1] Luigi Buglione, Some thoughts on quality models: evolution and perspectives, *Acta IMEKO*, vol. 4, no. 3, article 12, September 2015, identifier: IMEKO-ACTA-04 (2015)-03-1.
- [2] Demarco T., *Controlling Software Projects: Management, Measurement & Estimation*, Yourdon Press, 1982.
- [3] Sommerville, I. (2011). "Chapter 24: Quality Management". *Software Engineering* (9th ed.). Addison-Wesley. pp. 651–680. ISBN 9780137035151.
- [4] Mitchell Grant, (2020).Sustainability - <https://www.investopedia.com/terms/s/sustainability.asp> accessed: 6 May 2020
- [5] Gimenez, C., Sierra, V., & Rodon, J. (2012). Sustainable operations: Their impact on the triple bottom line. *International Journal of Production Economics*, 140(1), 149-159.
- [6] Calero, C. Piattini, M. *Introduction to Green in software engineering*. *Green Softw. Eng.* 2015, 3–27.
- [7] Pulse of the Profession 2020 (2020). *Ahead of the Curve: Forging a Future Focused Culture*. Project management institute. available. <https://www.pmi.org/learning/thought-leadership/pulse/pulse-of-the-profession-2020> visited: 10th May 2020.
- [8] Pulse of the Profession 2019 (2019). *The future of work, leading the way with PMTQ*. Project management institute (PMI). available: <https://www.pmi.org/learning/thought-leadership/pulse/pulse-of-the-profession-2019>, visited: 10 May 2020.
- [9] Clancy, T. *The Standish Group Report*, Retrieved 10 MAY 2020, from <http://www.projectsmart.co.uk/reports.html>, Chaos report, 2018.
- [10] Pulse of the Profession 2017 (2017). 9th global project management survey, <https://www.pmi.org/learning/thought-leadership/pulse/pulse-of-the-profession-2017>. visited: 10th May 2020.
- [11] Discenza, R. & Forman, J. B. (2007). Seven causes of project failure: how to recognize them and how to initiate project recovery. Paper presented at PMI® Global Congress 2007—North America, Atlanta, GA. Newtown Square, PA: Project Management Institute.
- [12] McCall J.A., Richards P.K. & Walters G.F., *Factors in Software Quality*, Voll. I, II, III: Final Tech. Report, RADC-TR-77-369, Rome Air Development Center, Air Force System Command, Griffiss Air Force Base, NY, 1977.
- [13] Boehm B.W., Brown J.R., Kaspar H., Lipow H., MacLeod G.J. & Merritt M., *Characteristics of Software Quality*, Elsevier North-Holland, 1978.
- [14] Suman et al, A Comparative Study of Software Quality Models. (IJCSIT) *International Journal of Computer Science and Information Technologies*, Vol. 5 (4) , 2014, 5634-5638.
- [15] ISO/IEC, IS 9126:1991 - Information Technology - Software product evaluation – Quality characteristics and guidelines for their use.
- [16] Ranbireswar S. Jamwal, Deepshikha Jamwal & Devanand Padha, “Comparative Analysis of Different Software Quality Models”, 3rd National Conference, February 26 – 27, 2009.
- [17] ISO/IEC, IS 9126-1:2001 - Software engineering -- Product quality -- Part 1.
- [18] Grady, Robert B. 1992. *Practical Software Metrics for Project Management and Process Improvement*. Englewood Cliffs (NJ), USA: Prentice-Hall, p 282.
- [19] Kruchten, P. 2000. *The Rational Unified Process: An Introduction*, 2nd Ed. Boston (MA), USA: Addison-Wesley Professional. 320 p.
- [20] Hyatt, L., and L. Rosenberg, 1996 "A Software Quality Model and Metrics for Identifying Project Risks and Assessing Software Quality," NASA SATC, 1996.
- [21] Bansiya, Jagdish, Davis, Carl G., 2002. A hierarchical model for object-oriented design quality assessment. *IEEE Trans. Software Eng.* 28 (1), 4–17.
- [22] Kazman, R.; Nord, R. L.; & Klein, M. A Life-Cycle View of Architecture Analysis and Design Methods (CMU/SEI-2003-TN-026, ADA421679). Pittsburgh, PA: Software Engineering Institute, Carnegie Mellon University, 2003. <http://www.sei.cmu.edu/publications/documents/03.reports/03tn026.html>.
- [23] Mark C. Paulk (1994). A Comparison of ISO 9001 and the Capability Maturity Model for Software. *Software Capability Maturity Model Project*. Technical Report CMU/SEI-94-TR-12 ESC-TR-94-12 [https://resources.sei.cmu.edu/asset\\_files/TechnicalReport/1994\\_005\\_001\\_435267.pdf](https://resources.sei.cmu.edu/asset_files/TechnicalReport/1994_005_001_435267.pdf).
- [24] Iso, Cmmi and Agile : A Comparison <https://www.vizteams.com/blog/iso-cmmi-and-agile-comparison/>.
- [25] Ajayi, W., Adekunle, Awodele, Akinsanya, Eze, & Seun, E. (2018). Software Development Top Models, Risks Control and Effect on Product Quality. *Global journal of computer science and technology*.
- [26] Vasile Selegean, (2015). Quality Assurance in Agile-SCRUM environment. *today software magazine*, issue-33. 2015. <https://www.todaysoftmag.com/article/1355/quality-assurance-in-agile-scrum-environment> visited: 12 May 2020.
- [27] Dave Litten. (2017). PRINCE2 quality control. *Prince2 Primer*. <https://www.prince2primer.com/prince2-2017-quality-control/> visited: 12 May 2020.
- [28] Simon Buehring, (2020). PRINCE2 benefits, advantages and disadvantages. <https://www.knowledgetrain.co.uk/project-management/prince2/prince2-benefits>. visited: 12 May 2020.
- [29] QATestLab, (2018). QA Activities in Waterfall Process. *Knowledge Center*. <https://qatestlab.com/resources/knowledge-center/waterfall-process/>. visited: 12 May 2020.
- [30] M. A. Akbar et al., "Improving the quality of software development process by introducing a new methodology—AZ-model", *IEEE Access*, vol. 6, pp. 4811-4823, Dec. 2017.
- [31] Nicola Jones, (2018). The information factories -Data centers are chewing up vast amounts of energy — so researchers are trying to make them more efficient. *Nature* 561, 163–166; 2018 available: <https://media.nature.com/original/magazine-assets/d41586-018-06610-y/d41586-018-06610-y.pdf> visited: 13 May 2020.
- [32] Andrae, A.S.G.; Edler, T. On Global Electricity Usage of Communication Technology: Trends to 2030. *Challenges* 2015, 6, 117-157.
- [33] Nicola Jones, (2018). How to stop data centres from gobbling up the world’s electricity. *nature*. 13th september 2018. <https://www.nature.com/articles/d41586-018-06610-y>. visited: May 15 2020.
- [34] Ericsson Energy and Carbon Report 2015. <https://www.ericsson.com/assets/local/about-ericsson/sustainability-and-corporate-responsibility/documents/ericsson-energy-and-carbon-report.pdf> (Accessed 13 May 2020).



[35] Malmodin, J., Bergmark, P. and Lundén, D. (2013), The future carbon footprint of the ICT and E&M sectors. Proceedings of the First International Conference on Information and Communication.

[36] Calero,C.; Piattini, M.IntroductiontoGreeninsoftwareengineering. GreenSoftw. Eng. 2015,3–27.

[37] Elbahrawy,A.; Alessandretti ,L. ;Kandler, A.; Pastor-satorras, R.Evolutionary dynamics of the cryptocurrency market. R. Soc. Open Sci. 2017, 1–16.

[38] David Bicknell. 2012. 8 ways to make your software applications more energy efficient. <https://www.computerweekly.com/blog/Green-Tech/8-ways-to-make-your-software-applications-more-energy-efficient> (visited 15 May 2020).

[39] Zahereel Ishwar Abdul Khalib, R. Badlishah Ahmad, Ong Bi Lynn, "Energy Efficient Scheduling Algorithm for Soft Real Time System with High Deadline Meeting Rate on Overload", Applied Mechanics and Materials, vol. 699, pp. 840, 2014.

[40] Vaquerizo-Hdez, D., Muñoz, P., R-Moreno, M. D., & F Barrero, D. (2017). A Low Power Consumption Algorithm for Efficient Energy Consumption in ZigBee Motes. Sensors (Basel, Switzerland), 17(10), 2179. <https://doi.org/10.3390/s17102179>.

[41] V.K. Chawlaa, A.K. Chanda, S. Angra and G.R.Chawla. (2018). "The sustainable project management: A review and future possibilities" Journal of Project Management 3 (2018) 157–170.

[42] Martens, M. L., & Carvalho, M. M. (2017). Key factors of sustainability in project management context: A survey exploring the project managers' perspective. International Journal of Project Management, 35(6), 1084-1102.

[43] Khodadadzadeh, T. (2016). Green building project management: obstacles and solutions for sustainable development. Journal of Project Management, 1(1), 21-26.

[44] Silvius, A. J., & Schipper, R. P. (2014). Sustainability in project management: A literature review and impact analysis. Social Business, 4(1), 63-96.

[45] Aarseth, W., Ahola, T., Aaltonen, K., Økland, A., & Andersen, B. (2017). Project sustainability strategies: A systematic literature review. International Journal of Project Management, 35(6), 1071-1083.

[46] Kivilä, J., Martinsuo, M., & Vuorinen, L. (2017). Sustainable project management through project control in infrastructure projects. International Journal of Project Management, 35(6), 1167-1183.

[47] International Organization for Standardization (2010). ISO 26000. Guidance on Social Responsibility, Geneva.

[48] Stavrou, P. (2018) Addiction to Video Games: A Case Study on the Effectiveness of Psychodynamic Psychotherapy on a Teenage Addict Struggling with Low Self-Esteem and Aggression Issues. Psychology, 9, 2436-2456. doi: 10.4236/psych.2018.910140.

[49] Austin,(2019).Smart devices may be hazardous to your health Contributor,THE TRUTH WILL OUT. CIO. JULY 2019, <https://www.cio.com/article/3411951/smart-devices-may-be-hazardous-to-your-health.html>. visited: 16 may 2020.

[50] Christoph, B. Sustainability and longevity: Two sides of the same quality? CEUR Workshop Proc. 2014, 1216, 1–6.

[51] Becker, C.; Chitchyan, R.; Duboc, L.; Easterbrook, S.; Penzenstadler, B.; Seyff, N.; Venters, C.C. Sustainability Design and Software: The Karlskrona Manifesto. In Proceedings of the 2015 IEEE/ACM 37th IEEE International Conference on Software Engineering, Florence, Italy, 16–24 May 2015.

[52] Oyediji S, Seffah A, Penzenstadler B. A Catalogue Supporting Software Sustainability Design. *Sustainability*. 2018; 10(7):2296.

**Appendix-1**

| This is no model that emphasis on sustainability, re-usability of project resources, and power-saving features |   |
|--|---|
| Quality models   | Quality attributes  |
| McCall's quality mode  | Product version: Maintainability (Simplicity, Conciseness, Self-descriptiveness, Modularity); Flexibility (Self-descriptiveness, Expandability, Generality); Testability (Simplicity, {instrumentation, Self-descriptiveness, Modularity}<br>Product operations: Correctness (Traceability, Completeness, Consistency); Efficiency (Execution, efficiency, Storage, efficiency); Reliability (Consistency, Accuracy, Error, tolerance); Integrity {Access control, Access audit); Usability (Operability, Training, Communicativeness}                        |
| Boehm's quality  | Utility: Reliability (Self Containedness, Accuracy, Completeness, Robustness/Integrity, Consistency); Efficiency (Accountability, Device Efficiency, Accessibility); Human Engineering (Robustness/Integrity, Accessibility, Communicativeness)Portability: Device independence; Self Containedness<br>Maintainability: Testability{ Accountability, Communicativeness, Self Descriptiveness, Structuredness}; Understandability (Consistency, Structuredness, Conciseness, Legibility); Modifiability (Structuredness, Augment ability}                      |
| Dromey's Quality model   | Correctness: Functionality, Reliability<br>Intemal, Maintainability, Efficiency, Reliability<br>Contextual: Maintainability, Reusability, Portability, Reliability<br>Descriptive: Maintainability, Efficiency, Reliability, Usability  |
| FURPS Quality Model  | Functionality: feature sets, capabilities, security<br>Usability. human factors, aesthetics, consistency, Correctness,<br>Reliability: severity of failure, recoverability, predictability, accuracy<br>Performance: speed, efficiency, availability, accuracy, throughput, response time, recovery time, and resource usage<br>Supportability: testability, extensibility, adaptability, maintainability, compatibility, serviceability, installability, and localizability.   |
| ISO 9126 quality model   | Functionality: Suitability, Accuracy, Interoperability, Security, Functionality Compliance<br>Reliability, Maturity, Fault Tolerance, Recoverability, Reliability Compliance<br>Usability: Understandability, Leamability, Operability, Attractiveness, Usability Compliance<br>Efficiency: Time Behavior Resource Utilization, Efficiency Compliance<br>Maintainability: Analyzability, Changeability, Stability, Testability, Maintainability Compliance<br>Portability: Adaptability, Installability, Co-existence, Replaceability, Portability Compliance |

# Forecasting the Global Horizontal Irradiance based on Boruta Algorithm and Artificial Neural Networks using a Lower Cost

Abdulatif Aoihan Alresheedi<sup>1</sup>, Mohammed Abdullah Al-Hagery<sup>2</sup>

Department of Computer Science, College of Computer, Qassim University, Buraydah, Saudi Arabia<sup>1,2</sup>  
BIND Research Group, College of Computer, Qassim University, Buraydah, Saudi Arabia<sup>2</sup>

**Abstract**—More solar-based electricity generation stations have been established markedly in recent years as new and an important source of renewable energy. That is to ensure a more efficient, reliable integration of solar power to overcome several challenges such as, the future forecasting, the costly equipment in the metrological stations. One of the effective prediction methods is Artificial Neural Networks (ANN) and the Boruta algorithm for optimal attributes selection, to train the proposed prediction model to obtain high accurate prediction performance at a lower cost. The precise goal of this research is to predict the Global Horizontal Irradiance (GHI) by building the ANN model. Also, reducing the total number of GHI prediction attributes/features consequently reducing the cost of devices and equipment required to predict this important factor. The dataset applied in this research is real data, collected from 2015-2018 by solar and meteorological stations in KSA. It provided by King Abdullah City for Atomic and Renewable Energy (KA CARE). The findings emphasize the achievement of accurate predictions of solar radiation with a minimum cost, which is considered to be highly important in KSA and all other countries that have a similar environment.

**Keywords**—Global horizontal irradiance; artificial neural networks; feature selection; boruta algorithm; cost reduction; machine learning

## I. INTRODUCTION

Alternative energy sources increasingly form the future of the world's energy system. This is due to fossil fuel resource limitations as well as their negative side effects on climate change and environmental pollution. The objective of sustainable power supply can be achieved by renewable energy sources, such as solar power, which still unused and it is characterized by high variability in availability and production. Rapid changes in solar power output are one of the negative consequences of rapid changes in weather conditions. The intermittent nature of renewable energy sources might hinder electrical utilities from effectively utilizing them.

Higher penetrations of solar energy into the electrical grid cause a more variable power output than with higher penetrations of wind [1]. Moreover, higher penetrations of alternative sources lead to power system technical operation and design issues, such as systems protection, systems control, power factor quality, and optimal operation of power systems [2]. Also, to manage the variability and the uncertainty in solar power, that is due to high penetrations of renewables,

adjustments to the power systems operation are needed including adding new ancillary services [3]. Thus, the economic feasibility of renewable energy sources is negatively affected by the expensive costs of these adjustments and requirements.

Many potential solutions can manage technical issues caused by short-term uncertainty in solar power (up to seven days ahead) [4]. For instance, increasing the level of demand-side participation, increasing the level of coordination to balance the allocation, and deploying more flexible but often also costlier energy storage systems. Still, the prediction of global solar radiation is one of the most efficient and economical ways to integrate more solar power, especially at current levels of integration. These forecasts can be utilized by balancing authorities to operate electric power systems more efficiently and reliably. In the literature, several forecasting approaches have been embraced [5]. Among these, Machine Learning (ML) algorithms are currently the most common methods to predict solar energy, because prediction is an important step in designing and assessing photovoltaic systems technically and economically.

Predictions of solar irradiance using ML algorithms were proposed in several studies, after the advent of fast computing capabilities as well as systems able to store massive data sets [6]–[9]. The ML methods include ANN, support vector regression (SVR), decision tree regression, and K-Nearest Neighbors, other methods [10], [11]. ANN is considered to be the most powerful ML method because of its capability to intrinsically deal with the nonlinear nature of solar and meteorological data. In recent studies, ANNs resulted in a lower mean absolute percentage error.

ANN was used, for example, to predict the electric load of Tai's power system [12]. Besides, in Salerno, Italy, two ANN models were developed to predict GHI and direct normal irradiance (DNI) in an hourly manner [13]. In the two later studies, the ANN model resulted in good accuracy. Also, ANN ensemble methods were applied by Alobaidi et al. to forecast the solar radiation variables utilizing satellite images. Five locations in the state of the United Arab Emirates were selected to apply the developed ANN prediction models [14]. As an application to forecasting one day-ahead solar radiation in a grid-connected-PV system, Mellit and Pavan developed ANNs that use mean daily solar radiation as well as air temperature as inputs into the prediction system [15]. Also, Lam et al.

employed the ANN to forecast the daily GHI in 40 different cities in China utilizing the observed duration of sunshine, and that study targeted areas with various thermal climates as well as sub-areas [16]. Moreover, a combination model of numerical weather prediction, using a hybridized autoregressive moving average and ANN algorithms were developed for forecasting the short-term GHI as presented [17].

An advanced embedded feature selection algorithm is known as the Boruta algorithm was implemented in this paper to choose from 13 available attributes in the dataset the most significant attributes. To the best of the author's knowledge, the feature selection approach employed by this research has not been applied to this issue before. This adds a great contribution to this paper. Since the analysis involves big datasets, authors embrace the powerful machine learning algorithm of ANN as a means of the computing system. The type of ANN used in this paper is a Multilayer-Feed-forward Back Propagation (MFBP) Network. In short, this paper introduces a novel, intelligent, a hybrid framework consisting of the ANN algorithm to conduct the training and testing processes and Boruta algorithm to select the most important features to be inputted into the ANN model. As a result, this paper introduces an ML-based model to predict GHI for each location of interest-based on the optimal number of attributes and it is an extension of limited work in [18], which focused partially on Buridah city in KSA.

The rest of this paper is organized as follows: Section 2 highlights the framework of the developed methodology including the data used and the used ML algorithms. Section 3, presents an overview of the feature selection technique used. Section 4 discusses the validation measures and metrics used to evaluate the accuracy of the proposed model. Besides, the analysis and discussion of the results were placed in Section 5. Finally, Section 6 explains the conclusions and future work.

## II. LITERATURE REVIEW

There are many research works concentrated in the field of energy, electricity generation, solar energy prediction methods based on ML algorithms, and other methods were based on mathematical models.

Several ML algorithms employed in the energy field for prediction purpose, such as the Support Vector Machine (SVM). This is because of the ability of SVM to model nonlinearity exists in time-series metrological data. Utilizing SVR applications to predict GHI was used and the results reveal feasibility and accurate prediction performance [14]. Also, the study achieved by [19] developed an SVM model based on a firefly algorithm to predict GHI [20]. Performance comparisons between the developed model and ANN and genetic programming models were created, and the results demonstrated that the enhanced SVR model has a better prediction accuracy. Also, another wavelet-based SVR model was developed in to predict GHI in different cities in Australia.

The research work by [21] concluded that the prediction accuracy in SVR models is directly proportional to the size of training data when SVM applied to predict electric load. Similarly, the SVR model accuracy considerably relies on selecting the optimal set of parameters. A proper determination of the optimal set of SVR models' parameters is not an easy

task. To solve this problem, several advanced optimization techniques have been used such as particle Swarm Optimization algorithm, Immune algorithms [22]. Besides, Genetic algorithm, for example, has been used to optimally select the SVR model parameters to forecast electricity market prices [23]. In Saudi Arabia, past studies on models of solar radiation applied various computational methods, most of which belonged to methods of empirical or artificial intelligence. The researches by [24], [25] carried out forecasting of the average of GHI per annum with good accuracy. A nonlinear Angstrom-type model was used in their study and then was compared to Bulut and Büyükalaca's trigonometric function model in [26]. Also, a related study achieved in Oman [27] for measuring various features such as the temperature, humidity, and solar radiation. The study provided statistical results compared to the NASA SSE Model.

A geostatistical methodology was used by [28] to predict GHI in the Kingdom of Saudi Arabia. This study had the purpose of producing a geographically persistent mapping system of solar irradiance, and also for every single month of the year, to draw the solar irradiance contour maps. In a mission originating from 1994 to April 2000, solar radiation measurements were taken over twelve locations in 12 cities [29]. Utilizing features of latitude, longitude, altitude, the number of months, and sunshine duration, Mohandes and Rehman used a predictive approach to forecast GHI anywhere across Saudi Arabia [30], [31]. The experiment used the 35 stations solar radiation data to test the accuracy of the prediction model where the outcomes of the forecasted values were near to the observed values to some extent. Benghanem et al. employed ANN-based prediction models to forecast daily averages of global horizontal irradiance for five years' period through the use of National Renewable Energy Laboratory repositories. Using recent data sets given by KA CARE. But in our study we add the features selection technique, to improve the results.

Almaraashi used automated fuzzy logic systems aiming at forecasting the next day's solar radiation [32]. To the best of the author's knowledge and even though many ML-based models have been introduced in the Kingdom of Saudi Arabia, no automated methods of feature selection have been examined to forecast short-term solar radiation in Saudi Arabia. In addition to the mentioned works, there are some recent researches are concentrating on the GHI forecasting but with different datasets and different strategies.

For instance, in northeast Iraq [33], a research study accomplished on a Satellite Datasets to obtain a more accurate and precise method for forecasting hourly GHI. The proposed method established based on the ANN and another training algorithm called "Levenberg Marquardt" algorithm. The obtained results showed a very high accuracy. Besides, in South Africa [34], a research study carried out for discussing probabilistic of forecasting the GHI before 24 hours, using two machine learning methods and the data collected during the period from 2009 to 2010. The study gave excellent results but not exceeded by 95.5%.

Moreover, in Croatia [35], a research study concentrated on several models that are used for estimating solar radiation.

These models assessed based on seven meteorological stations dataset, where the models studied, compared and evaluated to find out the best accuracy. As well, in southern Finland [36], a research study carried out a GHI forecasting using a data set of weather satellite imagery, using a mathematical modelling method. The results obtained show very good accuracy. In this regards, it noticed that in many countries, the number of GHI measurement locations is sufficient as in KSA and insufficient in other countries as in Korea [37], [38], where the satellite images can be helpful sources for getting the GHI over a wide area space, in these cases, usually predicted by secondary parameters such as readily obtainable climatic variables.

On the other hand, regarding efficient attributes selection and efficient data preprocessing, the feature selection techniques are utilized for minimizing and preparing data with high dimensionality for ML-based problems. Such techniques are usually categorized into either supervised algorithms, requiring information about labels, or unsupervised algorithms, that operate without a need for information about labels. The challenge is that the solar radiation intensity is influenced by a large number of parameters. Thus, by removing redundant attributes, dimensionality reduction algorithms, as well as feature selection, might positively affect the prediction accuracy of developed forecasting models. Furthermore, the need for an optimized feature space grows when broad degrees of uncertainty is involved in the considered application. Several studies have identified and discussed the need to have a feature selection approach to be embraced before forecasting GHI.

For example, Salcedo-Sanz et al. [39], examined the usage of a species-optimizing coral reef algorithm to gain a decreased collection of important features to forecast GHI. Also, Yadav et al. implemented a set of features to some specific input predictors and observed the parameters of latitude and longitude are having the slightest impact on the forecasting of solar radiation [40]. Hedar et al. applied a programming-based algorithm of adaptive memory to minimize the space of input features of a fuzzy classifier for global solar radiation [41]. They found that among nine attributes, DHI, DNI, and relative humidity have the best dependence degree values.

This paper concentrates on the prediction of the GHI in two different regions in Saudi Arabia, by building Neural Networks models whose input variables are optimally and systemically selected by the features selection algorithm named Boruta, which was used for the first time to improve the GHI forecasting results.

### III. METHODOLOGY

The overall methodology steps are listed as follows:

- Data-preprocessing tasks of all datasets used in this research are carried out ahead of the training, validation, and testing process of the proposed data-based model.
- The model trained and validated utilizing the all-feature set is established.

- In a similar way to (2), the model trained and validated utilizing the most important eight-feature set determined by the Boruta Algorithm is also built.
- Moreover, the model trained and validated utilizing the most important five-feature set by Boruta Algorithm is also developed.
- The predicted values of GHI by the model with different features are carried out through the testing process.
- To evaluate the performance accuracy of the developed forecasting model, the observed and the predicted values of the GHI are compared by a set of four evaluation metrics.

The Boruta is used to pick the most important variables among a wide range of meteorological variables that could impact solar radiation in the future. A prediction processed are then independently utilized, based on the Boruta's five and eight most important variables. Even though the targeted feature-selection-based model can be designed using any number of variables, we have fixed the number of features chosen to provide a fair comparison at the various locations of this study between the forecasting developed. Strictly speaking, in this case, a smaller set of features (that is, 5 and 8) must be made in advance among the 13 features.

#### A. Data Collection

To build the proposed forecasting model proposed in this paper, massively big observed solar and meteorological datasets are collected from KA CARE that provided more than 25 solar and metrological variables. The total number of observations (records) of the collected data is 35735 in the Qassim dataset, while in Jeddah city is 35856 observations before the pre-processing step. After the cleaning process, the data were reduced at the Qassim region to be 17892 and in Jeddah to be 19467 observations. The GHI observations under consideration are in 1-hour time resolutions for the period from March 1, 2013, out to June 30, 2017, and they are collected from the interest locations of Jeddah and Qassim. Furthermore, the corresponding 1-hour intervals weather variables are gathered. The whole dataset ( $x$ ) is splatted into three subsets namely: the training dataset,  $x_{training}$ , the cross-validation dataset,  $x_{cross-validation}$ , and the testing dataset,  $x_{testing}$ , such that  $X = x_{training} \cup x_{cross-validation} \cup x_{testing}$ . In this research, the ratio of the training, cross-validation, and testing datasets are 6:2:2, respectively.

At this time, the variables involved in this analysis include 13 independent variables: month of the year (M), day of the month (D), an hour of the day (H), air temperature (T) in ( $^{\circ}$ C), relative humidity (RH) in (%), surface pressure (P) in (hPa), wind speed at 3 meters (WS) in (m/s), Wind Direction (WD) at 3 meters in ( $^{\circ}$ N), Peak Wind Direction (PWD) at 3 meters in ( $^{\circ}$ N), diffuse horizontal irradiance (DHI) in (Wh/m<sup>2</sup>), direct normal irradiance (DNI) in (Wh/m<sup>2</sup>), azimuth angle (AA) in ( $\hat{A}^{\circ}$ ), and solar zenith angle (SZA). The GHI in (Wh/m<sup>2</sup>), as a target variable, is measured with a Kipp & Zonen Pyranometer.

Future work will involve an increased number of the observed features of the data as inputs into a novel proposed model. The input variables to the prediction model can be summarized in Table I. The prediction model can be expressed, in the scope of the used predictors, as shown in equation (1).

$$GHI_{\text{predicted}} = f(M, D, H, T, RH, P, WS, WD, PWD, DHI, DNI, AA, SZA) \quad (1)$$

TABLE I. INPUT VARIABLES TO THE PROPOSED PREDICTION MODEL

| Input variable | Input variable abbreviation | Input variable explanation          | Input variable unit |
|----------------|-----------------------------|-------------------------------------|---------------------|
| $x_1^{(i)}$    | M                           | the month of the year               | Month               |
| $x_2^{(i)}$    | D                           | day of the month                    | day                 |
| $x_3^{(i)}$    | H                           | hour of the day                     | hour                |
| $x_4^{(i)}$    | T                           | air temperature                     | °C                  |
| $x_5^{(i)}$    | RH                          | relative humidity                   | %                   |
| $x_6^{(i)}$    | P                           | surface pressure                    | hPa                 |
| $x_7^{(i)}$    | WS                          | wind speed at 3 meters              | m/s                 |
| $x_8^{(i)}$    | WD                          | Wind Direction                      | °N                  |
| $x_9^{(i)}$    | PWD                         | the peak wind direction at 3 meters | °N                  |
| $x_{10}^{(i)}$ | DHI                         | diffuse horizontal irradiance       | Wh/m <sup>2</sup>   |
| $x_{11}^{(i)}$ | DNI                         | direct normal irradiance            | Wh/m <sup>2</sup>   |
| $x_{12}^{(i)}$ | AA                          | azimuth angle                       | Â°                  |
| $x_{13}^{(i)}$ | SZA                         | solar zenith angle                  | Â°                  |

### B. Data Preprocessing

After a thorough inspection of the used datasets from KA CARE, we notice that there are some errors, noise, redundant records. For that, some of the data cleaning steps are applied. These steps are very important to have high-quality datasets because unclean data can decrease the classification or regression model accuracies [42]. Fig. 1 shows the flowchart of the data preprocessing steps, and these steps can be summarized in the below sections.

1) *Organizing dataset*: It is a process of transforming the data received to a common format to improve visualizing and dealing with the data.

2) *Removing redundant data*: The performance of the forecasting algorithms depends mainly on the amount and accuracy of the used data. Using redundant data to train and test the algorithms will make the model computationally expensive as well as increase the time of executing the algorithms. The metrological datasets received from KA CARE contain some duplicated data for the same features. For example, for Air Temperature, there are data for actual temperature as well as data for uncertainty in Air Temperature. The uncertainty in Air Temperature adds no value and is redundant data. Therefore, it should be removed to avoid any complications in the data. Similarly, for wind direction, wind speed, DHI, DNI, GHI, peak wind speed, relative humidity, and barometric pressure.

3) *Monitoring data errors*: After organizing the data and from initial scanning, the data file shows that there are missing data entries in some features. This is shown in the data file as empty cells and MATLAB as (NaN), meaning Not a Number in a numerical file. In our final dataset, the entire day at any of the unrecorded features is removed when creating the forecasting model.

4) *Feature construction/selection*: Where new attributes (features) are constructed and added from the given set of attributes to help the mining process. The format of the date of KA CARE is as follows: 01/01/2014, 12:00:00 AM (MM/DD/YYYY HH:mm:SS), where MM: month, DD: day, YYYY: year, HH: hour, mm: minute and SS: second. Day, Month, and Hour are used as features for training and testing in the forecasting algorithms.

These attributes are needed to be re-constructed to improve the mining process because the date format is not suitable for advanced mining. Therefore, splitting these variables into different columns is necessary.

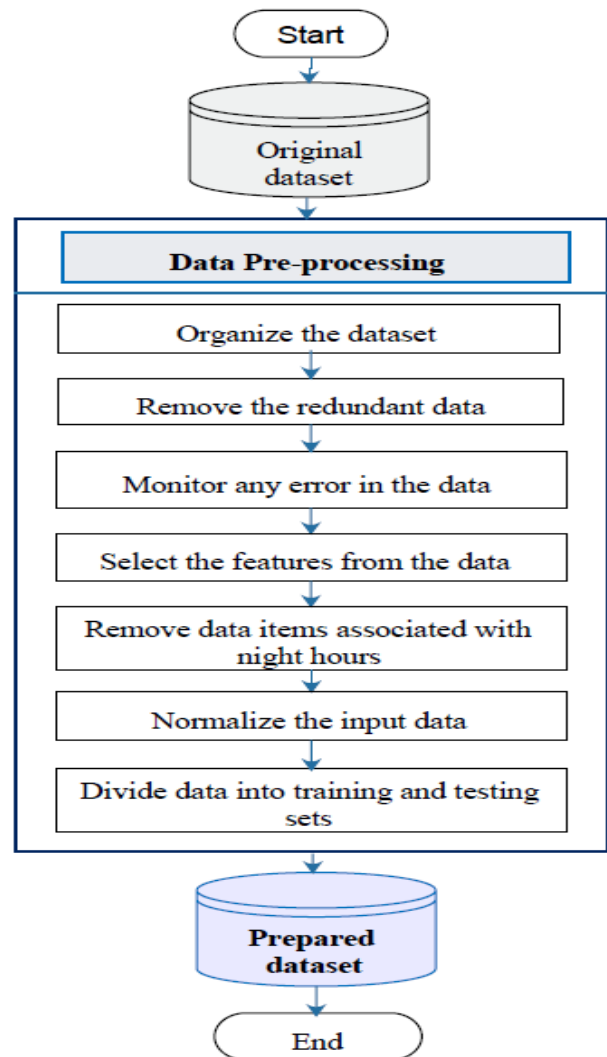


Fig. 1. Diagram of Data Preprocessing Steps.

The MATLAB code used to separate the data cells into single cells for the month, day, and hour. However, one of the problems encountered during preparing the data is inconsistency in the Date format. At each year and in each month of the year, the format of date from day 1 to day 12 is (Day, Month, Hour) while the remaining days' format is (Month, Day, Hour). This required creating a MATLAB code to change the format from day 1 to day 12 to make it smooth with the rest of the months' days (Month, Day, Hour). Accomplishing this task should be automatic because doing this manually is a very difficult task and time-consuming since we are dealing with a very huge dataset.

5) *Removing GHI night hours:* The main goal of this paper is to forecast the value of GHI. During the night, there is no solar radiation and the Pyranometer (a device used to measure GHI value) recorded zero values at night hours. These hours add no values to the forecasting model and removing them is very necessary. Therefore, all the night hours of the GHI on each day and all the corresponding features associated with it are removed from the dataset.

6) *Dividing data into training and testing:* After completing the aforementioned five steps, the data are divided into training and testing set in the ratio of 80% for the training process and 20% for the testing process.

The Pareto principle is a common rule of thumb to divide the dataset into two sub-sets; training and testing data. This is also called the 80/20 rule. The training and testing dataset are selected randomly. The objective of selecting the data randomly is to make our model be trained based on a variety of weather observations. These training and testing data are then fixed, and all the forecasting algorithms are encountering the same training and testing data input.

7) *Input data normalization:* Input data scaling, also known by normalization, is a very critical practical implementation when applying ANNs. The importance of this practical consideration is mainly to avoid the possible domination of attributes with greater numeric values upon those attributes with smaller ones. Another significant feature is to overcome numerical difficulties during computation processes. Because of the dependency of the kernel values on the inner products of the attribute vectors, attributes with large values cause numerical problems. In this research, each attribute is linearly normalized to the range: 0 to 1, using equation (2).

$$x_i^n = \frac{x_i - x_{min}}{x_{max} - x_{min}} \quad (2)$$

Where  $x_i$  is the actual value of the feature vector;  $x_{min}$  and  $x_{max}$  are the minimum and the maximum values corresponding to the actual dataset;  $x_i^n$  is the normalized value associated with  $x_i$ .

### C. Multilayer Feed-forward Back Propagating Neural Networks

Non-linear relationships between independent and dependent variables can be captured by Artificial Intelligence (AI) methods. One of the powerful non-linear forecasting

algorithms used here is an ANN. ANNs mimic how human nervous systems interpret information. Being ANNs are capable of modelling non-linear processes without a need to assume the relationship form between the input and output variables is considered one of the main advantages of this technique. The type of the ANN used here is shown in Fig. 2, it is a Multi-Layer Perception (MLP) [43]. To conduct the training process, the Back-Propagation Algorithm (BP) is selected in this research because it is one of the most common ANN algorithms [44]. As Fig. 2 depicts, the usual architecture of ANNs formed with three main layers. First, the input layer,  $[x_1, x_2, \dots, x_N]^T$ , which is composed of an N-dimensional input vector. After that, the hidden layer,  $[h_1, h_2, \dots, h_M]^T$ , which includes a nonlinear activation function known as the activation function. Finally, the output layer,  $[y_1, y_2, \dots, y_L]^T$ , which contains a linear function. Inside hidden layers, the outputs of nodes, also known as neurons, which represent the basic component of any ANN, can be calculated as below:

$$z_j = \sum_{i=0}^N v_{ij}x_i, j = 1, 2, \dots, M, i = 1, 2, \dots, N \quad (3)$$

$$h_j = f(z_j), j = 1, 2, \dots, M \quad (4)$$

In which,

- $z_j$  is the value of the activation function of the  $j$ th node associated with the hidden layer.
- $v_{ij}$  is the weight that connects the input  $i$ , in the input layer, with the node  $j$  in the hidden layer.
- $f$  is known as the transfer function of the neurons, often a sigmoid function is selected  $f(x) = \frac{1}{1 + \exp(-x)}$ .
- $h_j$  is the output value of the node  $j$  in the hidden layer.

In the output layer, the values of the output nodes can be calculated by using the equations (5) and (6).

$$z_l = \sum_{i=0}^M w_{jl}h_j, l = 1, 2, \dots, L, j = 1, 2, \dots, M \quad (5)$$

$$y_l = f(z_l), l = 1, 2, \dots, L \quad (6)$$

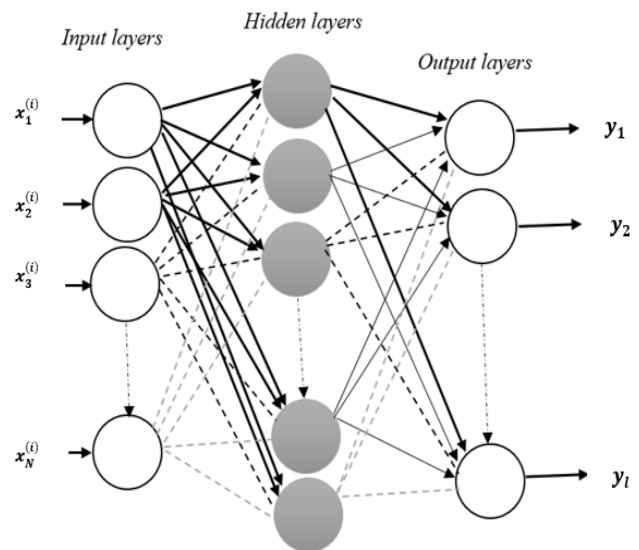


Fig. 2. The MFBP Network Model.

In which,

- $z_l$  is the value of the activation function associated with the  $l$ th node in the output layer.
- $w_{jl}$  the weight that connects the node  $j$  in the hidden layer with the node  $l$  in the output layer.
- $f$  is a sigmoid activation function.
- $y_l$  is the value of the activation function of the node  $l$  in the output layer.

Through experimenting with different choices, the required number of hidden layers, and the number of nodes in each layer were selected based on the optimal value that provides the best training prediction performance is reached. In this article, training networks for the model with three different sets of features were utilized by the MLP with the BP algorithm, while the Levenberg-Marquardt approach was the training function. One input layer, one hidden layer, and one output layer are used. The hidden layer in the ANN was constructed with 14, 8, and 5 neurons (nodes) for the model using the different sets of features; All-Feature, Eight-Feature, and Five-Feature. The input and the output of the training and testing dataset are similar for the model with the three sets.

#### D. Feature Selection Algorithm

In ML applications, feature selection, also known as variable selection, is frequently a crucial phase towards building a highly accurate ML-based model [45]. There are good reasons that support the use of feature selection algorithms. Nowadays, new datasets for practical model designing are usually described with a very high number of variables. Most of these variables are often irrelevant to the classification problem, and their significance is not established beforehand.

In dealing with data with too large feature sets, several disadvantages will appear. Practically, it found that dealing with massively large feature sets causes algorithms to slow down. Another reason is even more critical is the decrease in prediction accuracy is shown in many ML algorithms when dealing with higher than optimal sets of variables. Therefore, it is desirable for practical reasons to select the possibly smallest set of features that returns the best possible prediction results. This problem, known as the minimal-optimal problem, has been considerably researched where plenty of algorithms were developed to come up with manageable-size sets of features. Nevertheless, this very practical objective echoes another very important issue, which is the recognition of all features that are relevant to the classification problem under certain conditions. This is the so-called all-relevant problem. It can be very useful in itself to find all relevant features, rather than just non-redundant ones. This is especially necessary if one is interested in understanding processes related to the topic of interest, rather than simply building a predictive black-box model.

A good discussion on why it is important to find all relevant features is given by [46]. All relevant feature selection problems are more difficult to handle than normal minimal-optimal alternatives. To determine whether the variable is important or unimportant, therefore, we need a powerful

criterion to do so. The filtering methods can be used to select relevant features [47].

However, filtering methods are not the optimal choice for feature selection implementation due to the lack of a direct correlation between a particular feature and the decision that this feature is unimportant in combination with other features. Hence, one is limited to wrapper methods of feature selection that are more challenging computationally than filtering alternatives. As a black box, the classifier is used in wrapper methods to return a feature ranking, so any classifier that can return feature ranking can be used. In a short, a classifier used in feature selection problems should be both computationally effective and easy, optimally without any user-set parameters.

To find the all-relevant features in the solar radiation prediction problem, this paper utilizes the so-called R package Boruta [48]. This package is publically available from the Comprehensive R Archive Network at <http://CRAN.R-project.org/package=Boruta>. In this algorithm, a wrapper approach that is built around a random forest classifier is used [49].

In Slavic mythology, Boruta is the God of the forest. Selecting all relevant features, rather than just the non-redundant features. This algorithm is an extended version of the concept that Stoppiglia et al. implemented in [50] to assess relevance by comparing the relevance of the real features with that of the random probes. Although it is used in this paper as a wrapper algorithm, the concept is originally proposed in the context of filtering. In short, a brief overview of the algorithm is giving in the following section.

1) *Boruta algorithm*: Boruta Algorithm is a wrapper method that is designated based on the random forest regression algorithm executed in [51]. This paper utilized the regression version of the Boruta algorithm and its origin the random forest algorithm. However, to explain the basics of both algorithms, this article sticks to the classification versions of these algorithms.

The classification algorithm of the random forest is considered relatively fast, can normally run without a need for parameter setting, and it delivers a numerical approximation of the feature's importance. It is considered under the category of ensemble methods, which executes classification by acting on multiple unbiased poor classifiers recognized as decision trees. Such trees are freely and independently established upon different samples of bagging extracted from the training dataset. A feature importance metric is gained as the classification accuracy loss induced by the feature values' random permutation between objects. This measure of the importance of a feature is separately calculated for all trees available in the forest. These trees utilize a given feature for the classification task. Afterwards, the accuracy loss's mean and standard deviation are worked out. Dividing a feature's accuracy loss by its standard deviation results in the so-called Z score that can alternatively be used as the measure of the importance of a feature.

Unfortunately, since the distribution of the random forest algorithm is not  $N(0,1)$ , the Z score cannot be interpreted as a

direct relation to the statistical significance of the feature importance given by the random forest algorithm [52]. In the Boruta algorithm, nonetheless, since the Z score takes into consideration the average accuracy loss fluctuations among trees in the forest, we use it as the measure of the importance of a feature. Because the Z score cannot be directly used to calculate the importance, some external reference is needed to assess whether the significance of any given feature is important, i.e., whether it is perceptible from the significance of random variations. To that point, the information system has been expanded with random design features. We create a corresponding 'shadow' feature for each original feature, whose values are acquired by mixing values across objects from the original feature. We then use all the features of this extended information system to perform regression and measure the value of all features. Because of random fluctuations, the value of a shadow feature can be nonzero. Thus, the shadow features set of importance is utilized as a guide to decide, which features are considered significant. The significance indicator differs due to the random forest classifier stochasticity. Moreover, that is very sensitive to non-important features being present in the information system (as well as shadow features). It also depends on the specific realization of shadow features. Thus, to obtain statistically valid results, we need to perform the process of re-shuffling.

In short, the Boruta method is based on the same principle that serves as the foundations of the classifier of the random forest algorithm, that is by introducing randomness to the information system and gathering results from the randomized sample ensemble one can the misleading effect of random fluctuations and correlations. Thus, this added randomness will give us a clearer picture of which attributes matter significantly. The steps in which the Boruta Algorithm is executed consist of the following:

- Add copies of all features (variables/predictors) to expand the information system. Even if the number of features in the original dataset is smaller than 5, always extend the information system by at least 5 shadow features.
- To eliminate their correlation with the target variable, shuffle the added features.
- To collect the measured Z scores, run a classifier of the random forest upon the extended information system.
- Figure out the maximum Z score amongst shadow features (MZSF), and after that give a hit to any better-scored feature than MZSF.
- A two-sided equality test with the MZSF is conducted for each feature of undetermined significance.
- Consider the features of significantly lower value than MZSF as 'unimportant' and delete them permanently from the information system.
- Consider the features of significantly higher value than MZSF as 'important'.
- Delete all shadow features.
- Repeat the process until the importance for all the features has been allocated, or the algorithm has exceeded the random forest runs previously set.

#### *E. The Hybrid Strategy Proposed for Forecasting*

In this section, the designed hybrid big data-driven strategy for short-term global solar radiation forecasting based on the ANN and Boruta Algorithm, which is applied to select the optimal set of features to be inputted to the ANN algorithm, is introduced. The whole structure of the proposed hybrid system as a GHI forecasting model, as shown in Fig. 3.



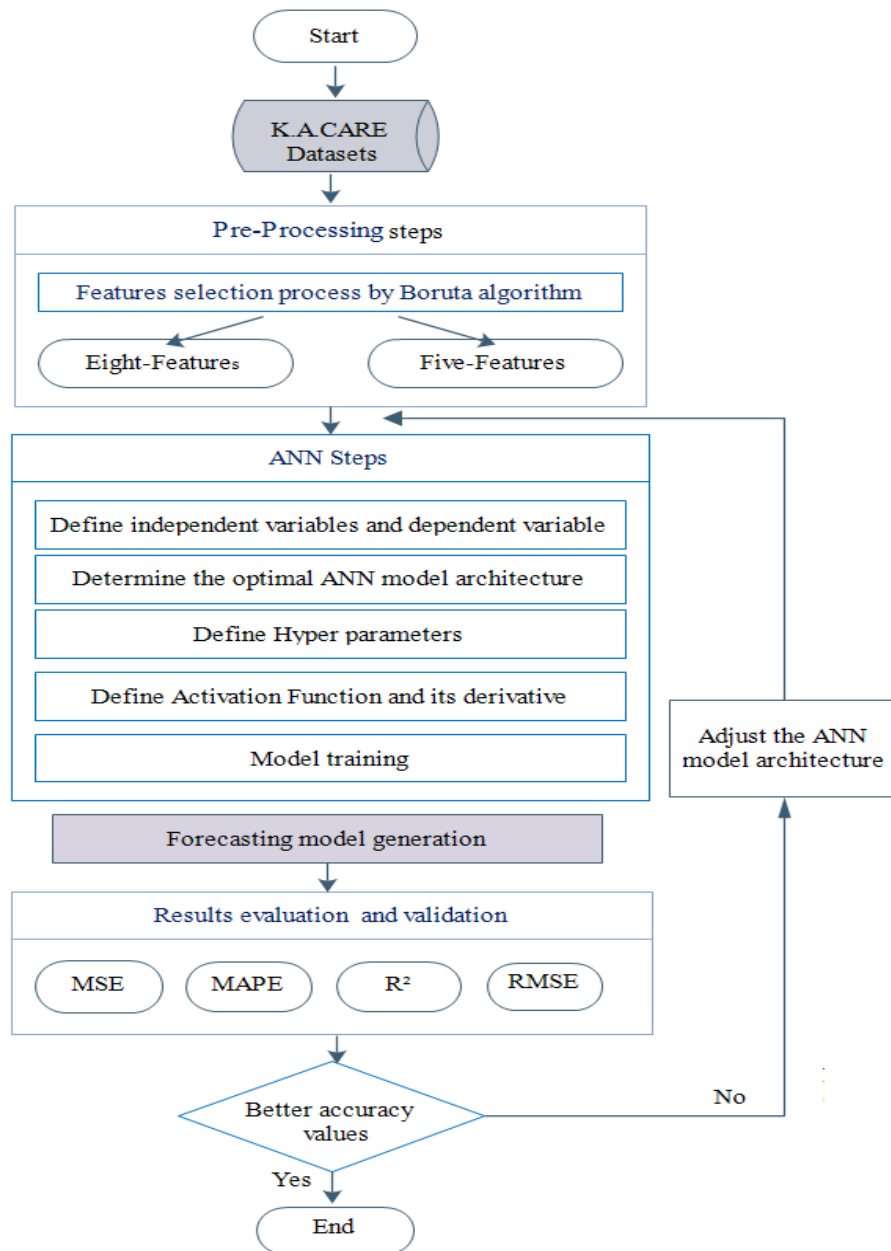


Fig. 3. The Proposed GHI Forecasting Model.

#### IV. EVALUATION MEASURES

Several statistical measures used to evaluate the prediction performance accuracy of the developed model. This research mainly considers three indicators, namely: mean absolute percentage error (MAPE), mean squared error (MSE), root mean squared error (RMSE), and goodness of fit ( $R^2$ ). Such measures are mathematically represented by the equations (7) to (10).

$$MAPE = \frac{1}{N} \sum_{i=1}^N \frac{|f_i - y_i|}{y_i} \times 100\% \quad (7)$$

$$MSE = \frac{1}{N} \sum_{i=1}^n (f_i - y_i)^2 \quad (8)$$

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (f_i - y_i)^2} \quad (9)$$

$$R^2 = \frac{\sum_{i=1}^N (f_i - \bar{y})^2}{\sum_{i=1}^N (y_i - \bar{y})^2} \quad (10)$$

Where N represents the number of data points involved in the analysis;  $y_i$  is the observed value of the target,  $f_i$  is the predicted value of the target;  $\bar{y}$  is the mean of the observed value of the target  $y_i$ . While MAPE is utilized to evaluate the model performance accuracy as a percentage, RMSE measures how the observed values deviate from the corresponding predicted values [13]. In regression problems, the  $R^2$  of a model describes how well the model fits a set of observations.  $R^2$  ranges from zero to the preferable number 1.

## V. RESULTS AND DISCUSSION

The forecasted model is employed to predict GHI at two selected sites in Saudi Arabia, namely Qassim and Jeddah. In this research, the prediction module utilized based on

- All-Features;
- Eight-Feature; and
- Five-Feature.

In this research, the model was implemented using MATLAB. The data was first cleaned and after that normalized to increase the performance of the forecasting and feature selection algorithms.

A set of eight and five features out of the available 13 features were identified to create the so-called Eight-Feature and Five-Feature forecasting model. The choice of these features was based on the output of the feature selection algorithm, Boruta. Eight and five features were selected to demonstrate the performance of the forecasting model with a

different number of features. Fig. 4 and Fig. 5 rank the features based on their importance to our outcome (GHI) at Qassim and Jeddah, respectively.

Fig. 4 and Fig. 5 show that in the Qassim region, the Eight-Feature model is built based on the following features: DNI, Zenith Angle, DHI, Hour, Month, Pressure, and Azimuth Angle. On the other hand, Jeddah's Eight-Feature forecasting model is based on the following features: DNI, Zenith angle, DHI, pressure, hour, month, Azimuth angle, and temperature. For the Five-Feature model, the Qassim region forecasting model is based on the following features: DNI, Zenith Angle, DHI, Hour, and Month, while Jeddah's Five-Feature forecasting model is created by using the following features: DNI, Zenith angle, DHI, pressure, and hour. The model developed was utilizing the MLP with a BP, while the Levenberg-Marquardt approach was the training function. One input layer, one hidden layer, and one output layer are used. The hidden layer in the ANN was constructed with 14, 8, and 5 neurons (nodes) for All-Feature, Eight-Feature, and Five-Feature model.

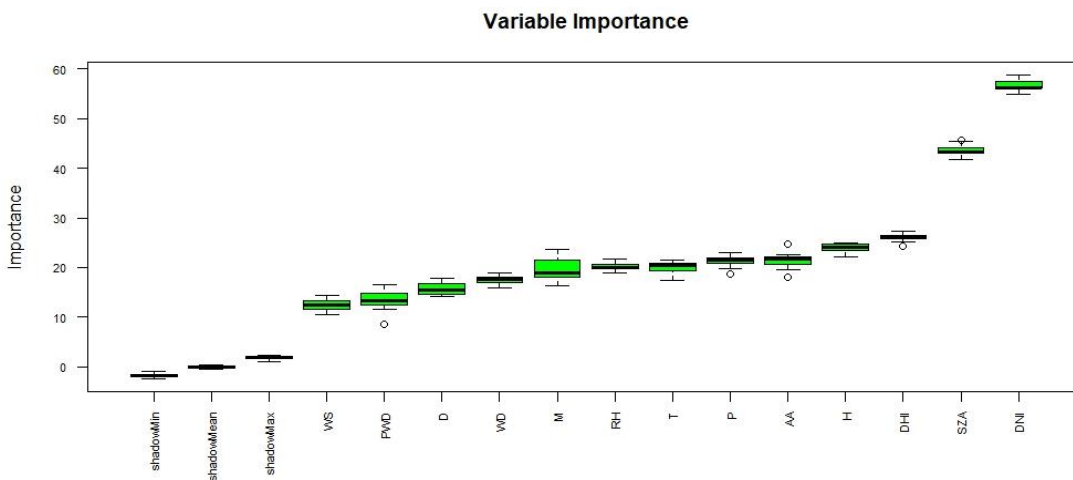


Fig. 4. Variable Importance for Weather Data for the Qassim Region using the Boruta Algorithm.

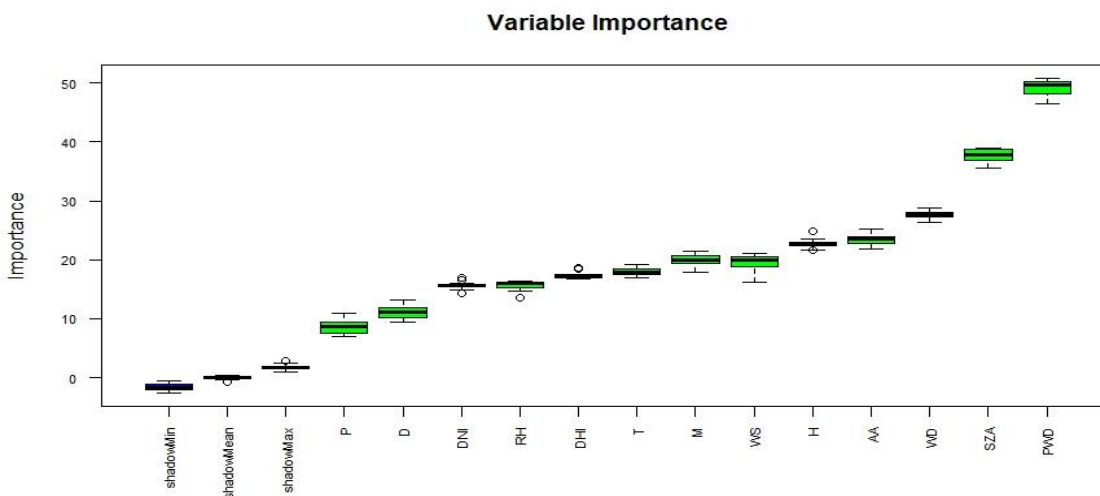


Fig. 5. Variable Importance for Weather Data for the Jeddah Region using the Boruta Algorithm.

The comparisons were conducted between the results of the All-Feature model with the results based on both the Eight-Feature and Five-Feature model in terms of their forecasting accuracy at Qassim and Jeddah sites. For this goal, the accuracy outputs of the testing model were tested based on the following metrics: MAPE, MSE, RMSE, and  $R^2$ . Table II and Table III compare the results of the model for the three different sets of features at Qassim and Jeddah.

MAPE determines the accuracy and errors ratio between measured and predicted data, while MSE and RMSE measure the relative error and expressed in ( $Watt/m^2$ ). Table II and Table III contain the hourly forecasting of GHI based on the study of the model at Qassim and Jeddah. In Qassim and from Table II, with the All-Feature model MAPE value found to be 13.496% (16.532% with Eight-Feature and 26.563% with Five-Feature).

The MSE and RMSE values of All-Feature model are  $1708.957 Watt/m^2$  ( $1756.145 Watt/m^2$  with Eight-Feature and  $2360.716 Watt/m^2$  with Five-Feature) and  $41.339 Watt/m^2$  ( $41.906 Watt/m^2$  with Eight-Feature and  $48.587 Watt/m^2$  with Five-Feature), respectively. The correlation scores of the forecasting model at Qassim was found to be 0.99124, 0.99167, and 0.9794 for All-Feature, Eight-Feature, and Five-Feature, respectively.

The results indicate that the proposed All-Feature model has the best performance compared to the eight and Five-Feature model. However, the Eight-Feature model performance is high in a way that can be compared with the All-Feature model, while Five-Feature can be considered the poorest model, still, it has satisfactory results. According to Table III that presents Jeddah results, the Five-Feature model can be considered as the best model followed by Eight-Feature and All-Feature model, respectively. Unlike Qassim site, Jeddah forecasting model with merely five and eight features prove the significance of using a feature selecting approach and how adding more feature may lead to overfitting forecasting model. In Jeddah, the MAPE value was determined to be 13.7013% with All-Feature (12.2024% with Eight-Feature and 9.6936% with Five-Feature).

TABLE II. RESULTS OF HOURLY FORECASTING GHI AT QASSIM

|               | MAPE % | MSE $Watt/m^2$ | RMSE $Watt/m^2$ | $R^2$ |
|---------------|--------|----------------|-----------------|-------|
| All Feature   | 13.50  | 1708.96        | 41.34           | 0.99  |
| Eight Feature | 16.50  | 1756.15        | 41.91           | 0.99  |
| Five Feature  | 26.60  | 2360.72        | 48.59           | 0.98  |

TABLE III. RESULTS OF HOURLY FORECASTING GHI AT JEDDAH

|               | MAPE % | MSE $Watt/m^2$ | RMSE $Watt/m^2$ | $R^2$ |
|---------------|--------|----------------|-----------------|-------|
| All Feature   | 13.70  | 994.37         | 31.54           | 0.99  |
| Eight Feature | 12.20  | 939.87         | 30.66           | 0.99  |
| Five Feature  | 9.70   | 913.84         | 30.23           | 0.99  |

The MSE and RMSE values of All-Feature model are  $994.369 Watt/m^2$  ( $939.869 Watt/m^2$  with Eight-Feature and  $913.84 Watt/m^2$  with Five-Feature) and  $31.533 Watt/m^2$  ( $30.6572 Watt/m^2$  with Eight-Feature and  $48.587 Watt/m^2$  with Five-Feature), respectively.

The correlation scores of the forecasting model at Qassim were found to be 0.99124, 0.99168, and 0.99184 for All-Feature, Eight-Feature, and Five-Feature, respectively. For further visualization, the measured GHI values are plotted against the output of the forecasting model by three different sets of Features at Qassim as in Fig. 6 and Jeddah as in Fig. 7. Where Fig. 6 confirms that the All-Feature model has high accuracy results at Qassim, and Fig. 9 shows how the Five-Feature model performs the better compare to All-Feature and Eight-Feature model in Jeddah. The model results of All-Feature, Eight-Feature, and Five-Feature are plotted all together with measures GHI values in Fig. 8 at Qassim and Fig. 9 at Jeddah for twenty random hours.

Fig. 8 shows that the All-Feature model has superior performance in tracking the original GHI value at Qassim, and Fig. 9 confirms the ability of the Five-Feature model in following the measured GHI values at Jeddah. In regards to the cost reduction by the generated model, Table IV illustrates the prices of devices and equipment required for the prediction purpose of the GHI. This table also shows the cost reduction by the forecasting model using eight and five features for the two regions.

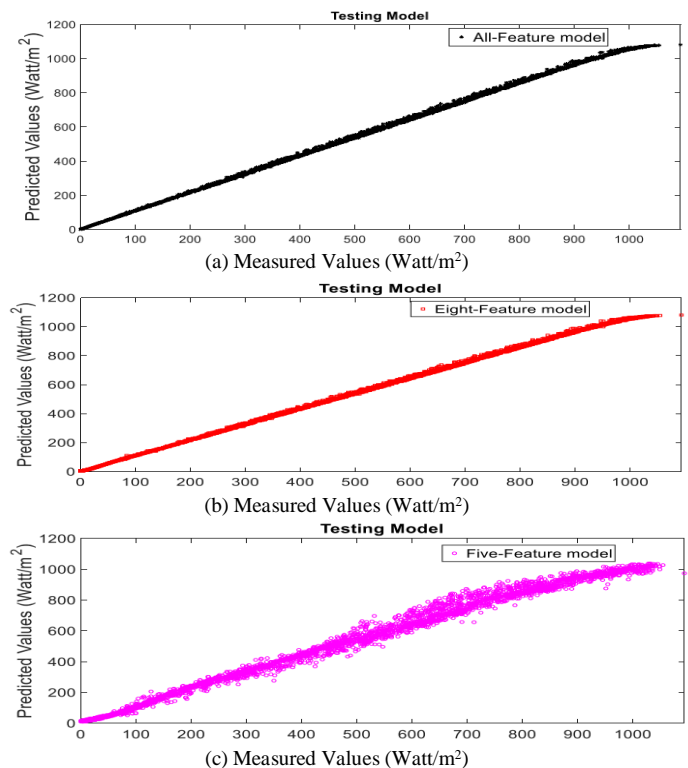


Fig. 6. Measured vs. Predicted Values of GHI for Qassim Region, (a) All-Feature Model, (b) Eight-Feature Model, (c) Five-Feature Model.

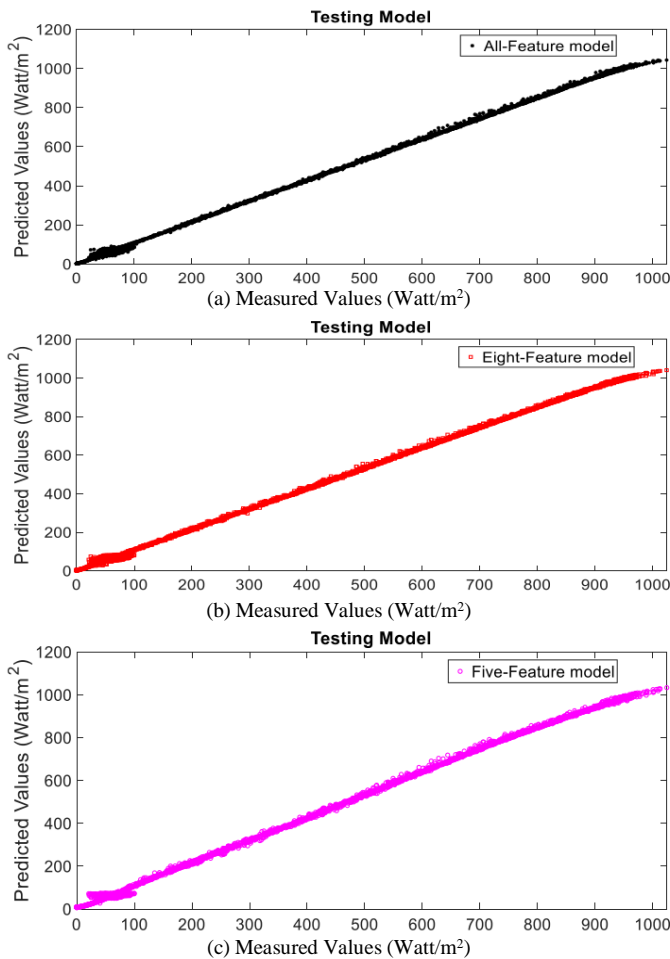


Fig. 7. Measured vs. Predicted Values of GHI for the Jeddah Region, (a) All-Feature Model, (b) Eight-Feature Model, (c) Five-Feature Model.

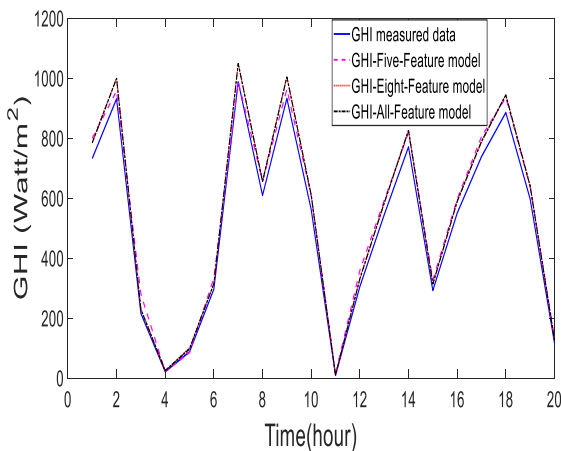


Fig. 8. The Forecasted GHI Values vs Measured GHI Values at Qassim.

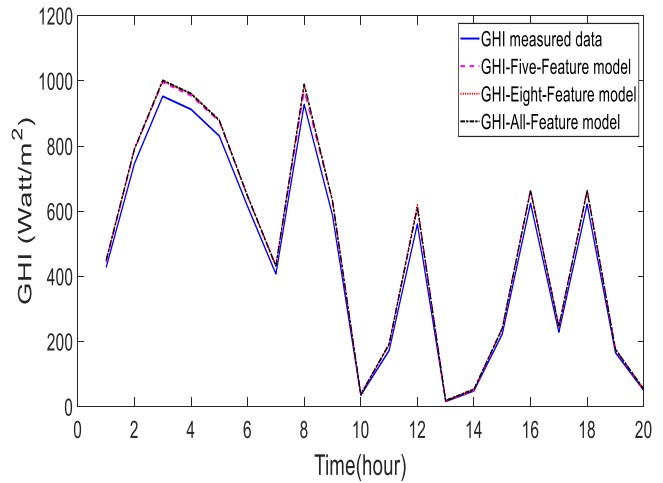


Fig. 9. The Forecasted GHI Values vs. Measured GHI Values at Jeddah.

The feature selection algorithm helped to decrease the cost of solar monitoring stations by reducing the number of features. For example, the less reduction in Jeddah for the Five-Features model, the cost decrease from 113990 RS to 60026 RS. The reduction percentage is 47%. As well in Jeddah using 8 features, the cost reduced from 113990 RS to 62065 RS, where the reduction percentage was equal to 46%. The rate of reduced cost is big although, it does not include the costs of maintenance, cables, and other accessories prices, On the other hand, in Qassim region, the cost reduced when using five features from 113990 to 76886 with a percentage =33%. As well when using 8 features in this region, the cost reduced from 113990 RS to 109333 RS, with a reduction percentage equal only 4%, as shown in Table IV.

At the site of Jeddah, and according to the findings of the feature selection algorithm, using the five-feature model resulted in the best prediction performance of the GHI compared to the prediction values of the model used a larger number of eight or all features. Also, the rate of cost reduction in Jeddah was very high, as presented in Table IV, where the cost reduction values for All, eight, and five attributes for that model were calculated.

On the other hand, in the Qassim area, the best prediction values were obtained by using the model with all-feature although the model with few attributes gives good results, where the value  $R^2$  gives approximately the same values for All-Features. This means the eight and five features also satisfying good results can apply in the future with a lower cost than the cost of all features, although the total costs reduced were small value it remains positive. Consequently, the feature selection algorithm helps to decrease the cost of solar monitoring stations. The reduced costs do not include specialists, maintenance, cables, and accessories prices. This, in turn, gives special importance to the research finding.

TABLE IV. COST REDUCTION FOR THE GENERATED MODEL IN THE TWO REGIONS IN SAUDI RIYALS (RS)

| I                   | Features                              | Unit price | All Features | Qassim Region  |               | Jeddah Region  |               |
|---------------------|---------------------------------------|------------|--------------|----------------|---------------|----------------|---------------|
|                     |                                       |            |              | Eight Features | Five Features | Eight Features | Five Features |
| 1                   | Month Of The Year(M)                  | -          | √            | Eight          | Five          | √              | -             |
| 2                   | Day Of The Month(D)                   | -          | √            | -              | -             | -              | -             |
| 3                   | Hour Of The Day (H)                   | -          | √            | √              | √             | √              | √             |
| 4                   | Air Temperature (T)                   | 730        | √            | √              | √             | √              | -             |
| 5                   | Relative Humidity (RH)                | 730        | √            | √              | -             | -              | -             |
| 6                   | Surface Pressure (P)                  | 3013       | √            | -              | -             | -              | -             |
| 7                   | Wind Speed At 3 Meters (WS)           | 1309       | √            | √              | -             | √              | -             |
| 8                   | Wind Direction (WD)                   | 1309       | √            | -              | -             | √              | √             |
| 9                   | Peak Wind Direction At 3 Meters (PWD) | 1309       | √            | -              | -             | √              | √             |
| 10                  | Diffuse Horizontal Irradiance(DHI)    | 35037      | √            | -              | -             | -              | -             |
| 11                  | Direct Normal Irradiance (DNI)        | 13145      | √            | √              | √             | -              | -             |
| 12                  | Azimuth Angle (AA)                    | 28704      | √            | √              | √             | √              | √             |
| 13                  | Solar Zenith Angle (SZA)              | 28704      | √            | √              | √             | √              | √             |
| Total               |                                       | -          | 113990       | 109333         | 76886         | 62065          | 60026         |
| Cost Reduction Rate |                                       | -          | -            | 4%             | 33%           | 46%            | 47%           |

## VI. CONCLUSIONS

The use of an advanced embedded feature selection algorithm and ANN is addressed in this paper to forecast the hourly solar radiation at two sites in the Kingdom of Saudi Arabia. The data from two stations; Qassim and Jeddah in Saudi Arabia, it was obtained to examine the prediction performance of the developed model. The five and eight most important variables among a wide range of metrological variables that could impact solar radiation in the future were optimally and systematically determined by employing a recent feature selection technique named as Boruta algorithm. For the comparison reasons of the model results with different features. The all-feature model was used to assess the benefits of using a feature selection method. The 13 input variables are the maximum number of features considered for developing a data-driven forecasting model.

At the site of Jeddah, and according to the findings of the feature selection algorithm, using the five-feature model resulted in the best prediction performance compared to the prediction values of the model when used a larger number of eight or all features. Also, the rate of cost reduction in Jeddah was very high. In the Qassim area, the best prediction values were obtained by using the all-feature model although the model of other features with few attributes is good where the value  $R^2$  gives approximately the same values for All-Features. This means the 8 and 5 features also satisfying very good results can apply in the future with a lower cost than the cost of all features, but, it noted that the total costs reduced were small value. Using feature selection methods may successfully exploit the larger interdependent variables relevant to hourly global horizontal irradiance prediction without sacrificing predictive efficiency. The findings, therefore, emphasize the importance of using feature selection techniques when using the model for computational intelligence to achieve accurate predictions of solar radiation. On the other hand, the

cost rate of the GHI prediction was reduced for the generated model, as presented in Table IV above, where the cost reduction values of the model using all, eight and five attributes were calculated. Consequently, From the discussed results, it found that the feature selection algorithm helps to decrease the cost of instruments and equipment required for solar monitoring stations for a high rate. Although, the costs that were reduced do not include the cost of specialists, maintenance, cables, and accessories prices. This, in turn, gives strength and special importance to the research finding. Besides, the lower-cost models can be used in future to collect new data for the coming years for forecasting.

## VII. FUTURE WORK

The research results are leading all researchers who have the same interest to achieve important extensions in the future to the current finding, it can include the following:

- 1) Expanding the samples of the study, using other ML tools, and going deeper into the data analysis based on other selection features methods.
- 2) Studying more locations across Saudi Arabia to address the geographical effects.
- 3) Investigating more ML algorithms and comparing prediction performances.
- 4) Considering different types of feature selection methods.
- 5) Going more into the dataset analysis deeply to find other important insights that can also help in KSA community services in the future.
- 6) Considering different test locations with different climate conditions, to investigate their effects on the performance of the prediction model enhanced by feature selection techniques can form another research potential in the scope of the solar prediction.

#### ACKNOWLEDGMENT

The authors would like to express their great thanks to King Abdullah City for Atomic and Renewable Energy for providing the required datasets for this research.

#### REFERENCES

- [1] D. Lew et al., "Sub-Hourly Impacts of High Solar Penetrations in the Western United States," 2012.
- [2] M. Sandhu and T. Thakur, "Issues, Challenges, Causes, Impacts and Utilization of Renewable Energy Sources - Grid Integration," *J. Eng. Res. Appl.*, vol. 4, no. 3, pp. 636–643, 2014, [Online]. Available: [http://www.ijera.com/papers/Vol4\\_issue3/Version1/DH4301636643.pdf](http://www.ijera.com/papers/Vol4_issue3/Version1/DH4301636643.pdf).
- [3] B. . Hernandez, "The Religiosity and spirituality scale for youth : The Developmeny and initial validation," *Anesthesiology*, vol. 115, no. 3, p. A13, 2011, doi: 10.1097/ALN.0b013e3182318466.
- [4] A. Tuohy et al., "Solar Forecasting: Methods, Challenges, and Performance," *IEEE Power Energy Mag.*, vol. 13, no. 6, pp. 50–59, 2015, doi: 10.1109/MPE.2015.2461351.
- [5] M. Hossain, S. Mekhilef, M. Danesh, L. Olatomiwa, and S. Shamsirband, "Application of extreme learning machine for short term output power forecasting of three grid-connected PV systems," *J. Clean. Prod.*, vol. 167, pp. 395–405, 2017, doi: 10.1016/j.jclepro.2017.08.081.
- [6] L. Martín, L. F. Zarzalejo, J. Polo, A. Navarro, R. Marchante, and M. Cony, "Prediction of global solar irradiance based on time series analysis: Application to solar thermal power plants energy production planning," *Sol. Energy*, vol. 84, no. 10, pp. 1772–1781, 2010, doi: 10.1016/j.solener.2010.07.002.
- [7] A. Alzahrani, J. W. Kimball, and C. Dagli, "Predicting solar irradiance using time series neural networks," in *Procedia Computer Science*, 2014, vol. 36, pp. 623–628, doi: 10.1016/j.procs.2014.09.065.
- [8] A. Sharma and A. Kakkar, "Forecasting daily global solar irradiance generation using machine learning," *Renewable and Sustainable Energy Reviews*, Elsevier Ltd, vol. 82, no. 5, pp. 2254–2269, 2018, doi: 10.1016/j.rser.2017.08.066.
- [9] I. A. Ibrahim and T. Khatib, "A novel hybrid model for hourly global solar radiation prediction using random forests technique and firefly algorithm," *Energy Convers. Manag.*, vol. 138, pp. 413–425, 2017, doi: 10.1016/j.enconman.2017.02.006.
- [10] S. Ghimire, R. C. Deo, N. J. Downs, and N. Raj, "Self-adaptive differential evolutionary extreme learning machines for long-term solar radiation prediction with remotely-sensed MODIS satellite and Reanalysis atmospheric products in solar-rich cities," *Remote Sens. Environ.*, vol. 212, pp. 176–198, 2018, doi: 10.1016/j.rse.2018.05.003.
- [11] R. Kumar, R. K. Aggarwal, and J. D. Sharma, "Comparison of regression and artificial neural network models for estimation of global solar radiations," *Renewable and Sustainable Energy Reviews*, Elsevier Ltd, vol. 52, pp. 1294–1299, 2015, doi: 10.1016/j.rser.2015.08.021.
- [12] W.-Y. Chang, "Short-Term Load Forecasting Using Radial Basis Function Neural Network," *J. Comput. Commun.*, vol. 03, no. 11, pp. 40–45, 2015, doi: 10.4236/jcc.2015.311007.
- [13] C. Renno, F. Petit, and A. Gatto, "Artificial neural network models for predicting the solar radiation as input of a concentrating photovoltaic system," *Energy Convers. Manag.*, vol. 106, pp. 999–1012, 2015, doi: 10.1016/j.enconman.2015.10.033.
- [14] M. H. Alobaidi, P. R. Marpu, T. B. M. J. Ouarda, and H. Ghedira, "Mapping of the solar irradiance in the UAE using advanced artificial neural network ensemble," *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, vol. 7, no. 8, pp. 3668–3680, 2014, doi: 10.1109/JSTARS.2014.2331255.
- [15] A. Mellit, A. H. Arab, N. Khorissi, and H. Salhi, "An ANFIS-based forecasting for solar radiation data from sunshine duration and ambient temperature," 2007, doi: 10.1109/PES.2007.386131.
- [16] J. C. Lam, K. K. W. Wan, and L. Yang, "Solar radiation modelling using ANNs for different climates in China," *Energy Convers. Manag.*, vol. 49, no. 5, pp. 1080–1090, 2008, doi: 10.1016/j.enconman.2007.09.021.
- [17] C. Voyant, M. Muselli, C. Paoli, and M. L. Nivet, "Numerical weather prediction (NWP) and hybrid ARMA/ANN model to predict global radiation," *Energy*, vol. 39, no. 1, pp. 341–355, 2012, doi: 10.1016/j.energy.2012.01.006.
- [18] P. a Gutu, "Hybrid Artificial Neural Networks : Models ," vol. 9, no. 2, pp. 177–184, 2011.
- [19] L. Olatomiwa, S. Mekhilef, S. Shamsirband, K. Mohammadi, D. Petković, and C. Sudheer, "A support vector machine-firefly algorithm-based model for global solar radiation prediction," *Sol. Energy*, vol. 115, pp. 632–644, 2015, doi: 10.1016/j.solener.2015.03.015.
- [20] R. C. Deo and M. Şahin, "An extreme learning machine model for the simulation of monthly mean streamflow water level in eastern Queensland," *Environ. Monit. Assess.*, vol. 188, no. 2, pp. 1–24, Feb. 2016, doi: 10.1007/s10661-016-5094-9.
- [21] M. Mohandes, "Support vector machines for short-term electrical load forecasting," *Int. J. Energy Res.*, vol. 26, no. 4, pp. 335–345, 2002, doi: 10.1002/er.787.
- [22] L. M. Saini, S. K. Aggarwal, and A. Kumar, "Parameter optimisation using genetic algorithm for support vector machine-based price-forecasting model in National electricity market," *IET Gener. Transm. Distrib.*, vol. 4, no. 1, pp. 36–49, 2010, doi: 10.1049/iet-gtd.2008.0584.
- [23] A. Hepbasli and Z. Alsuhaibani, "A key review on present status and future directions of solar energy studies and applications in Saudi Arabia," *Renewable and Sustainable Energy Reviews*, vol. 15, no. 9, pp. 5021–5050, 2011, doi: 10.1016/j.rser.2011.07.052.
- [24] A. A. El-Sebaei, A. A. Al-Ghamdi, F. S. Al-Hazmi, and A. S. Faidah, "Estimation of global solar radiation on horizontal surfaces in Jeddah, Saudi Arabia," *Energy Policy*, vol. 37, no. 9, pp. 3645–3649, 2009, doi: 10.1016/j.enpol.2009.04.038.
- [25] J. C. Lam, K. K. W. Wan, and L. Yang, "Solar radiation modelling using ANNs for different climates in China," *Energy Convers. Manag.*, vol. 49, no. 5, pp. 1080–1090, 2008, doi: 10.1016/j.enconman.2007.09.021.
- [26] H. Bulut and O. Büyükalaca, "Simple model for the generation of daily global solar-radiation data in Turkey," *Appl. Energy*, vol. 84, no. 5, pp. 477–491, 2007, doi: 10.1016/j.apenergy.2006.10.003.
- [27] H. A. Kazem, "Solar Radiation, Temperature and Humidity Measurements in Sohar-Oman," *Int. J. Comput. Appl. Sci.*, vol. 1, no. 1, pp. 15–20, 2016, doi: 10.24842/1611/0003.
- [28] S. Rehman and S. G. Ghorri, "Spatial estimation of global solar radiation using geostatistics," *Renew. Energy*, vol. 21, no. 3–4, pp. 583–605, 2000, doi: 10.1016/S0960-1481(00)00078-1.
- [29] A. A. El-Sebaei, A. A. Al-Ghamdi, F. S. Al-Hazmi, and A. S. Faidah, "Estimation of global solar radiation on horizontal surfaces in Jeddah, Saudi Arabia," *Energy Policy*, vol. 37, no. 9, pp. 3645–3649, 2009, doi: 10.1016/j.enpol.2009.04.038.
- [30] M. Mohandes, S. Rehman, and T. O. Halawani, "Estimation of global solar radiation using artificial neural networks," *Renew. Energy*, vol. 14, no. 1–4, pp. 179–184, 1998, doi: 10.1016/S0960-1481(98)00065-2.
- [31] M. Benganem, A. Mellit, and S. N. Alamri, "ANN-based modelling and estimation of daily global solar radiation data: A case study," *Energy Convers. Manag.*, vol. 50, no. 7, pp. 1644–1655, Jul. 2009, doi: 10.1016/j.enconman.2009.03.035.
- [32] M. Almarashi, "Short-term prediction of solar energy in Saudi Arabia using automated-design fuzzy logic systems," *PLoS One*, vol. 12, no. 8, pp. 1–16, doi: 10.1371/journal.pone.0182429.
- [33] B. Ameen, H. Balzter, C. Jarvis, and J. Wheeler, "Modelling hourly global horizontal irradiance from satellite-derived datasets and climate variables as new inputs with artificial neural networks," *Energies*, vol. 12, no. 1, 2019, doi: 10.3390/en12010148.
- [34] P. Mpfumali, C. Sigauke, A. Bere, and S. Mulaudzi, "Day ahead hourly global horizontal irradiance forecasting—application to South African data," *Energies*, vol. 12, no. 18, pp. 1–28, 2019, doi: 10.3390/en12183569.
- [35] T. Betti, I. Zulim, S. Brkić, and B. Tuka, "A Comparison of Models for Estimating Solar Radiation from Sunshine Duration in Croatia," *Int. J. Photoenergy*, vol. 2020, 2020, doi: 10.1155/2020/9605950.
- [36] V. Kallio-Myers, A. Riihelä, P. Lahtinen, and A. Lindfors, "Global horizontal irradiance forecast for Finland based on geostationary weather satellite data," *Sol. Energy*, vol. 198, no. 1, pp. 68–80, 2020, doi: 10.1016/j.solener.2020.01.008.

- [37] Y. H. Koo, M. Oh, S. M. Kim, and H. D. Park, "Estimation and mapping of solar irradiance for Korea by using COMS MI satellite images and an artificial neural network model," *Energies*, vol. 13, no. 2, 2020, doi: 10.3390/en13020301.
- [38] J. Fan, X. Wang, F. Zhang, X. Ma, and L. Wu, "Predicting daily diffuse horizontal solar radiation in various climatic regions of China using support vector machine and tree-based soft computing models with local and extrinsic climatic data," *J. Clean. Prod.*, vol. 248, p. 1-14, 2020, doi: 10.1016/j.jclepro.2019.119264.
- [39] S. Salcedo-Sanz, S. Jiménez-Fernández, A. Aybar-Ruiz, C. Casanova-Mateo, J. Sanz-Justo, and R. García-Herrera, "A CRO-species optimization scheme for robust global solar radiation statistical downscaling," *Renew. Energy*, vol. 111, pp. 63-76, 2017, doi: 10.1016/j.renene.2017.03.079.
- [40] A. K. Yadav, H. Malik, and S. S. Chandel, "Application of rapid miner in ANN based prediction of solar radiation for assessment of solar energy resource potential of 76 sites in Northwestern India," *Renewable and Sustainable Energy Reviews*, Elsevier Ltd, vol. 52, pp. 1093-1106, Aug. 2015, doi: 10.1016/j.rser.2015.07.156.
- [41] A. R. Hedar, A. E. Abdel-Hakim, and M. Almarashi, "Granular-based dimension reduction for solar radiation prediction using adaptive memory programming," in *GECCO 2016 Companion - Proceedings of the 2016 Genetic and Evolutionary Computation Conference*, Jul. 2016, pp. 929-936, doi: 10.1145/2908961.2931648.
- [42] P. Jeatrakul, K. W. Wong, and C. C. Fung, "Data cleaning for classification using misclassification analysis," *J. Adv. Comput. Intell. Informatics*, vol. 14, no. 3, pp. 297-302, 2010, doi: 10.20965/jaciii.2010.p0297.
- [43] M. Mohandes, "Support vector machines for short-term electrical load forecasting," *Int. J. Energy Res.*, vol. 26, no. 4, pp. 335-345, 2002, doi: 10.1002/er.787.
- [44] S. Leva, A. Dolara, F. Grimaccia, M. Mussetta, and E. Ogliari, "Analysis and validation of 24 hours ahead neural network forecasting of photovoltaic output power," *Math. Comput. Simul.*, vol. 131, pp. 88-100, 2017, doi: 10.1016/j.matcom.2015.05.010.
- [45] R. Kohavi and G. H. John, "Wrappers for feature subset selection," *Artif. Intell.*, vol. 97, no. 1-2, pp. 273-324, 1997, doi: 10.1016/s0004-3702(97)00043-x.
- [46] R. Nilsson, J. M. Peña, J. Björkegren, and J. Tegnér, "Consistent feature selection for pattern recognition in polynomial time," *J. Mach. Learn. Res.*, vol. 8, pp. 589-612, 2007.
- [47] I. Guyon and A. Elisseeff, "An introduction to variable and feature selection," *Journal of Machine Learning Research*, vol. 3, pp. 1157-1182, 2003, doi: 10.1162/153244303322753616.
- [48] R. C. Team, "The R Project for Statistical Computing," [Http://www.R-Project.Org/](http://www.R-Project.Org/), pp. 1-12, 2013, [Online]. Available: <https://www.r-project.org/>.
- [49] A. Ng and K. Soo, "0.1 Die Weisheit der Crowd," *Data Sci. ist das Eig.*, vol. 45, pp. 5-32, 2018, doi: 10.1007/978-3-662-56776-0\_10.
- [50] H. Stoppiglia, G. Dreyfus, R. Dubois, and Y. Oussar, "Ranking a random feature for variable and feature selection," *J. Mach. Learn. Res.*, vol. 3, pp. 1399-1414, 2003.
- [51] J. Hebebrand, "Editorial: Contents of this issue," *Obesity Facts*, vol. 3, no. 6, pp. 343-344, 2010, doi: 10.1159/000323281.
- [52] M. INUIGUCHI, "Rough Sets and Current Trends in Computing 2004," *Syst. Control Inf.*, vol. 48, no. 11, p. 473, 2004, doi: 10.11509/isciesci.48.11\_473.

# Towards Computational Models to Theme Analysis in Literature

Abdulfattah Omar  
Department of English  
College of Sciences and Humanities  
Prince Sattam Bin Abdulaziz University

**Abstract**—The recent years have witnessed the development of numerous computational methods that have been widely used in humanities and literary studies. In spite of their potentials of such methods in terms of providing workable solutions to different inherent problems within these domains including selectivity, objectivity, and replicability, very little has been done on thematic studies in literature. Almost all the work is done through traditional methods based on individual researchers' reading of texts and intuitive abstraction of generalizations from that reading. These approaches have negative implications to issues of objectivity and replicability. Furthermore, it is challenging for such traditional methods to deal effectively with the hundreds of thousands of new novels that are published every year. In the face of these problems, this study proposes an integrated computational model for the thematic classifications of literary texts based on lexical clustering methods. As an example, this study is based on a corpus including Thomas Hardy's novels and short stories. Computational semantic analysis based on the vector space model (VSM) representation of the lexical content of the texts is used. Results indicate that the selected texts were thematically grouped based on their semantic content. It can be claimed that text clustering approaches which have long been used in computational theory and data mining applications can be usefully used in literary studies.

**Keywords**—*Computational models; computational semantics; lexical clustering; lexical content; philological methods; Thomas Hardy; Vector Space Model (VSM)*

## I. INTRODUCTION

An important development in literary studies over the past few decades is the increasing application of scientific methods in analysis of literary works [1-5]. It has been argued that the use of such scientific methods can assist in preventing the formation of false theories of criticism and the generation of unreliable thematic classifications [6, 7]. The present study is intended as a contribution to that development. The study seeks to propose a computational model that helps readers and critics of literary texts in an objective, replicable, therefore scientific way through exploring the thematic relationships of texts in a conceptually coherent way. The study is based on the novels and short stories of Thomas Hardy as an example. Thomas Hardy is one of the most important figures in the history of the English novel and he has ever sustained readers' interest in the themes and topics he tackled. Thomas Hardy was a Victorian poet and novelist and is considered by many critics as a main component of the English cultural heritage [8-10].

In spite of the proliferation of computational technology and the articulation of an explosive production of electronically encoded information of all kinds, computational methods have been very little used in humanities in general and literature in particular [11-13]. The wide cultural gap between the literary critic and computational research communities is the most obvious reason. The study is an attempt towards bridging the gap between traditional literary criticism and computational methods. The study employs experimentally replicable data representation and clustering methods. The greatest advantage of these methods is that they are completely objective in the sense that the results obtained are independent of the person applying the method.

The remainder of this article is organized as follows. Section 2 is a brief survey of the approaches to thematic studies of literary works. Section 3 outlines the methods and procedures of the study. It describes document clustering methods are used to classify the selected works in a thematically coherent way. Section 4 reports the results of the proposed methods and explores the thematic interrelationships between the texts. Section 5 is conclusion. It summarizes the main findings and suggests propositions that may be generalized to other literary texts and genres.

## II. LITERATURE REVIEW

Theme analysis of literary texts is one of the oldest and most established disciplines in literary studies. Critics have been generally concerned with identifying the themes within literary texts. It was thought that part of the critic's job is to understand the deep meanings conveyed by authors, make observations about literary texts in order to construct the expression of themes in these works [14-16].

Although theme analysis is very old in literature studies, the issue of the way themes are defined is still controversial in literary criticism. There is no single agreed upon approach to theme analysis in literature. Over the years, there is no consensus among critics on the best ways of interpreting texts and deriving thematic concepts. It is true to claim that theme analysis is still controversial and problematic in literary criticism studies [17].

With the development of different literary theories including Marxism, Modernism, and feminism, theme analysis has been widely considered a reflection of these theories [18, 19]. Critics have been more concerned with identifying the relationship between author and work as reflected on themes of



race, class, and gender. In this way, thematic concepts of novelists and authors are usually confined and restricted to the critic's engagement with a given theory or selections from a text or some texts.

The issue of theme analysis with its complexities and controversies has its implications to the thematic studies of Thomas Hardy's literary texts. The thematic classification of Hardy's prose fiction ranges from a broad general classification of his novels and short stories to a discussion of a single thematic aspect in one, some, or all his writings [10, 20-25]. The main observation about almost the critical studies on the thematic structures of Hardy's work is that critics have been generally concerned with what Hardy himself classified as Major Works. Despite the rich thematic concepts exhibited in Hardy's prose fiction works exhibit rich thematic concepts, the majority of the thematic discussions of Hardy have been flawed in limiting their discussions to the series of novels and short stories he wrote between 1871 and 1895 [26].

It can be claimed thus that the work on Thomas Hardy's prose fiction is widely selective. Some critics focus on what is referred to Wessex novels. They think of Hardy's works as a cry for the lost beauty of the English countryside. Evidently, many commentators have characterized Hardy as a regional novelist, attribute this regionalist focus to his fascination with Wessex, an old English kingdom covering an area that provided the fictionalised setting for Hardy [27, 28]. Balanced against this argument, others insist that Hardy was a Victorian social critic since his writings depict the sufferings of England's working class and society's responsibility for their tragic fates. Through this process, Hardy is seen to have been preoccupied with improving conditions in society. These concerns mark Hardy out as a realistic writer who took on the role of expressing the joys and woes of the victims of the merciless harshness of their lives [25, 29].

One major problem with studies in this tradition is that they ignore much of the thematic richness in Hardy's works. In the face of this limitation, this study suggests the use of empirical approaches and new technologies. These should have the impact of developing a comprehensive and more detailed structuring of Hardy's thematic concepts.

### III. METHODS AND PROCEDURES

For developing a computational model for deriving taxonomies of thematic concepts in literary texts, document clustering theory is adopted. Document clustering theory has been widely used in data mining and information retrieval (IR) applications [30, 31]. Document clustering methods are generally used for grouping similar texts together [32, 33]. The hypothesis is that texts grouped together are more likely to have the same theme [34-36]. Document clustering methods have been proved effective in grouping and categorizing unstructured text data and exploring. In such processes, similar texts are separated together in distinct groups or clusters. Accordingly, document clustering methods can be usefully used in the domains of theme analysis in literary studies.

There are numerous document clustering methods. For the purposes of the study, vector space clustering is used. VSC is one of the earliest computer-based clustering methods [33, 37].

It is thought however that it is appropriate for the study. The rationale is that the study is concerned with build thematic structures of the texts based on their lexical semantics. VSC is thus appropriate for the purposes of the study. In VSC, documents can be grouped into distinct classes based on their lexical content [30, 38, 39]. In this regard, it is assumed that VSC is appropriate for the purposes of the study. VSC is used for organizing the novels and short stories of Thomas Hardy into distinct classes based on the lexical content of these texts. Herein, lexical clustering (one of VSC methods) is used.

Conventional lexical-clustering algorithms treat text fragments as a mixed collection of words, with a semantic similarity between them calculated based on the term of how many the particular word occurs within the compared fragments. Whereas this technique is appropriate for clustering large-sized textual collections, it operates poorly when clustering small-sized texts such as sentences [40].

In so doing, a corpus including all the selected texts is designed. The tradition of building a corpus for text clustering applications has always been based on the assumption that the corpus is both large and representative of the research domain. An important question in the context of this study is what size the corpus should be in order to support objective and reliable generalizations about Thomas Hardy's prose fiction. The corpus on which this analysis is based consists of all the known (published and unpublished) prose fiction texts of Hardy.

As a first step for data representation, the corpus was confined to what is referred to a bag of words. It was also decided that the corpus to be built of only the content words. All function words were thus removed. The hypothesis is that they do not usually carry semantic meaning; thus, they cannot be considered as distinctive features. The corpus should include only and all the distinctive features [41]. Content words can act as strong predictors of the topic(s) or content of a document [42]. Moreover, the experimental results of document classification indicate that content-word representation gives good results in identifying the content of a document and its latent structure [43-45]. Equally important, most studies seem to agree that content-word representation has been proved to give much better results than any other approach to clustering. This study considers content words to be indicators of semantic content. In other words, the analysis identifies all the morphological variants of a given stem as just one lexical type. It can be observed that variant word forms with similar semantic conceptions can be treated as equivalent. To take an example, the words 'marry', 'marries', 'married', and 'marriage' deal with a single semantic concept, which is necessarily different from, for example, *dogs* and *cats*. The analysis thus reduces all these variant forms to just one form, i.e. *marry*.

In order for the texts to be amenable for computational analysis, texts were mathematically represented using vector space model (VSM). The reason is that it is conceptually simple as well as it is convenient for computing semantic similarity within documents.

A data matrix was created including all the 62 selected texts and lexical types (45,298 variables) included in this study.

An initial observation about the corpus is that the 62 texts vary substantially in size, ranging from 002 Kb to 389 Kb. One major problem with a corpus of the kind is that documents will be clustered based on size rather than lexical content and semantic similarity. The row vectors of were thus normalized to compensate for variations in length. Mean document length method was used for the purpose. This had the effect that the documents were equally represented in the matrix and their lexical frequency profiles could be meaningfully clustered.

In order to extract the most distinctive lexical variables, term frequency inverse document frequency (TFIDF) was used. Based on the TFIDF of the Hardy matrix, the highest 200 variables were decided to be the most distinctive lexical features within the corpus. It was also clear that the texts are best categorized into four distinct classes as seen in Fig. 1.

Cluster analysis was then used to find meaningful clusters in the data. Cluster analysis was used to generate a centroid-based lexical clustering structure that captures and describes the semantic similarities of the selected texts in the data matrix based on the lexical resource. The hypothesis is that texts grouped together should have common thematic features and based on the lexical semantic properties of the variables of these texts; it is easy to assume the recurrent themes in these

texts. For visualization of the clustering structure, hierarchical cluster analysis is used. Hierarchical cluster analysis is one of the main statistical approaches that is used for finding distinct classes or groups based on the shared and common features. Despite the development of different clustering algorithms, hierarchical cluster analysis or hierarchical clustering remains one of the most widely unsupervised clustering algorithms in clustering applications to find discrete groups with varying degrees of (dis)similarity in a data set represented by a (dis)similarity matrix [46]. The selected texts fall into four main clusters as shown in Fig. 2.

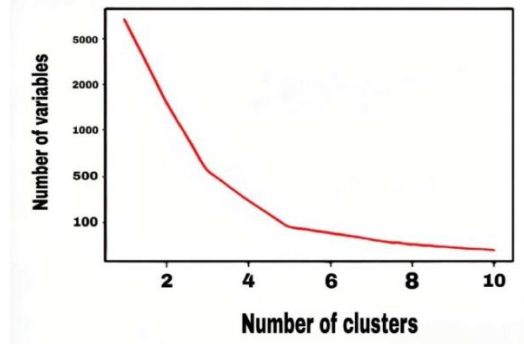


Fig. 1. A TFIDF Analysis of the Matrix H62, 200.

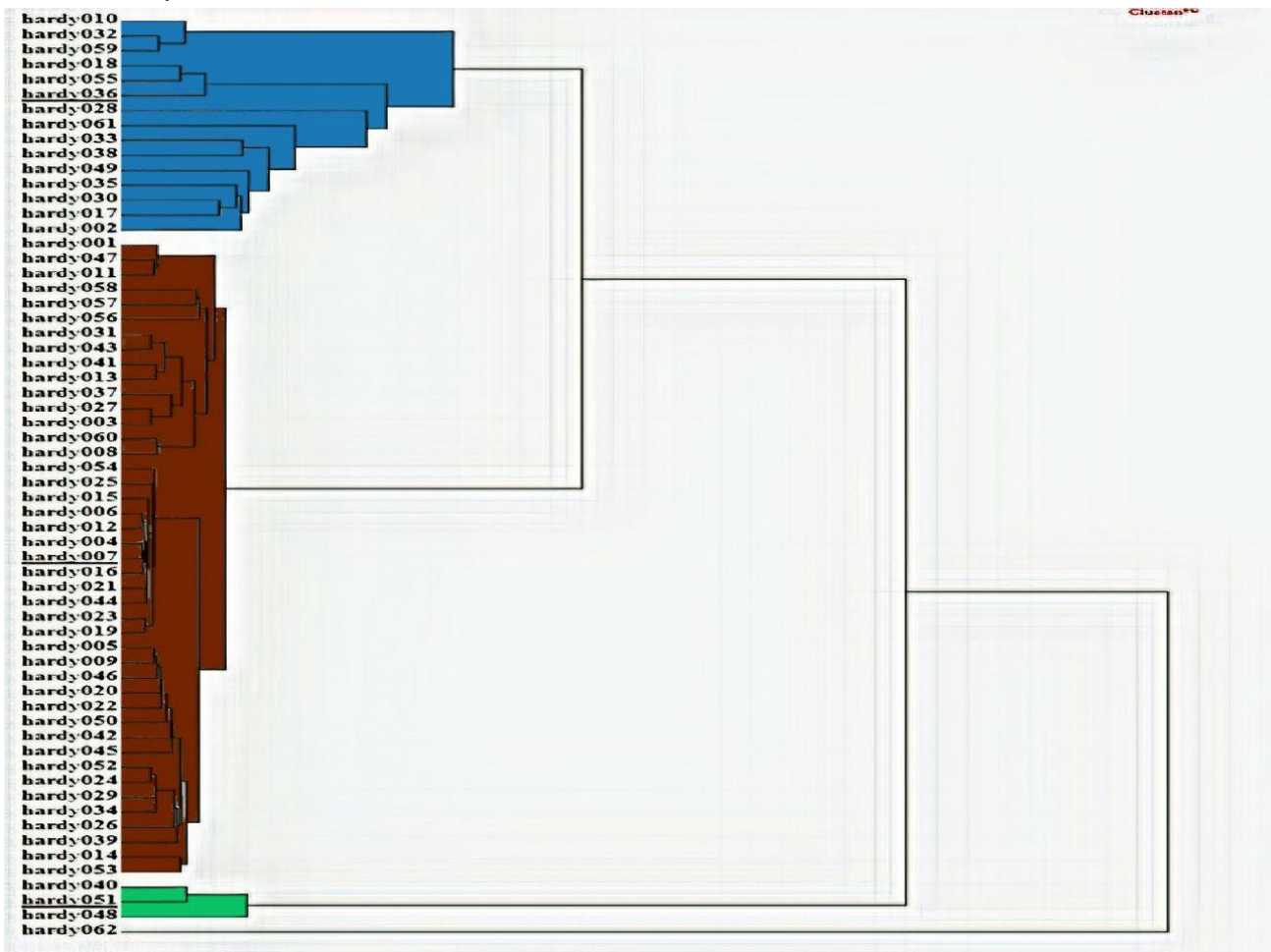


Fig. 2. A Clustering Structure of Hardy's Matrix using Euclidean Distance and Ward Linkage Clustering.

The result is a centroid-based lexical-clustering structure that can be used in any application in which the relationship between patterns is expressed in terms of pairwise semantic similarities [40]. In our case, this clustering structure is used for building hypotheses and making generalizations about the thematic relations of the texts in each group or cluster.

For validity purposes, PCA was used. The validity of cluster analysis results is an important requirement since different cluster structures may lead to completely different interpretations of the same data and thus generating contradictory hypotheses about the data. The purpose of clustering validation here thus is to see whether the same analytical methods applied to an alternative representation of the data gives identical or at least similar results. The alternative data representation was generated by principal component analysis, a dimensionality reduction method whereby H62, 200 was reduced to a dimensionality of 61, yielding the matrix H62, 50 [47, 48]. The result is that there is a full agreement between the results of the clustering structures.

In order to identify the most distinctive lexical features of each group, the columns of the matrix were rearranged in order of descending variance. Centroid vectors for the clusters A, B, C, and D were constructed by taking the means of the vectors in the matrix that constitute Groups A, B C, and D in accordance with the function

$$V_i = \frac{\sum_{i=1 \dots m} H_{ij}}{m}$$

Where

$V_j$  is the  $j$ th element of the centroid vector (for  $j = 1 \dots$  the number of columns in H),

H is the data matrix, and

$m$  is the number of row vectors in the cluster in question

The resulting vectors Group A<sub>centroid</sub>, Group B<sub>centroid</sub>, Group C<sub>centroid</sub> and Group D<sub>centroid</sub> were compared to show how, on average, the three groups differ on each of the extracted lexical variables, the aim being to identify the variables on which they differ most and thereby the thematic characteristics of each group can be inferred.

#### IV. ANALYSIS AND DISCUSSIONS

The aim in this section is to see whether the clustering structures thus far validated are meaningful. Given that the texts were clustered on the basis of lexical frequency vectors, this implies that each cluster has a characteristic lexical frequency profile which distinguishes it from the others [49]. By doing so, it should be possible to identify the most important variables for each group, and, on the basis of the lexical semantics of these items, to infer thematic characteristics of the respective groups [50, 51].

According to the computation of the quantitative findings and an intuitive understanding of the texts, each of these groups displays the distinctive lexical variables that make them thematically distinct. The frequent use of words like ‘duke’, ‘baron’, ‘duchess’, ‘knight’, ‘estate’, and ‘squire’ in Group A is

a good indication that this group is particularly concerned with aristocratic life and class differences. It can be suggested that this group touches on many aspects of class difference, adventure, romance, matrimony and mismatched unions and the conflicts they bring. Parallel to these themes, the Napoleonic era appears as a recurring theme in many of the texts of this group, as reflected by the frequent use of words like ‘Napoleon’, ‘France’, ‘French’, and ‘war’. This quantitative finding is supplemented by an intuitive reading of the texts and is also supported by critical assessments of the texts included here. Gilmartin and Mengham [52] argue that The Poor Man and the Lady and A Group of Noble Dames feature one of the most recurrent themes in Hardy’s books: that of cross-class relationships or marriages. The texts included here discuss issues of elopement, failure in marriage, and illegitimate children. This finding also agrees in principle with Hardy’s classification of his own works since the majority of the texts included here he classified under the category of Romance and Fantasies.

The hierarchical clustering structure, which is based on pure mathematical methods, supports Hardy’s tendency to group similar short stories together. Six texts of this group are included in his volume of short stories, A Group of Noble Dames. It also includes The Doctor’s Legend, which was first collected in Noble Dames when it was published in serial form in Harper’s Weekly and the Graphic in late 1890 [53]. Purdy [54] comments that the text appeared later in the collection A Group of Noble Dames under the title Barbara, which is thematically similar to The Legend. As such, the results of this analysis agree with the thematic structure that Hardy defined for his books.

Although the texts included here can be placed under the heading of Romance and Fantasies, as Hardy classified them, the element of social criticism persists through almost all of the texts. In The Poor Man and the Lady and the stories of A Group of Noble Dames, discussion of social problems is clear. Hardy is concerned with the problem of mismatched unions in a very class-conscious society. This argument is supported by Brady [55], who writes: In its subject matter, however, A Group of Noble Dames has interesting links with Hardy’s earlier work. The book is one of his many attempts, beginning with The Poor Man and the Lady, to portray the fascination and the difficulty of sexual alliances that cross class boundaries [55]. The texts involved in this group highlight the historical development of Hardy as a novelist and it is clear that Hardy was preoccupied with social issues throughout his career as a novelist and prose writer. This is supported by Dalziel [56], who stresses the essential continuity of Hardy’s thinking on social issues from the beginning to the end of his career as a writer.

The majority of the texts in this group as a whole are thus thematically related around romance and adventure. This does not contradict, however, the inclusion of texts like A Tradition of Eighteen Hundred and Four, Anna, or A Committee Man of “The Terror”, which are all about the political upheavals that took place in England and France as a result of the French Revolution and English Civil War—the main thematic frame in the first story is adventure while in the other two stories it is romance. Gilmartin and Mengham [52] argue that in spite of

the fact that the story is concerned with the theme of English-French conflicts: "it exhibits many of the expected features of a Christmas story (being written for the annual Harper's Christmas); it is meant to give a frisson of fear to those within the story who are sheltering from the rain and cold by the inn's fireside, and also to the readers of the periodical sitting by the Christmas hearth" [52].

The largest group, Group B, includes 43 texts out of the matrix's 62 rows, and is concerned with the English countryside; domestic life (as reflected in words such as 'river', 'cabbage', 'village', 'horse', 'mare', 'farmer', 'mill', 'tub', 'heath', 'cloth', 'sky', 'vicar', 'cover', 'passage', 'stream', 'hut', 'lane', and 'rain'); and struggle, outrage, and the frustrations of the poor ('public', 'money', 'children', 'work', 'fact', 'trade', 'bureau', and 'penny'). A common theme of contemporary social life can be suggested. Nevertheless, each subclass displays characteristic thematic features. One subclass which can be defined as Group B1, for instance, is tragedy, which correlates with ideas of social promotion/hostility and struggle. This subclass includes texts referred to by many critics as Hardy's major works. These include: *Far From the Madding Crowd*; *The Return of the Native*; *The Woodlanders*; *The Mayor of Casterbridge*; *Tess of the D'Urbervilles*; *Jude the Obscure*; and *The Trumpet Major*. The texts included here reflect Hardy's sense of disdain for the fashionable world and mock the social mores of the age. The texts talk generally about heroes and heroines who aspire for a better life and their attempts to achieve social promotion; as well as how they discover the falsity of their lives. They cannot escape the miserable conditions in which they live and are destined to suffer. Fate is an important factor in their suffering. The combination of social elements with these tragedies suggests the theme of social tragedy. Love is a recurrent theme in the other subcluster. The realistic representation, however, is always there. This is represented in texts such as *The Romantic Adventures of a Milkmaid* and *The Trumpet-Major*. Given that the texts represent different historical stages of Hardy's career, it may be claimed that the social element is heavily emphasized from the beginning of his career as a novelist up until he gave up writing novels. Unlike Hardy's classification of his own works, hierarchical cluster analysis along with qualitative analysis results point to social indicators influencing his career as a novelist. The social dimension is never absent in his writing.

There is also a correlation between the texts included here and Hardy's vision of Wessex and the English countryside. Many of the texts are set in that imaginary world of Wessex, the name of an Anglo-Saxon kingdom that covered a large area of south and southwest England prior to the Norman Conquest. It may also be claimed that woman and feminist issues are central themes in the texts of this group. Thomas Hardy was keen on describing Victorian hypocrisy in relation to women's issues. *Tess* highlights the rampant sexual assault and exploitation of the age. The novel also reflects Hardy's disapproval of the Victorians' obsession with female virginity. Fanny Robin in *Far from the Madding Crowd* is another example. When Troy refuses to marry her and abandons her, she tries to pick up her life as best she can. Finally, she becomes unable to work and is left without any money. As a

result, she and her child die of need and starvation. This offers another typical example of the suffering of women in the Victorian age.

Texts in Group C seem to form a distinct thematic relationship. The three short stories in this group, *What the Shepherd Saw* (Hardy048), *The Duke's Reappearance* (Hardy051), and *The Duchess of Hamptonshire* (Hardy040), are concerned with the idea of hidden or unrevealed death. This idea is repeated in the three texts where problems of jealousy and suspicion in marriage lead to death. The main idea of each of these three texts is that there is a beautiful married woman who belongs to the elite. Her husband, as a man of high position, feels jealous and decides to take revenge against the person who he thinks to be her lover, because of the disgrace such an illicit relationship causes him. Finally, Group D includes just one text which is *The Unconquerable* (Hardy062). The most important variables of this group are 'book', 'linger', 'occupation', 'measure', 'copying', 'bold', 'quaint', 'style', 'architecture', 'graveyard', 'figure', 'draughtsman', 'antique', 'masonry', and 'rose'. Correlating this cluster with the bibliographical data, it emerged that the text was written by Hardy in collaboration with his wife Florence Dugdale-Hardy. This can be an indication that it has unique lexical features that makes it distinct from other texts.

On the basis of the foregrounding discussions, it can be claimed that clustering structures are meaningful. Each cluster or Group has its distinctive lexical profile that distinguishes it from other groups or clusters. It may be claimed that cluster analysis points to significant facts regarding the novels and short narratives of Hardy. This cluster analysis relates some works to each other in ways not found in the established criticisms of Hardy. In Hardy's classification of his works, *The Return of the Native* is classified under the category of Novels of Environment and Character, while *The Hand of Ethelberta* is classified under the category of Novels of Ingenuity. However, here the two texts are clustered together in Group B. The dominant realistic approach of the works of Hardy may be one reason that many critics, who have attempted the thematic classification of Hardy's work, have not thought about connections and similarities between the two texts. In our case, we suggest that these two texts are related to each other in terms of their dealing with class consciousness. In his introduction to the New Wessex Edition of *The Hand of Ethelberta*, for example, Gittings [57] underestimates the novel, classifying it as 'the joker in the pack' of Hardy's novels. Widdowson [58], on the other hand, insists that *The Hand of Ethelberta* is not merely a romance, as Hardy classified it. He argues that the novel demonstrates Hardy's concern with the issue of class consciousness. He gives evidence that the text reflects bibliographical elements of Hardy's own life and draws parallels between the narrator of the story, who takes novel writing as a means for social promotion, and Thomas Hardy himself. Widdowson [58] comments that in all his novels, especially in *The Hand of Ethelberta*, Hardy appears concerned with the idea of class consciousness.

Equally important, the clustering structures provide ways of classifying the novels and short stories of Hardy according to genre. The idea is that thematic classification has pointed to

tragic, historic, and fantasy elements in the texts. Consequently, as far as thematic interrelationships are concerned, it is clear that the texts exhibit obvious features for genre classification. This can be a starting point for a comprehensive genre classification of the novels and short stories of Hardy and texts can be classified under the main categories of tragedy, comedy, romance, epic, fantasy, history, and pastoral histories, etc. One advantage of such a classification is that it can narrow down the ways in which we think about them. Here, I give some examples. Those texts that critics have usually considered tragedies are included in Group B—The Woodlanders, Tess, and Jude, for instance are all included in just one group. These are modern social tragedies and in these novels, Hardy deals with the social factors that determine the tragic end of his protagonists.

## V. CONCLUSION

This study addressed the question whether thematic concepts can be identified in literary texts using computational models. In this regard, document clustering methods were used for grouping the selected texts into distinct classes based on their semantic similarity. Results indicate clearly that document clustering methods can be usefully used for generating distinct and meaningful classes that express some thematic concepts such as class status, sex, marriage, love, romance, and the English countryside. The analytical results are objective in the sense that they are generated by mathematically based computational methods working on empirically derived data, and as such are not open to influences from any theoretical presuppositions that the researcher might have. Unlike results from the philological method, the computational results are replicable and therefore testable and scientifically respectable. This is not to overlook the subtle elaborations of literary criticism of Thomas Hardy over the past years. In point of fact, these elaborations generate a number of hypotheses which have not been empirically confirmed. Although the results of the present study largely agree with non-computational philological classifications of Hardy's texts, the contribution, however, is that the results obtained here are objective.

## ACKNOWLEDGMENTS

This project was supported by the Deanship of Scientific Research at Prince Sattam Bin Abdulaziz University under the research project No. 2020/02/11847.

## REFERENCES

- [1] C. Mullings, S. Kenna, M. Deegan, and S. Ross, *New Technologies for the Humanities*. De Gruyter, 2019.
- [2] G. Balossi, *A Corpus Linguistic Approach to Literary Language and Characterization: Virginia Woolf's The Waves*. John Benjamins Publishing Company, 2014.
- [3] I. Mani, Morgan, and C. Publishers, *Computational Modeling of Narrative*. Morgan & Claypool, 2013.
- [4] R. Siemens and S. Schreibman, *A Companion to Digital Literary Studies*. Wiley, 2013.
- [5] J. G. Shanahan, Y. Qu, and J. Wiebe, *Computing Attitude and Affect in Text: Theory and Applications*. Springer Netherlands, 2005.
- [6] M. L. Jockers and R. Thalken, *Text Analysis with R: For Students of Literature*. Springer International Publishing, 2020.
- [7] W. van Peer and S. Zyngier, *Directions in Empirical Literary Studies: In Honor of Willie Van Peer*. John Benjamins Publishing Company, 2008.
- [8] N. Page, *Oxford Reader's Companion to Hardy*. Oxford University Press, 2000.
- [9] M. Bevis, *The Oxford Handbook of Victorian Poetry*. OUP Oxford, 2013.
- [10] P. Mallett and S. E. Maier, *Thomas Hardy in Context*. Cambridge University Press, 2013.
- [11] J. Burrows, "Textual Analysis," in *A Companion to Digital Humanities*, S. Schreibman, R. Siemens, and J. Unsworth, Eds. Oxford: Blackwell, 2004, pp. 88-97.
- [12] M. K. Gold and L. F. Klein, *Debates in the Digital Humanities*. University of Minnesota Press, 2016.
- [13] D. L. Hoover, J. Culpeper, and K. O'Halloran, *Digital Literary Studies: Corpus Approaches to Poetry, Prose, and Drama*. Taylor & Francis, 2014.
- [14] T. Pugh and M. E. Johnson, *Literary Studies: A Practical Guide*. Taylor & Francis, 2013.
- [15] P. P. Headrick, *The Wiley Guide to Writing Essays About Literature*. Wiley, 2013.
- [16] M. L. Kamil, P. B. Mosenthal, P. D. Pearson, and R. Barr, *Handbook of Reading Research*. Taylor & Francis, 2014.
- [17] F. Mulhern, *Contemporary Marxist Literary Criticism*. Taylor & Francis, 2014.
- [18] R. Wellek and A. Warren, *Theory of Literature*. Dalkey Archive Press, 2020.
- [19] B. Kachuck, "Feminist Social Theories: Theme and Variations," *Sociological Bulletin*, vol. 44, no. 2, pp. 169-193, 1995.
- [20] J. L. Bownas, *Thomas Hardy and Empire: The Representation of Imperial Themes in the Work of Thomas Hardy*. Taylor & Francis, 2016.
- [21] R. G. Cox, *Thomas Hardy; the critical heritage (Critical heritage series)*. New York: Barnes & Noble, 1970, pp. xlvii, 473 p.
- [22] J. Dillion, *Thomas Hardy: Folklore and Resistance*. Palgrave Macmillan UK, 2016.
- [23] R. Nemesvari, *Thomas Hardy, Sensationalism, and the Melodramatic Mode*. New York: Palgrave Macmillan US, 2011.
- [24] P. Vigar, *The Novels of Thomas Hardy: Illusion and Reality*. Bloomsbury Publishing, 2014.
- [25] K. Wilson, *A Companion to Thomas Hardy*. Wiley, 2010.
- [26] K. Ireland, *Thomas Hardy, Time and Narrative: A Narratological Approach to his Novels*. Palgrave Macmillan UK, 2014.
- [27] P. Brantlinger and W. Thesing, *A Companion to the Victorian Novel*. Wiley, 2008.
- [28] J. Hodson, *Dialect and Literature in the Long Nineteenth Century*. Taylor & Francis, 2017.
- [29] J. King, "Thomas Hardy: Tragedy Ancient and Modern," in *Tragedy in the Victorian Novel* Cambridge: Cambridge University Press, 1978, pp. 97-126.
- [30] W. Wu, H. Xiong, and S. Shekhar, *Clustering and Information Retrieval*. Springer 2013.
- [31] A. N. Srivastava and M. Sahami, "Text Mining Classification, Clustering, and Applications," (*Data Mining and Knowledge Discovery Series*. Chapman and Hall, 2009, p. ^pp. Pages.
- [32] C. C. Aggarwal and C. K. Reddy, *Data Clustering: Algorithms and Applications*. CRC Press, 2016.
- [33] W. Wu, H. Xiong, and S. Shekhar, *Clustering and Information Retrieval*. Springer US, 2003.
- [34] F. M. G. França and A. F. de Souza, *Intelligent Text Categorization and Clustering*. Springer Berlin Heidelberg, 2008.
- [35] A. K. Somani, R. S. Shekhawat, A. Mundra, S. Srivastava, and V. K. Verma, *Smart Systems and IoT: Innovations in Computing: Proceeding of SSIC 2019*. Springer Singapore, 2019.
- [36] G. Chakraborty, M. Pagolu, and S. Garla, *Text Mining and Analysis: Practical Methods, Examples, and Case Studies Using SAS*. SAS Institute, 2014.
- [37] K. Riesen and H. Bunke, *Graph Classification And Clustering Based On Vector Space Embedding*. World Scientific Publishing Company, 2010.
- [38] J. Kogan, *Introduction to Clustering Large and High-Dimensional Data*. Cambridge: Cambridge University Press, 2007.

- [39] H. Moisl, Cluster Analysis for Corpus Linguistics. De Gruyter, 2015.
- [40] K. Abdalgader, "Centroid-Based Lexical Clustering," in Recent Applications in Data Clustering, H. Pirim, Ed.: IntechOpen, 2018, pp. 378-403.
- [41] D. Glynn and J. A. Robinson, Corpus Methods for Semantics: Quantitative studies in polysemy and synonymy. John Benjamins Publishing Company, 2014.
- [42] K. Luyckx, Scalability Issues in Authorship Attribution. UPA, 2011.
- [43] M. L. Eaton, Multivariate Statistics: A Vector Space Approach (Institute of Mathematical Statistics. Lecture notes-monograph series). Beachwood, Ohio: Institute of Mathematical Statistics, 2007, pp. viii, 512 p.
- [44] A. Gani, A. Siddiq, S. Shamshirband, and F. Hanum, "A survey on indexing techniques for big data: taxonomy and performance evaluation," Knowledge and information systems, vol. 46, no. 2, pp. 241-284, 2016.
- [45] T. Hofmann, "Probabilistic latent semantic indexing," in ACM SIGIR Forum, 2017, vol. 51, no. 2, pp. 211-218: ACM.
- [46] TomTullis and B. Albert, Measuring the User Experience: Collecting, Analyzing, and Presenting Usability Metrics, Second Edition ed. Elsevier 2013.
- [47] A. KASSAMBARA, Practical Guide To Principal Component Methods in R: PCA, M(CA), FAMD, MFA, HCPC, factoextra. CreateSpace Independent Publishing Platform, 2017.
- [48] H. Bozdogan and A. K. Gupta, Multivariate Statistical Modeling and Data Analysis: Proceedings of the Advanced Symposium on Multivariate Modeling and Data Analysis May 15–16, 1986. Springer Netherlands, 2012.
- [49] A. A. Omar, "Addressing Subjectivity in Thematic Classification of Literary Texts: Using Cluster Analysis to Derive Taxonomies of Thematic Concepts in the Thomas Hardy's Prose Fiction," Proceedings of the Chicago Colloquium on Digital Humanities and Computer Science, vol. 1, no. 2, 2010.
- [50] A. A. Omar, "Feature Selection in Text Clustering Applications of Literary Texts: A Hybrid of Term Weighting Methods," International Journal of Advanced Computer Science and Applications, vol. 11, no. 2, pp. 99-107, 2020.
- [51] A. A. Omar, "On the Digital Applications in the Thematic Literature Studies of Emily Dickinson's Poetry," International Journal of Advanced Computer Science and Applications, vol. 11, no. 6, pp. 361-365, 2020.
- [52] S. Gilmartin and R. Mengham, Thomas Hardy's Shorter Fiction: A Critical Study. Edinburgh: Edinburgh University Press, 2007, pp. x, 144 p.
- [53] P. Dalziel, "Thomas Hardy: The Excluded and Collaborative Stories." Oxford: Clarendon Press, 1992, p.^pp. Pages.
- [54] R. L. Purdy, Thomas Hardy : A Bibliographical Study. [S.l.]: Oxford University Press, 1979.
- [55] K. Brady, The Short Stories of Thomas Hardy. New York: St. Martin's Press, 1982, pp. xii, 235.
- [56] P. Dalziel, "Hardy's Unforgotten 'Indiscretion': The Centrality of an Uncontrolled Work," Review of English Studies, vol. XLIII, no. 171, pp. 347-366, August 1, 1992 1992.
- [57] R. Gittings, "An Introduction to The Hand of Ethelberta." New York: St. Martin's Press, 1978, p.^pp. Pages.
- [58] P. Widdowson, On Thomas Hardy : late essays and earlier. Basingstoke: Macmillan, 1998, pp. x, 218.

APPENDIX NO. 1: TEXTS AND NAME CODES

| Title                                      | Code     | Title                                 | Code     |
|--|----------|---------------------------------------|----------|
| A Laodicean                                | hardy001 | The First Countess of Wessex          | hardy032 |
| A Pair of Blue Eyes                        | hardy002 | Barbara of the House of Grebe         | hardy033 |
| An Indiscretion in the Life of an Heiress  | hardy003 | The Marchioness of Stonehenge         | hardy034 |
| Desperate Remedies                         | hardy004 | Lady Mottisfont                       | hardy035 |
| Far from the Madding Crowd                 | hardy005 | The Lady Icenway                      | hardy036 |
| Jude the Obscure                           | hardy006 | Squire Petrick's Lady                 | hardy037 |
| Tess of the D'Urbervilles                  | hardy007 | Anna, Lady Baxby                      | hardy038 |
| The Hand of Ethelberta                     | hardy008 | The Lady Penelope                     | hardy039 |
| The Mayor of Casterbridge                  | hardy009 | The Duchess of Hamptonshire           | hardy040 |
| The Poor Man and the Lady                  | hardy010 | The Honourable Laura                  | hardy041 |
| The Well-Beloved                           | hardy011 | A Changed Man                         | hardy042 |
| The Return of the Native                   | hardy012 | The Waiting Supper                    | hardy043 |
| The Trumpet-Major                          | hardy013 | Alicia's Diary                        | hardy044 |
| The Woodlanders                            | hardy014 | The Grave by the Handpost             | hardy045 |
| Two on a Tower                             | hardy015 | Enter a Dragoon                       | hardy046 |
| Under the Greenwood Tree                   | hardy016 | A Tryst At An Ancient Earthwork       | hardy047 |
| The Three Strangers                        | hardy017 | What The Shepherd Saw                 | hardy048 |
| The Three Strangers                        | hardy018 | A Committee-Man of The Terror         | hardy049 |
| A Tradition of Eighteen Hundred and Four   | hardy019 | Master John Horseleigh, Knight        | hardy050 |
| The Melancholy Hussar of The German Legion | hardy020 | The Duke's Reappearance               | hardy051 |
| The Withered Arm                           | hardy021 | A Mere Interlude                      | hardy052 |
| Fellow-Townsmen                            | hardy022 | The Romantic Adventures of a Milkmaid | hardy053 |
| Interlopers At The Knap                    | hardy023 | How I Built Myself a House            | hardy054 |
| The Distracted Preacher                    | hardy024 | Destiny and a Blue Cloak              | hardy055 |
| An Imaginative Woman                       | hardy025 | The Thieves Who Couldn't Help         | hardy056 |
| The Son's Veto                             | hardy026 | Our Exploits at West Poley            | hardy057 |
| For Conscience' Sake                       | hardy027 | Old Mrs. Chundle                      | hardy058 |
| A Tragedy of Two Ambitions                 | hardy028 | The Doctor's Legend                   | hardy059 |
| On the Western Circuit                     | hardy029 | The Spectre of the Real               | hardy060 |
| To Please His Wife                         | hardy030 | Blue Jimmy: The Horse Stealer         | hardy061 |
| The Fiddler of the Reels                   | hardy031 | The Unconquerable                     | hardy062 |

# Pynq-YOLO-Net: An Embedded Quantized Convolutional Neural Network for Face Mask Detection in COVID-19 Pandemic Era

Yahia Said

Electrical Engineering Department, College of Engineering  
Northern Border University, Arar, Saudi Arabia<sup>1</sup>  
Laboratory of Electronics and Microelectronics (LR99ES30)  
Faculty of Sciences of Monastir, University of Monastir, TUNISIA<sup>2</sup>

**Abstract**—The recent Coronavirus COVID-19 is a very infectious disease that is transmitted through droplets generated when an infected person coughs, sneezes, or exhales. So, people must wear a face mask to reduce the power of the transition of this virus. Governments around the world have imposed the use of face masks in public spaces and supermarkets. In this paper, we propose to build a face mask detection system based on a lightweight Convolutional Neural Network (CNN) and the YOLO object detection framework to implement it on an embedded low power device. The object detection framework was designed using a single Convolutional Neural Network for object detection in real-time. To make the YOLO framework suitable for embedded implementation, we propose to build a lightweight Convolutional Neural Network and quantize it by using a single bit for weight and 2 bits for activations. The proposed network called Pynq-YOLO-Net was implemented on the Pynq Z1 platform. The computation was divided between the software and the hardware. The features extraction part was executed on the hardware device and the output part was executed on the software. This configuration has allowed reaching real-time processing with a very good detection accuracy of 97% when tested on the combination of collected datasets.

**Keywords**—Face mask detection; Coronavirus COVID-19; YOLO; Convolutional Neural Network (CNN); embedded device; Pynq Z1 board

## I. INTRODUCTION

According to the World Health Organization (WHO) [1], the COVID-19 is causing a world crisis because of its fast infection and the absence of a cure. This new virus is considered as the fastest infecting virus all over time. Based on the latest statistics [2], more than 25 million in 190 countries are infected, and more than 844000 deaths, until the writing of this paper. As a protection step, the authorities oblige people to wear face masks in public spaces to reduce the transmission impact of the virus. Many peoples are ignoring the rules and do not wear face masks. So, it is important to detect those peoples that do not use facemasks and warn them about the importance of this step to stay uninfected by the coronavirus. Thus, an automatic face mask detector must be installed in public streets, supermarkets, and all public service agencies. Based on surveillance systems of all public spaces, it is possible to

process visual data and detect peoples that do not use face masks.

Recently, the performance of computer vision applications has been boosted to high-level thanks to the use of deep learning [3]. The deep learning is based on a deep neural network with tens of hidden layers. Deep learning models can learn directly from input data without any handcrafted features. For image processing, the Convolutional Neural Network (CNN) is the most used. It was inspired by the biological nervous system, and based on mathematics and informatics representation, it mimics the vision cortex of an animal. The CNN was successfully deployed to solve many computer vision applications such as traffic light detection and recognition [4], [5], medical image segmentation [6], indoor object detection and identification [7], [8], scene identification [9], face detection and recognition [10].

CNN models are characterized by their high performance and intensive computation. The feature extraction part (convolution layers, activation layers, pooling layers) uses the most computation effort and the output part (fully connected layers) uses the most of memory storage. Until today, Graphical Processing Units (GPU) are considered as the best target platform for the deployment of CNN. But GPUs need a lot of power and expansion. At this end, CNN must be optimized for low power platforms such as Field-Programmable Gate Arrays (FPGA). Many techniques were proposed to optimize CNN for low power implementation. The quantization technique is a very useful technique which aims to reduce the number of bits used for the representation of the weights and the activations on CNN. Many works have been proposed in this context with different methodologies. Doyun et al. [11] proposed a quantization algorithm based on the generalized gamma distribution. The proposed algorithm was tested with different representation and the achieved result were courageous. As reported in [11], the performance of the algorithm can be improved by tuning the parameters of the quantizer. A Kernel Density Estimation based Non-uniform Quantizer was proposed in [12]. In this work, a 4-bits representation of the weights and activations were used. The proposed quantization algorithm was tested on the ImageNet dataset and it was very effective in compressing the model without a big loss in performance. In [13], a quantization

algorithm based on trainable scaling factors and a nested-means clustering strategy was proposed. To quantize the weights, the nested-means clustering strategy was deployed to achieve high parameter compression. To quantize the activations, a linear quantization technique was used which take into account the statistical priorities of the batch normalization technique. There are many variants of the quantization technique and each of them has its impact on the model compression.

In this work, we propose a lightweight Convolutional Neural Network and quantize it for implementation on an edge device. The proposed CNN is composed of a convolution layer, 3 lightweight blocks, and a regression layer for output. The detection technique is based on You Look Only Once (YOLO) framework [14] which is designed to achieve real-time processing with good detection accuracy. The YOLO framework treats the detection task as a regression problem. For our case, it is perfect since we look for detecting if the person is wearing a face mask or not. That is a binary classification problem with a focus on the prediction of the bounding box used to locate the face mask. So, to speed up the processing, we ignored the classification and we focus on the localization task. The proposed network was called Pynq-YOLO-Net.

The proposed CNN was tested on the Pynq board which a hybrid board (software/hardware) equipped with an ARM processor and an FPGA in a single ship. This configuration allows taking the advantage of the FPGA blocks alongside the CPU. The Pynq board can be programmed using a high-level programming language (Python) or hardware description language (VHDL/Verilog).

The motivation behind optimizing the CNN for embedded implementation is to make it available for all surveillance systems without the need for high-performance computers and to reduce the power consumption of those systems. Also, it can be implemented in mobile devices such as smartphones and smart cameras.

The main contributions of this work are the following: (1) design a lightweight Convolutional Neural Network targeting embedded device; (2) the proposed CNN was quantized to fit in the Pynq board; (3) implementation of the proposed CNN inference for face mask detection on the Pynq board.

The rest of the paper is organized as follows. Section 2 was reserved to discuss related works about face mask detection methods. The proposed approach was described and discussed in Section 3. In Section 4, the experiment and results were presented and discussed. The paper was concluded in Section 5.

## II. RELATED WORKS

Recently, the detection of the face mask was an important application for reducing the transmission of the COVID-19. Building an automatic face mask detector is a challenging task and many works were proposed to achieve high results.

Loey et al. [15] proposed a hybrid system for face mask detection. The proposed system is composed of a CNN for feature extraction and decision trees, Support Vector Machine

(SVM), and an ensemble algorithm for the detection. The transfer learning technique was applied to the ResNet 50 model [16] to finetune it for face mask detection. The proposed system was trained and evaluated on 3 datasets, the Real-World Masked Face Dataset (RMFD) [25], the Simulated Masked Face Dataset (SMFD) [25], and the Labeled Faces in the Wild (LFW) [26]. The proposed system has achieved a high accuracy of more than 99% but it was very complex and hard to train. In addition, the proposed system is computationally intensive and cannot be used for real-time processing.

The Single Shot Multi-box Detector (SSD) [17] was proposed for face mask detection in public spaces [18]. The SSD model was pre-trained on the MSCOCO dataset for object detection and finetuned on a custom-made dataset for face mask detection. The MobileNet V2 model [19] was used as a backbone for the SSD to limit the computation complexity. The proposed model was implemented on a Raspberry PI 4 equipped with a quad-core ARM processor and 4GB of RAM. An accuracy of 85% was achieved when testing the model on the custom-made dataset. This work was a good step for implementing facemasks on embedded devices. But the Raspberry PI 4 is considered as a software device and its power consumption is too high compared to low power devices.

Jiang et al. [19] proposed the use of the RetinaNet model [20] for face mask detection. The RetinaNet was finetuned for face mask detection through the transfer learning technique. Two backbones were tested, the Resnet and the MobileNet. In addition, a new technique was added to the RetinaNet to reject predictions with low confidences and the high intersection of a union. The RetinaNet was pre-trained on the ImageNet [21] dataset and then fine-tuned on the face mask dataset. The proposed RetinaFaceMask has achieved good results with both backbones while the best results were achieved using the ResNet model. The achieved result was good but the RetinaFaceMask was not suitable for implementation on low power devices because of its computation intensively and the need for large storage memory.

IN [22], a CNN model was proposed to detect if a person wears a face mask or not. Also, the proposed network was used to detect if the mask is correctly worn or not. The proposed CNN has a simple architecture with a convolution layer, an activation layer, a pooling layer, a fully connected layer, and a softmax layer. The proposed approach was designed to detect faces and face masks separately. The proposed CNN model was trained using publicly available datasets, Masked Face Detection Dataset (MFDD) [25], Real-world Masked Face Recognition Dataset (RMFRD) [25], and Simulated Masked Face Recognition Dataset (SMFRD) [25]. The reported results are good in terms of accuracy and speed. The main disadvantage of the proposed model is the need for high-performance computers and large memory usage.

All the mentioned methods are designed to be implemented on a high-performance computer with a very high-power consumption. In this work, we propose to quantize a lightweight CNN for implementation on low power devices with a focus on high performance. In the next section, we will



present the proposed approach and detailing the different optimization applied to achieve an embedded implementation.

### III. PROPOSED APPROACH

In this section, we will describe the proposed lightweight CNN model and the compression techniques applied to make this model fit in the resource constraint of a low power device while maintaining high performance and real-time processing.

Recently, many techniques are proposed to build lightweight CNN models. The most important technique in the use of Bottlenecks instead of normal convolution layers. In this work, we adopt the Bottlenecks concept proposed by the MobileNet v2 model [19]. The main contribution of the Inverted Residuals and Linear Bottlenecks is the use of depthwise convolution and point convolutions instead of simple convolution layers with adding a residual connection. The depthwise convolution is similar to the normal convolution layer but the main difference that depthwise convolution reserves the number of channels and does not compress them. For normal convolution it the number of input channels is  $n$  then the number of output channels is 1 but for depthwise convolution, if the number of the input channels is  $n$  then the number of output channels is  $n$ . The pointwise convolution is a normal convolution layer with a kernel size of  $1 \times 1$ . The combination of a depthwise convolution with a pointwise convolution makes the same functionality of a normal convolution layer but 9 times faster as they claim [19]. Also, the separation of the filtering and combining functionalities allow the implementation of more than one activation layer and batch normalization layer which results in enhancing the performance of the model and reducing the computation complexity. The proposed Inverted Residuals and Linear Bottlenecks are presented in Fig. 1. Another technique was deployed by the MobileNet model was the use of strided convolution layers and eliminate the use of pooling layers. As proves in [23], the use of strided convolution layers instead of max-pooling layers is more efficient to build CNN models for embedded implementation and helps to enhance the accuracy of those models.

In this work, we propose to use a convolution layer with a kernel size of  $3 \times 3$  and a stride of 2, three inverted residual bottlenecks, and three linear bottlenecks as a backbone for the YOLO framework. The YOLO framework takes an input image, applies a feature extraction through a backbone based on a CNN model to generate an output grid of  $N \times N$  dimension. For each cell of the obtained grid, it predicts only one object with the parameters of the bounding box (the  $x$ ,  $y$  coordinates, the height, the width, and the confidence score) and the class

probability. For the Pascal VOC dataset [24] the YOLO framework generated  $7 \times 7$  grid and used 2 bounding boxes (B) for 20 classes (C). The architecture of the YOLO framework is presented in Fig. 2.

In this work, we propose to eliminate the calculation of the class probability because we are solving a binary classification problem. The YOLO framework computes the score confidence for each predicted bounding box. Since there is one object to detect, we consider the confidence score the class probability. This step allows reducing the calculation at the output layer and speeds up the processing speed.

Besides, more optimizations were applied to the YOLO framework to reduce the computation complexity and to enhance the performance. First, the input image was resized to a power of 2 sizes. Thus, we propose to use  $128 \times 128$  size. After applying the features extraction module, the YOLO framework generates  $8 \times 8$  grid. Using an input image with a size power of 2 facilitates the implementation of the convolutional layers on the hardware device because it is easier to perform multiplications by using only shifting registers. Second, the use of an all convolution backbone allows enabling the data reuse technique to reduce the communication between the external memory and use only the on-chip memory to compute the convolutions. Third, the pruning technique was applied to eliminate weak connections and to avoid the overfitting problem in the finetuning step. Finally, the model was quantized by replacing 32 bits floating-point representation by a 2-bits fixed-point representation for the activations (A) and 1-bit fixed-point representation for the weights (W). For the input image and the output layer (grid), we used an 8-bits fixed-point representation. The backbone architecture and configuration were presented in Fig. 3.

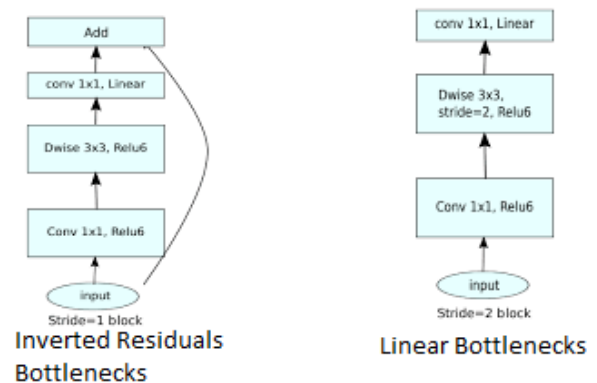


Fig. 1. Inverted Residuals and Linear Bottlenecks.

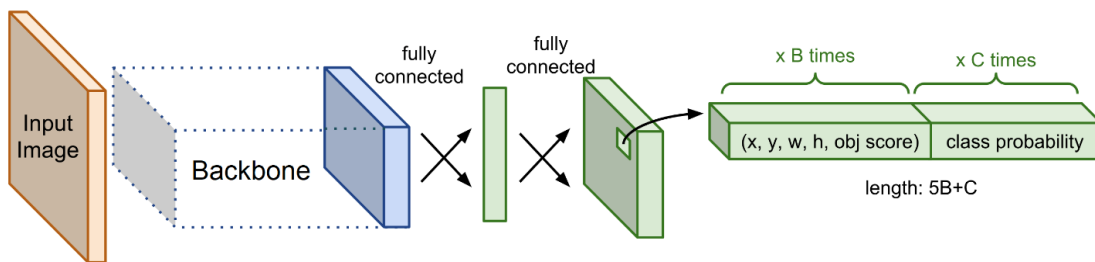


Fig. 2. YOLO Architecture.

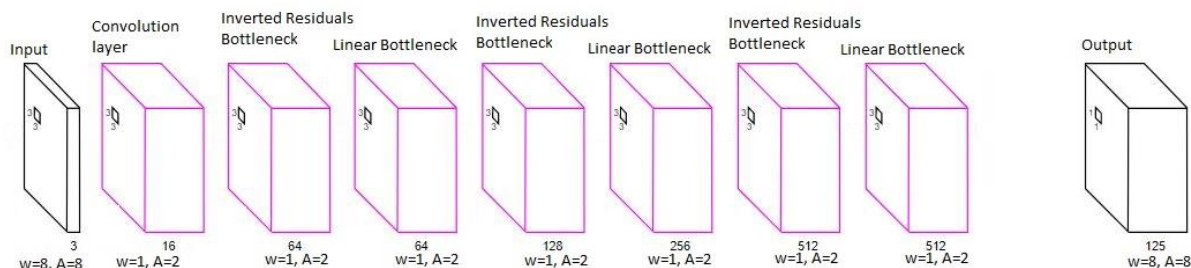


Fig. 3. Backbone Architecture for Embedded Implementation.

In this work, we propose to use the post-training quantization technique which is based on training and fine-tuning after reducing the presentation of the weights and activations. The YOLO was trained using 32 bits floating-point representation then it was reduced to the proposed representations and retrained again to recover the accuracy. In the experiments, we will report the accuracies obtained with different representations. The retraining process is very important to recover the degradation of the accuracy caused by the quantization technique.

The workflow of the proposed approach is divided into 5 steps. The first step is to develop the proposed model based on the YOLO framework. The second step is to train the proposed model using a specific dataset. The third step is to optimize the model for embedded implementation by applying the pruning technique and the quantization technique. The fourth step is the retraining of the model on the same dataset used in step 2 to recover the accuracy degraded by the optimization techniques. The final step is to implement the model on the Pynq Z1 board. The workflow of the proposed approach is illustrated in Fig. 4.

#### IV. EXPERIMENTS AND RESULTS

##### A. Training Data

To train the proposed model, we proposed to combine publicly available datasets to increase the amount of training data. The performance of CNN models increases with the amount of training data. The used datasets are presented in the following:

- Real-World Masked Face Dataset (RMFD) [25]: The images of the dataset were collected automatically from the internet for public famous figures with and without a mask. The dataset contains 5,000 images of 525 persons wearing masks, and 90,000 images of the same 525 persons without masks. The dataset was manually filtered and annotated using semi-automatic labeling tools.
- Masked Face Detection Dataset (MFDD) [25]: The dataset was designed for the detection of masked faces during the era of the coronavirus. This dataset combines existing face detection datasets with images collected from the internet. The collected images were manually annotated where the coordinate of the face with the mask was defined in addition to the condition of wearing a mask or not. It contains 24771 images for masked faces.

- Simulated Masked Face Recognition Dataset (SMFRD) [25]: to increase the amount of training data for face mask detection, an automatic wearing tool was designed to add masks to faces of existing face detection and recognition datasets such as LFW [26] and Webface [27] then the collected data was added to the MFDD. This dataset allows adding 500000 annotated faces of 10000 persons to the MFDD.
- MAsked FAces (MAFA) dataset [28]: it is a dataset that contains 30,811 images and 34,806 labeled masked faces. The dataset contains faces masked by the medical mask and others masked hand or other objects. This allows enhancing the generalization power by distinguishing between the mask that it must be detected and other masks.

The mentioned datasets were combined to build a very large dataset to increase the training data. Thus, it will enhance the performance of the trained model. Fig. 5 present examples of images from the collected datasets used for the training of the model.

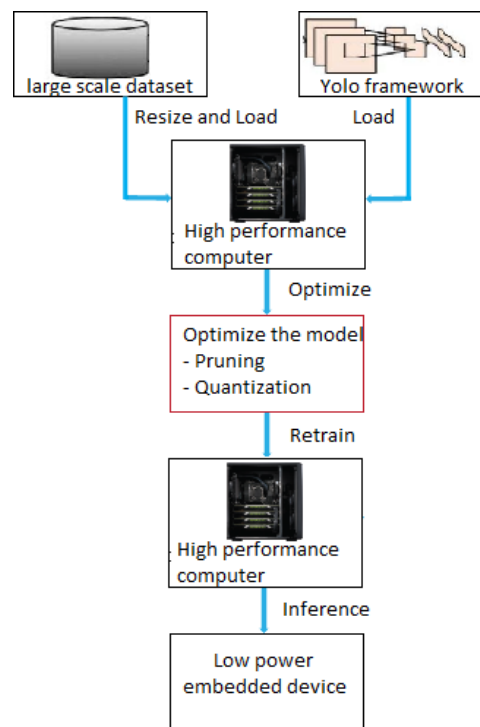


Fig. 4. Workflow of the Proposed Approach.



Fig. 5. Examples of Images from the Collected Datasets.

### B. Training and Evaluation

The proposed model was trained on the combination of the collated datasets using the gradient descent algorithm. The Adam optimizer was used as a learning algorithm which is a variant of the gradient descent algorithm with many advantages. The Adam Optimizer optimizes the weights and the learning rate accordingly to achieve a better minimum of the loss function.

To evaluate the proposed Pynq-YOLO-Net, we propose to use the precision and the recall as evaluation metrics. The precision presents the percentage of relevant results and the recall presents the percentage of relevant results correctly identified. The precision and the recall are computed as (1).

$$\text{Precision} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}}$$
$$\text{Recall} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}} \quad (1)$$

The performance of the Pynq-YOLO-Net was evaluated on the testing set which is 30% of the collected data. The Pynq-YOLO-Net achieved a Precision of 94.6% and a Recall of 95.8% for the first training process using the 32-bits floating-point representation. After compressing the model and retrained it on the same data, it achieves an accuracy of 90.7% of Accuracy and 92.3 of Recall. To find our model in the state-of-the-art, we compared against existing works. Table I present a comparison against state-of-the-art works. As shown in Table I, the proposed Pynq-YOLO-Net achieved better results than state-of-the-art works in its normal version. The compressed model achieves lower results but with the advantage of implementation on low power devices. The achieved results still good enough to generate trusted predictions.

### C. Inference

The inference of the Pynq-YOLO-Net was implemented on the Pynq Z1 board. The board was connected to the internet and it was connected via an ssh node to visualize the results on the computer screen. The Pynq Z1 board is presented in Fig. 6. An operating system based on Linux kernel was loaded to the

board via a pre-booted SSD card. The implementation of the Pynq-YOLO-Net on the Pynq Z1 board achieves a processing speed of 16 FPS. The achieved result can be considered as real-time processing speed. The energy consumption of the board does not exceed 5 watts.

The implementation of the proposed model on the Pynq Z1 board was divided into parts. The first part composed of the feature extraction, which is composed of the convolution layer and the bottlenecks, was implemented on the hardware to take advantage of the parallel processing of the programmable units.

The second part which is composed of the fully connected layers and the output layer, was implemented on the software because it needs more memory and less computation. The performance of the board is presented in Fig. 7.

The Pynq-YOLO-Net was tested using images that does not belong to the collected datasets to evaluate the generalization power of the model. Fig. 8 present an illustration of the obtained results. The model was very effective when tested on new images which prove that it have a good generalization power.

TABLE I. COMPARISON AGAINST STATE-OF-THE-ART WORKS

| Model                           | Precision (%) | Recall (%) |
|---------------------------------|---------------|------------|
| RetinaFaceMask [19]             | 93.4          | 94.5       |
| SSD [18]                        | 91            | 91         |
| Pynq-YOLO-Net (ours)            | 94.6          | 95.8       |
| Pynq-YOLO-Net compressed (ours) | 90.7          | 92.3       |



Fig. 6. Pynq Z1 Board.

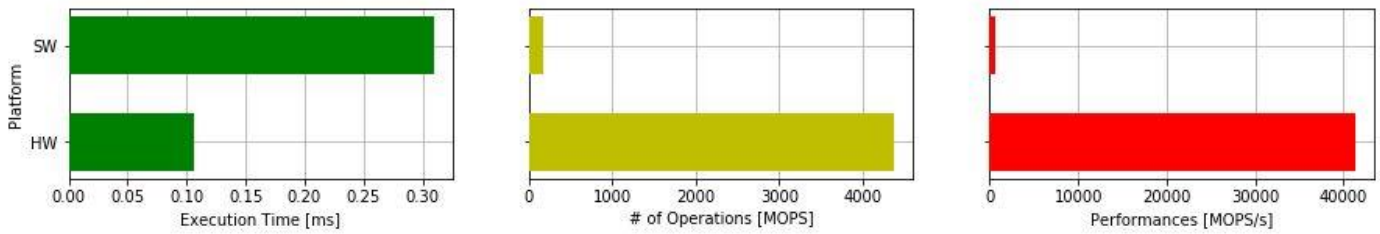


Fig. 7. Performance of the Pynq Z1 Board.



Fig. 8. Visualization of the Obtained Results of the Pynq-YOLO-Net.

#### D. Discussion

The reported results prove the efficiency of the proposed Pynq-YOLO-Net for implementation in low power devices. Starting by building a lightweight CNN is a very important step to reach embedded implementations. Also, the YOLO framework was a good choice since it was designed with a focus on speed. The model compression techniques used in this work allow to reduce the size of the model and speed up the processing speed without damaging the accuracy. The choice of the size of the input images was very effective for building the convolution layers on the hardware part of the board. All those factors were correlated together to achieve an embedded implementation of the proposed model.

#### V. CONCLUSIONS

The coronavirus COVID-19 is a very fast-spreading disease. It is important to protect ourselves from being infected by wearing masks and respecting social distances in public environments. In this paper, we propose to build a face mask detector in public spaces to detect if people are wearing masks or not. The proposed detector was based on the YOLO framework with a lightweight backbone. The proposed model, called Pynq-YOLO-Net, was designed to be implemented on the Pynq Z1 board. To achieve this implementation, some model compression techniques was applied such as pruning and quantization. Those techniques were very effective to reduce the model size and the computation complexity. The model was implemented on both hardware and software to accelerate the inference. The achieved performance has proved the efficiency of the proposed approach for mask detection in public spaces. As future work, the model will be implemented on video surveillance systems to be tested on real conditions.

#### ACKNOWLEDGMENT

The author wishes to acknowledge the help of Mr. Ayachi Riadh from Laboratory of Electronics and Microelectronics at University of Monastir for assistance with implementing the proposed model.

#### REFERENCES

- [1] Coronavirus disease (COVID-19) pandemic, Available at: <https://www.who.int/emergencies/diseases/novel-coronavirus-2019>, accessed on 29/08/2020.
- [2] COVID-19 CORONAVIRUS PANDEMIC, Available at: <https://www.worldometers.info/coronavirus/>, accessed on 29/08/2020.
- [3] Goodfellow, Ian, Aaron Courville, and Yoshua Bengio. Deep learning. Vol. 1. Cambridge: MIT press, 2016.
- [4] Ayachi, Riadh, Mouna Afif, Yahia Said, and Mohamed Atri. "Traffic signs detection for real-world application of an advanced driving assisting system using deep learning." *Neural Processing Letters* 51, no. 1 (2020): 837-851.
- [5] Ayachi, R., Y. E. Said, and M. Atri. "To perform road signs recognition for autonomous vehicles using cascaded deep learning pipeline." *Artificial Intelligence Advances* 1, no. 1 (2019): 1-58.
- [6] Wang, Guotai, Wenqi Li, Maria A. Zuluaga, Rosalind Pratt, Premal A. Patel, Michael Aertsen, Tom Doel et al. "Interactive medical image segmentation using deep learning with image-specific fine tuning." *IEEE transactions on medical imaging* 37, no. 7 (2018): 1562-1573.
- [7] Afif, Mouna, Riadh Ayachi, Yahia Said, Edwige Pissaloux, and Mohamed Atri. "An evaluation of retinanet on indoor object detection for blind and visually impaired persons assistance navigation." *Neural Processing Letters* (2020): 1-15.
- [8] Afif, Mouna, Riadh Ayachi, Edwige Pissaloux, Yahia Said, and Mohamed Atri. "Indoor objects detection and recognition for an ICT mobility assistance of visually impaired people." *Multimedia Tools and Applications* (2020): 1-18.
- [9] Afif, Mouna, Riadh Ayachi, Yahia Said, and Mohamed Atri. "Deep Learning Based Application for Indoor Scene Recognition." *Neural Processing Letters* (2020): 1-11.
- [10] Sun, Xudong, Pengcheng Wu, and Steven CH Hoi. "Face detection using deep learning: An improved faster RCNN approach." *Neurocomputing* 299 (2018): 42-50.
- [11] Kim, Doyun, Han Young Yim, Sanghyuck Ha, Changgwun Lee, and Inyup Kang. "Convolutional Neural Network Quantization using Generalized Gamma Distribution." *arXiv preprint arXiv:1810.13329* (2018).
- [12] Seo, Sanghyun, and Juntae Kim. "Efficient weights quantization of convolutional neural networks using kernel density estimation based non-uniform quantizer." *Applied Sciences* 9, no. 12 (2019): 2559.
- [13] Schindler, Günther, Wolfgang Roth, Franz Pernkopf, and Holger Fröning. "N-Ary Quantization for CNN Model Compression and Inference Acceleration." (2018).
- [14] Redmon, Joseph, Santosh Divvala, Ross Girshick, and Ali Farhadi. "You only look once: Unified, real-time object detection." In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 779-788. 2016.
- [15] Loey, Mohamed, Gunasekaran Manogaran, Mohamed Hamed N. Taha, and Nour Eldeen M. Khalifa. "A Hybrid Deep Transfer Learning Model

- with Machine Learning Methods for Face Mask Detection in the Era of the COVID-19 Pandemic." *Measurement* (2020): 108288.
- [16] He, Kaiming, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. "Deep residual learning for image recognition." In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770-778. 2016.
- [17] Liu, Wei, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C. Berg. "Ssd: Single shot multibox detector." In *European conference on computer vision*, pp. 21-37. Springer, Cham, 2016.
- [18] Yadav, Shashi. "Deep Learning based Safe Social Distancing and Face Mask Detection in Public Areas for COVID-19 Safety Guidelines Adherence." *International Journal for Research in Applied Science & Engineering Technology (IJRASET)*. 2020.
- [19] Jiang, Mingjie, and Xinqi Fan. "RetinaFaceMask: A Face Mask detector." *arXiv preprint arXiv:2005.03950* (2020).
- [20] Lin, Tsung-Yi, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. "Focal loss for dense object detection." In *Proceedings of the IEEE international conference on computer vision*, pp. 2980-2988. 2017.
- [21] Deng, Jia, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. "Imagenet: A large-scale hierarchical image database." In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248-255. Ieee, 2009.
- [22] Inamdar, Madhura, and Ninad Mehendale. "Real-Time Face Mask Identification Using Facemasknet Deep Learning Network." Available at SSRN 3663305 (2020).
- [23] Ayachi, Riadh, Mouna Afif, Yahia Said, and Mohamed Atri. "Strided convolution instead of max pooling for memory efficiency of convolutional neural networks." In *International conference on the Sciences of Electronics, Technologies of Information and Telecommunications*, pp. 234-243. Springer, Cham, 2018.
- [24] Everingham, Mark, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. "The pascal visual object classes (voc) challenge." *International journal of computer vision* 88, no. 2 (2010): 303-338.
- [25] Wang, Zhongyuan, Guangcheng Wang, Baojin Huang, Zhangyang Xiong, Qi Hong, Hao Wu, Peng Yi et al. "Masked face recognition dataset and application." *arXiv preprint arXiv:2003.09093* (2020).
- [26] Huang, Gary B., Marwan Mattar, Tamara Berg, and Eric Learned-Miller. "Labeled faces in the wild: A database for studying face recognition in unconstrained environments." 2008.
- [27] Yi, Dong, Zhen Lei, Shengcai Liao, and Stan Z. Li. "Learning face representation from scratch." *arXiv preprint arXiv: 1411.7923* (2014).
- [28] Ge, Shiming, Jia Li, Qiting Ye, and Zhao Luo. "Detecting masked faces in the wild with lle-cnns." In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2682-2690. 2017.

# Best Path in Mountain Environment based on Parallel Hill Climbing Algorithm

Raja Masadeh<sup>1</sup>

Computer Science department  
The World Islamic Sciences and  
Education University  
Amman, Jordan

Ahmad Sharieh<sup>2</sup>, Sanad Jamal<sup>3</sup>

Mais Haj Qasem<sup>4</sup>  
Computer Science department  
The University of Jordan  
Amman, Jordan

Bayan Alsaaidah<sup>5</sup>

Computer Science department  
Al-Balqa Applied University  
Al-Salt, Jordan

**Abstract**—Heuristic search is a search process that uses domain knowledge in heuristic rules or procedures to direct the progress of a search algorithm. Hill climbing is a heuristic search technique for solving certain mathematical optimization problems in the field of artificial intelligence. In this technique, starting with a suboptimal solution is compared to starting from the base of the hill, and improving the solution is compared to walking up the hill. The optimal solution of the hill climbing technique can be achieved in polynomial time and is an NP-complete problem in which the numbers of local maxima can lead to an exponential increase in computational time. To address these problems, the proposed hill climbing algorithm based on the local optimal solution is applied to the message passing interface, which is a library of routines that can be used to create parallel programs by using commonly available operating system services to create parallel processes and exchange information among these processes. Experimental results show that parallel hill climbing outperforms sequential methods.

**Keywords**—Hill climbing; heuristic search; parallel processing; Message Passing Interface (MPI)

## I. INTRODUCTION

Hill climbing algorithm based on the local optimal solution was proposed and applied to the Message Passing Interface (MPI), which is a library of routines that can be used to create parallel programs in C, C++, and Fortran 77 by using commonly available operating system services to create parallel processes and exchange information among these processes [1]. In this algorithm, the 10 closest points around the current point are scanned, and the cost needed to go from the current point to the next point is obtained by calculating the sum of the obstacles between the current point and the 10 other points. The MPI method is used to validate the performance of the hill climbing algorithm by using parallel and distributed computing systems compared with sequential methods [2].

Hill climbing is a heuristic search technique for solving certain mathematical optimization problems in the field of artificial intelligence [3]. In this technique, starting with a suboptimal solution is compared to starting from the base of the hill, and improving the solution is compared to walking up the hill. The solution is improved repeatedly until a certain condition is maximized and becomes optimal. This technique is mainly used to solve difficult problems computationally [4].

Heuristic search is an artificial intelligence search technique and a computer simulation of thinking that utilizes heuristic for its moves [5, 6]. Heuristic search is a search process that uses domain knowledge in heuristic rules or procedures to direct the progress of a search algorithm, is utilized to prune the search space, and is adopted in applications where a combinatorial explosion indicates that an exhaustive search is impossible [7].

The objective of heuristic search is to produce a solution in a reasonable time frame that is sufficient to solve the problem at hand. This solution may not be the best of all the solutions to this problem, or it may simply approximate the exact solution, but it is still valuable because finding it does not require a prohibitively long time. For large and complex problems, finding an optimal solution path can take a long time and a suboptimal solution that can be obtained rapidly may be useful. Various techniques for modifying a heuristic search algorithm, such as hill climbing, to allow a tradeoff between solution quality and search time have been investigated [8, 9].

The optimal solution of Hill Climbing technique can be achieved in polynomial time and it is one of the NP-Complete problem that the numbers of local maxima can be the cause of exponential computational time. To address these problems parallel and distributed computing systems can be applied to Hill Climbing algorithm. Parallel and distributed computing systems are high-performance computing systems that spread out a single application over many multi-core and multi-processor computers in order to rapidly complete the task. Parallel and distributed computing systems divide large problems into smaller sub-problems and assign each of them to different processors in a typically distributed system running concurrently in parallel. MPI are one these computing systems [10, 11].

The MPI is a standardized means of exchanging messages among multiple computers running a parallel program across a distributed memory. The MPI is generally considered to be the industry standard and forms the basis for most communication interfaces adopted by parallel computing programmers. The MPI is used to improve scalability, performance, multi-core and cluster support, and interoperation with other applications [12].

The rest of the paper is organized as follows. Section II reviews works that are closely related to the hill climbing algorithm. Sections III and IV present the methodology of the new proposed algorithm and its analysis. Section V presents the experimental results. Section VI provides the conclusion.

## II. RELATED WORK

Mathematicians and research scientists have found many applications that use heuristic search. The increased use of heuristic search in a wide variety of applications, such as science, engineering, economics, and technology, is due to the advent of personal and large-scale computers.

Rashid and Tao [13] presented an optimized hill climbing algorithm called parallel iterated local search with efficiently accelerated GPUs. They also tested the algorithm by using a typical case study of the graph bisection in computational science. The proposed algorithm minimizes the data transfer between two components to achieve the best performance. Then, the purpose of the parallelism control is to control the generation of the neighborhood for meeting the memory constraints and finding efficient mappings between neighborhood candidate solutions and GPU threads. The authors found through experiments that GPU computing not only accelerates the search process but also exploits parallelism to improve the quality of the obtained solutions for combinatorial optimization problem.

Jiang et al. [14] proposed an optimal power allocation (OPA) method to exert the maximum efficiency of parallel grid-connected inverters. They established the power model of every inverter and compared each model by using the equal power allocation (EPA) method. Theoretically, high overall system efficiency can be achieved using the OPA method. Then, the authors calculated the overall system efficiency with easily measurable electric parameters and realized online optimization by adopting an existing method, such as the hill climbing method. They tested the effectiveness of the proposed method by comparing it with the equivalent power allocation method. They found that the overall system efficiency of the OPA method is higher than that of the EPA method. Moreover, the system using the hill climbing method performs effectively in the dynamic process with short response time.

Kim et al. [15] performed component sizing of power sources of parallel hybrid vehicle by applying the golden section search and hill climbing algorithms. The golden section search algorithm was used in selecting a reduction gear ratio that connects the transmission to the electric motor by using the hill climbing search algorithm to find the optimal engine and electric motor sizes. The use of the hill climbing search algorithm reduces the number of simulations and simultaneously optimizes the capacity of the power source and the gear ratio of the torque coupler. The authors verified the validity of the component sizing results by comparing the global optimal solution obtained by the conventional technique with the solution obtained by the proposed optimization technique.

Robinson et al. [16] presented an improved algorithm for approximating the TSP on fully connected, symmetric graphs by utilizing the GPU. They improved an existing 2-opt hill

climbing algorithm with random restarts by considering multiple updates to the current path found in parallel. Their approach has a k-swap function, which allows k number of updates per iteration. The authors showed that their modifications result in a substantial speedup without a reduction in the quality of the result by applying the k-swap method. Their experimental results showed that common assumptions in obtaining good performance for the GPU are not always true, such as saturating the GPU with blocks. Instead, for problems in which the search space can be deterministically enumerated, the number of active blocks can be limited as determined by the hardware. This property allows for reduced memory allocation. A limited amount of memory can be used because each block can allocate the amount of memory upfront.

Chen et al. [17] proposed an automatic machine learning (AutoML) modeling architecture called Autostacker, which is a machine learning system with an innovative architecture for automatic modeling and a well-behaved efficient search algorithm. Autostacker improves the prediction accuracy of machine learning baselines by utilizing an innovative hierarchical stacking architecture and an efficient parameter search algorithm. Neither prior domain knowledge about the data nor feature preprocessing is needed. The authors reduced the time of AutoML by using a naturally inspired algorithm called PHC. They demonstrated the operation and performance of their system by comparing it with human initial trails and related state-of-the-art techniques. They also confirmed the scaling and parallelization ability of their system. The authors also automated the machine learning modeling process by providing an efficient, flexible, and well-behaved system. This system can be generalized into complicated problems and integrated with data and feature processing modules.

## III. METHODOLOGY

Hill climbing is a heuristic search technique for solving certain mathematical optimization problems in the field of artificial intelligence [18]. In this technique, starting with a suboptimal solution is compared to starting from the base of the hill, and improving the solution is compared to walking up the hill; the solution is improved repeatedly until some condition is maximized and becomes optimal, as illustrated in Fig. 1. This technique is mainly used for solving difficult problems computationally. It focuses on the current and immediate future states and does not maintain a search tree, thereby making it memory efficient [19].

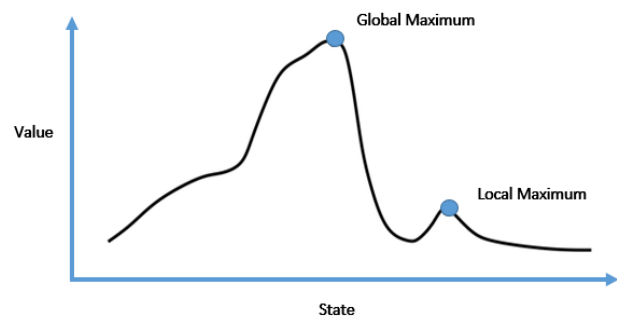


Fig. 1. Hill Climbing.

Hill climbing technique can be used to solve many problems, such as network flow, traveling salesman, and 8-Queens, in which the current state allows for an accurate evaluation function [20]. This technique does not suffer from space-related issues because it focuses on the current state in which previously explored paths are not stored; nonetheless, an optimal solution can be achieved in polynomial time. However, for NP-complete problems, computational time can be exponential based on the number of local maxima [21].

Hill climbing technique comprises several phases, which start by constructing a suboptimal solution considering the constraints of the problem, followed by improving the solution by step and enhancing the solution until no more improvement is possible [22]. This technique is performed following the steps below.

- 1) Define the current state as initial state.
- 2) Loop until the goal state is achieved or no more operators can be applied on the current state.
  - a) Apply an operation to the current state and obtain a new state.
  - b) Compare the new state with the goal state.
  - c) Quit if the goal state is achieved.
  - d) Evaluate the new state with heuristic function and compare it with the current state.
  - e) If the newer state is closer to the goal than the current state, then update the current state.

In the hill climbing algorithm, achieving the goal is equivalent to reaching the top of the hill.

In this research, a new hill climbing algorithm is proposed on the basis of local optimal solution. In this algorithm, the closest 10 points around the current point are scanned, and the cost needed to go from the current point to the next point is obtained by calculating the sum of the obstacles between the current point and the 10 other points. In this proposed algorithm we have designed the optimization technique as explained in the following equation:

$$Optimize \ Min \left( \sum_{Right \ Paths=1}^5 w_s * S + w_g * G + w_o * O, \sum_{Left \ Paths=1}^5 w_s * S + w_g * G + w_o * O \right)$$

Where  $w_s$  is the weight of the slop,  $S$  is the slop of the point,  $w_g$  is the weight of the gravity,  $G$  is the gravity at that point,  $w_o$  is the weight of the obstacles,  $O$  is the value of the obstacles.

These values will be counted for each path from right and from left. Then, the path which has the lowest cost will be selected.

The scanning approach is performed as follows:

In Table I, the current point is in red (23); each cell has a weight represented by a number. The scanning area is five paths to the right and five paths to the left.

Paths from left:

- 1) 23→16→16→12→16 (83).
- 2) 23→15→23→10→32 (103).
- 3) 23→15→16→56→20 (103).
- 4) 23→15→16→16→15 (85).
- 5) 23→15→16→16→23 (93).

Paths from Right:

- 1) 23→65→26→15→23 (152).
- 2) 23→23→15→20→32 (113).
- 3) 23→23→29→30→10 (115).
- 4) 23→23→29→28→32 (135).
- 5) 23→23→29→28→72 (175).

Thereafter, we decide which path should be taken depending on the minimum value among the total obstacle weights in each path. In the example above, the best path is number 1 because it has the minimum total value. Fig. 2 illustrates the Proposed Sequential Algorithm.

The above pseudocode is for sequential execution. To parallelize this algorithm, we must follow the following approach:

- 1) There will be a master node that will generate the matrix.
- 2) The master node must fill the matrix with the obstacle's weights based on the following equation so that the algorithm can calculate the cost.

$$Cost = w_s * S + w_g * G + w_o * O$$

$w_s$  is the weight of the slop,  $S$  is the slop of the point,  $w_g$  is the weight of the gravity,  $G$  is the gravity at that point,  $w_o$  is the weight of the obstacles,  $O$  is the value of the obstacles.

- 1) The master node will broadcast the matrix to all other nodes so they can work in parallel.
- 2) Each node will calculate its start region from bottom and its end region from top as shown in Fig. 3.
- 3) All the node will start working at the same time.
- 4) After they all finish, each node will send the best path for the master node.
- 5) Finally, the master node will decide which path is the best path based on what did it get from the other nodes.

TABLE I. PROPOSED ALGORITHM EXAMPLE

|    |    |    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|----|----|
| 2  | 10 | 20 | 26 | 23 | 23 | 15 | 16 | 10 | 20 | 23 |
| 51 | 32 | 30 | 15 | 65 | 32 | 16 | 23 | 56 | 15 | 65 |
| 32 | 23 | 12 | 54 | 72 | 28 | 29 | 30 | 21 | 16 | 72 |
| 16 | 32 | 15 | 32 | 64 | 23 | 95 | 65 | 12 | 32 | 64 |
| 2  | 10 | 20 | 26 | 23 | 23 | 15 | 16 | 10 | 20 | 23 |
| 51 | 32 | 30 | 15 | 65 | 32 | 16 | 23 | 56 | 15 | 65 |
| 54 | 72 | 28 | 29 | 23 | 23 | 15 | 16 | 16 | 23 | 23 |
| 32 | 64 | 23 | 95 | 65 | 32 | 16 | 23 | 29 | 30 | 65 |



```
Let Mat [100][100];
Set StartPoint;
Let PRow = 0;
Let Check = 0;
Sub GO(ByVal rowIndex As Integer, ByVal Colindex As Integer)
    If PRow = rowIndex Then
        check = 0
    If rowIndex >= 4 Then
        GO(rowIndex - 1, Colindex + 1)
    End If
    Exit Sub
End If
If PRow = rowIndex Then
    check += 1
End If
PRow = rowIndex
Dim arr As New ArrayList
Dim Rounds As Integer = 4
Dim TotalRounds As Integer = 4
For j As Integer = Colindex To Colindex + 4
    Dim counter As Integer = 0
    Dim sum As Double = 0
    For i As Integer = rowIndex To rowIndex - Rounds Step -
1
        sum += grd.Rows(i).Cells(j + counter).Value
        counter += 1
    Next
    For i As Integer = j - 1 To Colindex Step -1
        If j <> Colindex Then
            sum += grd.Rows(rowIndex).Cells(i).Value
        End If
    Next
    arr.Add(rowIndex - Rounds & "," & Colindex +
TotalRounds)

    arr.Add(sum)
    Rounds -= 1
Next
Rounds = 4
For j As Integer = Colindex To Colindex - 4 Step -1
    Dim counter As Integer = 0
    Dim sum As Double = 0
    For i As Integer = rowIndex To rowIndex - Rounds Step -1
        sum += grd.Rows(i).Cells(j - counter).Value
        counter += 1
    Next
    For i As Integer = j + 1 To Colindex
        If j <> Colindex Then
            sum += grd.Rows(rowIndex).Cells(i).Value
        End If
    Next
```

```
arr.Add(rowIndex - Rounds & "," & Colindex - TotalRounds)
arr.Add(sum)
Rounds -= 1
Next
Dim min As Integer = arr(1)
Dim Row As Integer = 0
Dim Col As Integer = 0
Dim Sign As Integer = txtDest.Text - Colindex
Dim Right As Integer = Math.Abs(txtDest.Text - (Colindex
+ TotalRounds))
Dim Left As Integer = Math.Abs(txtDest.Text - (Colindex -
TotalRounds))
Dim SelectedPath As Integer = 0
If Right < Left AndAlso Colindex + TotalRounds <
Colindex + TotalRounds + Right Then
    SelectedPath = Colindex + TotalRounds

ElseIf Left < Right AndAlso Colindex - TotalRounds >
Colindex - TotalRounds - Left Then
    SelectedPath = Colindex - TotalRounds
    min = arr(11)
End If
For i As Integer = 1 To arr.Count - 1 Step 2
    If txtSpiciifc.Text = 1 AndAlso SelectedPath <> 0 Then
        If arr(i) <= min AndAlso arr(i - 1).ToString.Split(",")(1) =
SelectedPath Then
            min = arr(i)
            Row = arr(i - 1).ToString.Split(",")(0)
            Col = arr(i - 1).ToString.Split(",")(1)
        End If
    Else
        If arr(i) <= min Then
            min = arr(i)
            Row = arr(i - 1).ToString.Split(",")(0)
            Col = arr(i - 1).ToString.Split(",")(1)
        End If
    End If
End If

Next

grd.Rows(Row).Cells(Col).Style.BackColor = Color.Red
coloring(Col, Colindex, Row, rowIndex)
Dim endT As TimeSpan = Now.TimeOfDay
TotalTime += (endT - start).Milliseconds
If Row - 4 >= 0 AndAlso Col + 4 < grd.Columns.Count
AndAlso Row > 0 Then
    GO(Row, Col)

End If
check = Not check
End Sub
```

Fig. 2. Sequential Pseudocode.

For example, after broadcasting the matrix to all node, each processor will choose start region and end region depending on its ID. Thus, we will divide the very last row between the processor based on the following equations to get the start region for each processor:

$$\text{Block Size} = \frac{\text{Number of starting points}}{\text{Number of processors}}$$

$$\text{Starts From} = \text{Processor ID} * \text{Block Size}$$

$$\text{Ends At} = \text{Starts From} + \text{Block Size} - 1$$

To get the end region for each processor, we will divide the very first row in the matrix between the processors based on the following equations:

$$\text{Block Size} = \frac{\text{Number of Ending points}}{\text{Number of processors}}$$

$$\text{Starts From} = \text{Processor ID} * \text{Block Size}$$

$$\text{Ends At} = \text{Starts From} + \text{Block Size} - 1$$

Fig. 3 illustrates this example with 10 points and 5 processors.

In the proposed system we have considered three approaches for finding the best path.

**Approach #1:** From All points below to unknown point above. In this approach the algorithm will start from each point below the hill and try to find a path depending on the discussed algorithm above, but the destination is not specified. So, the algorithm will suggest the path and will determine the destination point. To parallelize this approach each processor will start from the points in its region only.

**Approach #2:** From All points below to a specific point above. In this approach the algorithm will start from each point below the hill and try to find a path to a specific point above. To parallelize this approach each processor will start from the points in its region only.

**Approach #3:** From One point below to all points above. In this approach the algorithm will start from a specific point below the hill and try to find a path to all point above the hill. To parallelize this approach each processor will start from the specified point then it will use the End region to determine its destination based on the end region points only.

The evaluation of Hill Climbing technique used only at the current state, it does not suffer from computational space issues, where the source of its computational complexity arises from the time required to explore the problem space. The optimal solution of Hill Climbing technique can be achieved in polynomial time and it is one of the NP-Complete problem that the numbers of local maxima can be the cause of exponential computational time [23]. To address these problems proposed algorithm are applied on message passing interface (MPI) parallel and distributed computing systems with high-performance computing that spread out a single application over many multi-core and multi-processor computers to rapidly complete the task. MPI divide large problems into smaller sub-

problems and assign each of them to different processors in a typically distributed system running concurrently in parallel.

In this research, Proposed Hill Climbing techniques are tested on two methods. First method is sequential that accessed code by a single thread. This means that a single thread can only do code in a specific order, hence it being sequential. Second method is MPI that is a library of routines that can be used to create parallel programs in C, C++, and Fortran77 using commonly available operating system services to create parallel processes and exchange information among these processes, as shown in Fig. 4.

The design process of MPI includes vendors (such as IBM, Intel, TMC, Cray, and Convex), parallel library authors (involved in the development of PVM, and Linda), and application specialists. The final version for the draft standard became available in May of 1994 [8].

MPI is a standardized means of exchanging messages among multiple computers running a parallel program across a distributed memory to improve scalability, performance, multi-core and cluster support, and interoperability with other applications. These programs cannot use any MPI communication routine. The two basic routines are MPI\_Send, to send a message to another process, and MPI\_Recv, to receive a message from another process.

| End Regions   |     |             |     |             |     |             |     |             |    |
|---------------|-----|-------------|-----|-------------|-----|-------------|-----|-------------|----|
| Processor 1   |     | Processor 2 |     | Processor 3 |     | Processor 4 |     | Processor 5 |    |
| 23            | 68  | 789         | 45  | 21          | 32  | 12          | 32  | 32          | 1  |
| 213           | 32  | 12          | 3   | 5           | 65  | 6           | 56  | 5           | 45 |
| 1             | 23  | 56          | 6   | 65          | 48  | 78          | 48  | 65          | 12 |
| 33            | 321 | 12          | 32  | 15          | 35  | 82          | 25  | 25          | 22 |
| 65            | 65  | 84          | 654 | 987         | 12  | 3           | 123 | 45          | 8  |
| 23            | 321 | 45          | 94  | 3           | 321 | 123         | 978 | 56          | 23 |
| 798           | 23  | 546         | 32  | 15          | 32  | 56          | 12  | 32          | 12 |
| 32            | 32  | 12          | 3   | 15          | 32  | 12          | 31  | 32          | 32 |
| Processor 1   |     | Processor 2 |     | Processor 3 |     | Processor 4 |     | Processor 5 |    |
| Start Regions |     |             |     |             |     |             |     |             |    |

Fig. 3. Example of Parallelizing the Matrix.

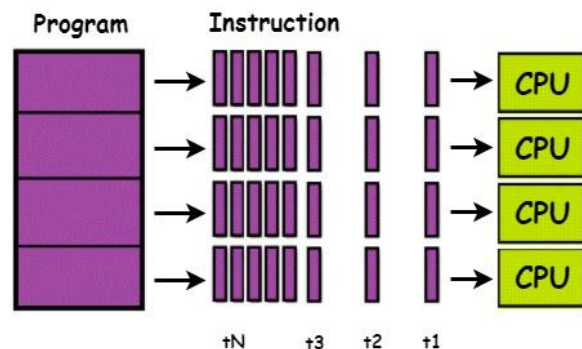


Fig. 4. MPI Parallel Processes.

Proposed algorithm run MPI code in IMAN1, Jordan's first and fastest high-performance Computing resource, funded by JAEC and SESAME. It is available for use by academia and industry in Jordan and the region. In our project, we worked in a Zaina server, an Intel Xeon-based computing cluster with 1G Ethernet interconnection as shown in Table II. The cluster is mainly used for code development, code porting, and synchrotron radiation application purposes. In addition, this cluster is composed of two Dell PowerEdge R710 and five HP ProLiant DL140 G3 servers.

TABLE II. ZAINA TECHNICAL DETAILS

|                    |   |
|--------------------|---|
| Server             | 7 Servers (Two Dell PowerEdge R710 and five HP ProLiant DL140 G3) |
| CPU per server     | Dell (2 X 8 cores Intel Xeon) HP (2 X 4 cores Intel Xeon)         |
| RAM per server     | Dell (16 GB) HP (6 GB)  |
| Total storage (TB) | 1 TB NFS Share  |
| OS                 | Scientific Linux 6.4  |

#### IV. PROPOSED HILL CLIMBING ALGORITHM ANALYSIS

In this section, analysis of the proposed hill climbing algorithm were discussed. Variables that included in all the analysis equation are giving as follows:

Let N = Number of rows in matrix;

Let Paths = Number of scanned paths each time;

Let Points = Number of points in each path;

First; sequential analysis for best paths and parallel analysis for best paths indicate that the proposed algorithm is cost optimal based on the following equation:

##### A. Sequential Analysis for Best Paths

$$TS = \frac{N}{Points} * (Paths * Points) = N * Paths \quad (1)$$

##### B. Parallel Analysis for Best Paths

Let P = Number of Processors

$$Tp = TComm + TComp \quad (2)$$

$$TComm = t_s \left( \frac{N}{P} \right) + t_w \left( \frac{\frac{N}{Points} * (Paths * Points)}{P} \right) = t_s \left( \frac{N}{P} \right) + t_w \left( \frac{N * (Paths)}{P} \right) \quad (3)$$

$$TComp = \left( \frac{N * (Paths)}{P} \right) \quad (4)$$

$$Tp = t_s \left( \frac{N}{P} \right) + t_w \left( \frac{N * (Paths)}{P} \right) + \left( \frac{N * (Paths)}{P} \right) \quad (5)$$

##### C. Total Parallel Overhead

$$T = P \left( t_s \left( \frac{N}{P} \right) + t_w \left( \frac{N * (Paths)}{P} \right) + \left( \frac{N * (Paths)}{P} \right) \right) - (N * Paths) = t_s(N) + t_w(N * Paths) \quad (6)$$

$$Speedup = \frac{N * Paths}{t_s \left( \frac{N}{P} \right) + t_w \left( \frac{N * (Paths)}{P} \right) + \left( \frac{N * (Paths)}{P} \right)} \quad (7)$$

$$Efficiency = \frac{N * Paths}{t_s(N) + t_w(N * Paths) + (N * Paths)} \quad (8)$$

$$Cost = P \left( t_s \left( \frac{N}{P} \right) + t_w \left( \frac{N * (Paths)}{P} \right) + \left( \frac{N * (Paths)}{P} \right) \right) \quad (9)$$

Second; sequential analysis for all to one or one to all and parallel analysis for all to one or one to all indicate that the proposed algorithm is cost optimal based on the following equation.

##### D. Sequential Analysis for All to One or One to All

$$TS = \frac{N}{Points} * \left( \frac{Paths}{2} * Points \right) = N * \frac{Paths}{2} \quad (10)$$

##### E. Parallel Analysis for All to One or One to All

Let P = Number of Processors

$$Tp = TComm + TComp \quad (11)$$

$$TComm = t_s \left( \frac{N}{P} \right) + t_w \left( \frac{\frac{N}{Points} * \left( \frac{Paths}{2} * Points \right)}{P} \right) = t_s \left( \frac{N}{P} \right) + t_w \left( \frac{N * \left( \frac{Paths}{2} \right)}{P} \right) \quad (12)$$

$$TComp = \left( \frac{N * \left( \frac{Paths}{2} \right)}{P} \right) \quad (13)$$

$$Tp = t_s \left( \frac{N}{P} \right) + t_w \left( \frac{N * \left( \frac{Paths}{2} \right)}{P} \right) + \left( \frac{N * \left( \frac{Paths}{2} \right)}{P} \right) \quad (14)$$

##### F. Total Parallel Overhead

$$T = P \left( t_s \left( \frac{N}{P} \right) + t_w \left( \frac{N * \left( \frac{Paths}{2} \right)}{P} \right) + \left( \frac{N * \left( \frac{Paths}{2} \right)}{P} \right) \right) - (N * \frac{Paths}{2}) = t_s(N) + t_w \left( N * \frac{Paths}{2} \right) \quad (15)$$

$$Speedup = \frac{N * \frac{Paths}{2}}{t_s \left( \frac{N}{P} \right) + t_w \left( \frac{N * \left( \frac{Paths}{2} \right)}{P} \right) + \left( \frac{N * \left( \frac{Paths}{2} \right)}{P} \right)} \quad (16)$$

$$Efficiency = \frac{N * \frac{Paths}{2}}{t_s(N) + t_w \left( N * \frac{Paths}{2} \right) + \left( N * \frac{Paths}{2} \right)} \quad (17)$$

$$Cost = P \left( t_s \left( \frac{N}{P} \right) + t_w \left( \frac{N * \left( \frac{Paths}{2} \right)}{P} \right) + \left( \frac{N * \left( \frac{Paths}{2} \right)}{P} \right) \right) \quad (18)$$

#### V. EXPERIMENTS AND RESULTS

This research uses different matrix sizes that contain points that need to go from the current point to the next point in a certain matrix. The cost of moving from the current point to the next point is calculated using the sum of the obstacles between the current point and all the 10 other points. The path is decided depending on the minimum value among the total obstacle weights in each path. The proposed algorithm is tested in sequential and parallel forms coded by MPI. The results are compared in terms of efficiency and speedup ratio.

First, simple hill climbing is used to calculate the time needed to find all the best paths from a specific point to another point with the least cost from that start point. Second, the proposed hill climbing algorithm is utilized to calculate the time required to find the best path from all the points below the

matrix to a specific point above the matrix. Finally, the proposed hill climbing algorithm is adopted to calculate the time needed to find the best path from a specific point below the matrix to all the points above the matrix.

The sequential results of the proposed hill climbing algorithm are tested with various matrix sizes. The algorithm is written in C++. The experimental results are calculated with the 1 core in MPI as shown in Table III.

The MPI results of the proposed hill climbing algorithm are tested using different numbers of cores and matrices. The results are effective and efficient when the number of cores is increased due to the large size of problems that need a high degree of parallelism. Table IV and Fig. 5 show the results for all the best paths. Fig. 6 illustrates all points below a point above, and Fig. 7 presents a point below all points above.

The comparison between MPI results and sequential methods indicates that MPI is always faster and more efficient than sequential methods for different matrix sizes.

Table V presents the calculation results for speedup ratio. Fig. 8 shows the comparison of speedup ratio for all best path results. Fig. 9 reveals the speedup ratio for a point below to all points above. Fig. 10 illustrates the speedup ratio for all points below to a point above.

TABLE III. SEQUENTIAL RUN TIME RESULTS

| One CPU     | From All Points Below to Unknown Point Above | From All Points Below to One Point Above | From One Point Below to All Points Above |
|-------------|--|--|--|
| 100 * 100   | 3.200 s                                      | 3.150 s                                  | 3.300 s                                  |
| 500 * 500   | 130.230 s                                    | 132.230 s                                | 134.230 s                                |
| 1000 * 1000 | 593.320 s                                    | 598.120 s                                | 601.300 s                                |

TABLE IV. MPI RUN TIME RESULTS

| 2 CPUs      | From All Points Below to Unknown Point Above | From All Points Below to One Point Above | From One Point Below to All Points Above |
|-------------|--|--|--|
| 100 * 100   | 2.910  | 2.890                                    | 3.210                                    |
| 500 * 500   | 105.600                                      | 106.900                                  | 106.230                                  |
| 1000 * 1000 | 342.250                                      | 341.530                                  | 342.680                                  |
| 4 CPUs      | From All Points Below to Unknown Point Above | From All Points Below to One Point Above | From One Point Below to All Points Above |
| 100 * 100   | 2.32   | 2.13                                     | 2.23                                     |
| 500 * 500   | 39.23  | 38.32                                    | 39.65                                    |
| 1000 * 1000 | 160.2  | 158.7                                    | 158.32                                   |
| 8 CPUs      | From All Points Below to Unknown Point Above | From All Points Below to One Point Above | From One Point Below to All Points Above |
| 100 * 100   | 1.51   | 1.35                                     | 1.56                                     |
| 500 * 500   | 18.5   | 18.23                                    | 18.9                                     |
| 1000 * 1000 | 76.45  | 75.32                                    | 76.81                                    |
| 16 CPUs     | From All Points Below to                     | From All Points Below to                 | From One Point Below to All              |

|             | Unknown Point Above                          | One Point Above                          | Points Above                             |
|-------------|--|--|--|
| 100 * 100   | 1.12   | 1.23                                     | 1.11                                     |
| 500 * 500   | 12.23  | 13.89                                    | 12.56                                    |
| 1000 * 1000 | 39.56  | 38.56                                    | 39.15                                    |
| 32 CPUs     | From All Points Below to Unknown Point Above | From All Points Below to One Point Above | From One Point Below to All Points Above |
| 100 * 100   | 0.927  | 0.978                                    | 0.968                                    |
| 500 * 500   | 6.2  | 6.3                                      | 6.51                                     |
| 1000 * 1000 | 19.65  | 20.3                                     | 20.3                                     |
| 64 CPUs     | From All Points Below to Unknown Point Above | From All Points Below to One Point Above | From One Point Below to All Points Above |
| 100 * 100   | 0.748  | 0.789                                    | 0.868                                    |
| 500 * 500   | 3.5  | 3.6                                      | 3.78                                     |
| 1000 * 1000 | 11.2   | 11.9                                     | 12.3                                     |
| 100 CPU     | From All Points Below to Unknown Point Above | From All Points Below to One Point Above | From One Point Below to All Points Above |
| 100 * 100   | 0.592  | 0.512                                    | 0.54                                     |
| 500 * 500   | 1.9  | 1.8                                      | 1.7                                      |
| 1000 * 1000 | 6.2  | 6.9                                      | 6.5                                      |

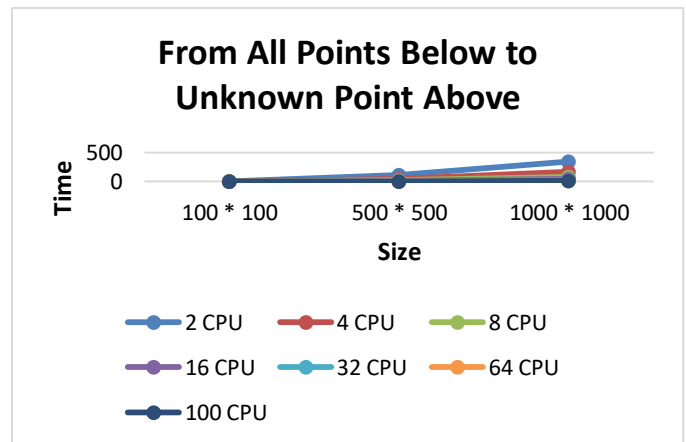


Fig. 5. From All Points below to unknown Point above.

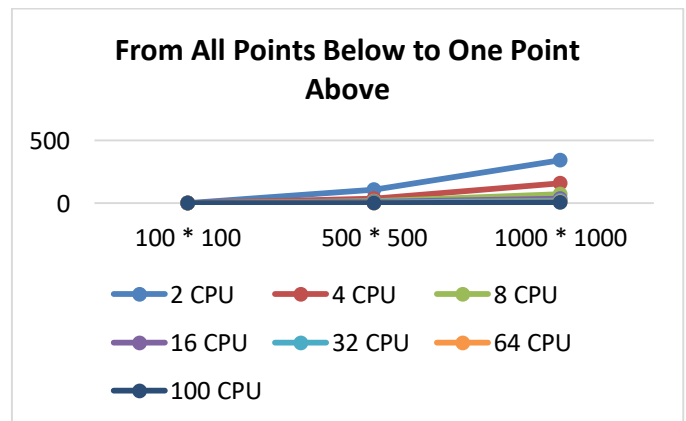


Fig. 6. From All Points below to One Point above Plotting.

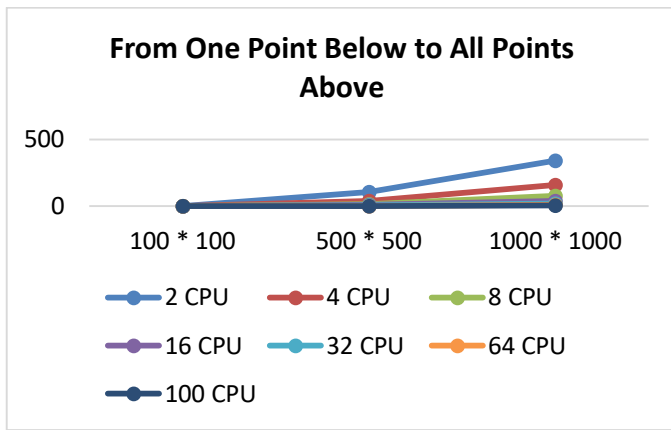


Fig. 7. From One Point below to All Points above Plotting.

TABLE V. SPEEDUP RESULTS

| 2 CPUs      | From All Points Below to Unknown Point Above | From All Points Below to One Point Above | From One Point Below to All Points Above |
|-------------|--|--|--|
| 100 * 100   | 1.100  | 1.090                                    | 1.028                                    |
| 500 * 500   | 1.233  | 1.237                                    | 1.264                                    |
| 1000 * 1000 | 1.734  | 1.751                                    | 1.755                                    |
| 4 CPUs      | From All Points Below to Unknown Point Above | From All Points Below to One Point Above | From One Point Below to All Points Above |
| 100 * 100   | 1.379  | 1.479                                    | 1.480                                    |
| 500 * 500   | 3.320  | 3.451                                    | 3.385                                    |
| 1000 * 1000 | 3.704  | 3.769                                    | 3.798                                    |
| 8 CPUs      | From All Points Below to Unknown Point Above | From All Points Below to One Point Above | From One Point Below to All Points Above |
| 100 * 100   | 2.119  | 2.333                                    | 2.115                                    |
| 500 * 500   | 7.039  | 7.253                                    | 7.102                                    |
| 1000 * 1000 | 7.761  | 7.941                                    | 7.828                                    |
| 16 CPUs     | From All Points Below to Unknown Point Above | From All Points Below to One Point Above | From One Point Below to All Points Above |
| 100 * 100   | 2.857  | 2.561                                    | 2.973                                    |
| 500 * 500   | 10.648                                       | 9.520                                    | 10.687                                   |
| 1000 * 1000 | 14.998                                       | 15.511                                   | 15.359                                   |
| 32 CPUs     | From All Points Below to Unknown Point Above | From All Points Below to One Point Above | From One Point Below to All Points Above |
| 100 * 100   | 3.452  | 3.221                                    | 3.409                                    |
| 500 * 500   | 21.005                                       | 20.989                                   | 20.619                                   |
| 1000 * 1000 | 30.194                                       | 29.464                                   | 29.621                                   |
| 64 CPUs     | From All Points Below to Unknown Point Above | From All Points Below to One Point Above | From One Point Below to All Points Above |
| 100 * 100   | 4.278  | 3.992                                    | 3.802                                    |
| 500 * 500   | 37.209                                       | 36.731                                   | 35.511                                   |

| 1000 * 1000 | 52.975                                       | 50.262                                   | 48.886                                   |
|-------------|--|--|--|
| 100 CPU     | From All Points Below to Unknown Point Above | From All Points Below to One Point Above | From One Point Below to All Points Above |
| 100 * 100   | 5.405  | 6.152                                    | 6.111                                    |
| 500 * 500   | 68.542                                       | 73.461                                   | 78.959                                   |
| 1000 * 1000 | 95.697                                       | 86.684                                   | 92.508                                   |

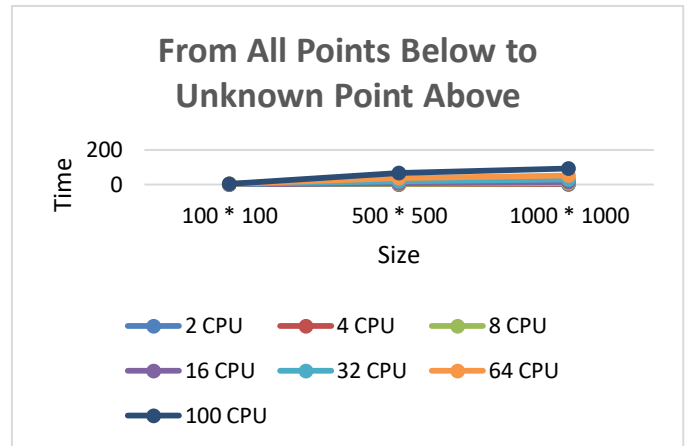


Fig. 8. From All Points below to unknown Point above Speedup Plotting.

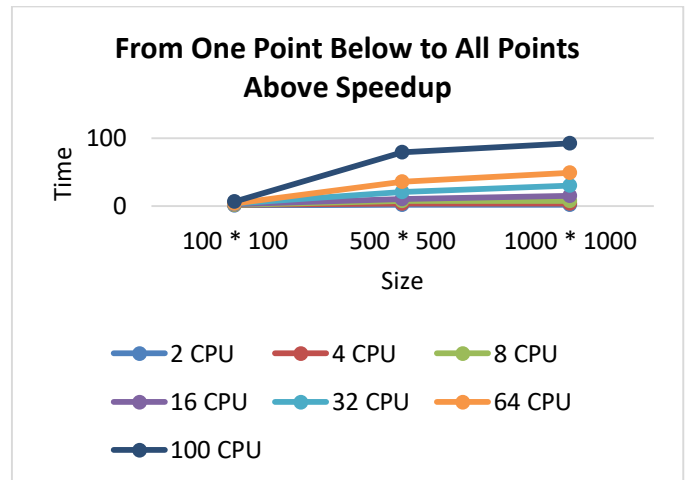


Fig. 9. From One Point below to All Points above Speedup Plotting.

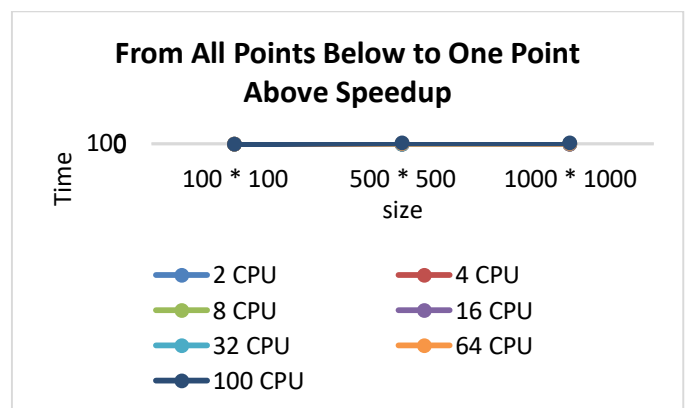


Fig. 10. From All Points below to One Point above Speedup Plotting.

Table VI shows the calculation results for parallel efficiency. Fig. 11 presents the comparison of parallel efficiency for all best path results. Fig. 12 reveals the parallel efficiency for all points below to unknown point above. Fig. 13 illustrates the parallel efficiency for a point below to all points above.

TABLE VI. EFFICIENCY RESULTS

|                | From All Points Below to Unknown Point Above | From All Points Below to One Point Above | From One Point Below to All Points Above |
|----------------|--|--|--|
| <b>2 CPUs</b>  |  |  |  |
| 100 * 100      | 0.550  | 0.545                                    | 0.514                                    |
| 500 * 500      | 0.617  | 0.618                                    | 0.632                                    |
| 1000 * 1000    | 0.867  | 0.876                                    | 0.877                                    |
| <b>4 CPUs</b>  |  |  |  |
| 100 * 100      | 0.345  | 0.370                                    | 0.370                                    |
| 500 * 500      | 0.830  | 0.863                                    | 0.846                                    |
| 1000 * 1000    | 0.926  | 0.942                                    | 0.950                                    |
| <b>8 CPUs</b>  |  |  |  |
| 100 * 100      | 0.265  | 0.292                                    | 0.264                                    |
| 500 * 500      | 0.880  | 0.907                                    | 0.888                                    |
| 1000 * 1000    | 0.970  | 0.993                                    | 0.979                                    |
| <b>16 CPUs</b> |  |  |  |
| 100 * 100      | 0.179  | 0.160                                    | 0.186                                    |
| 500 * 500      | 0.666  | 0.595                                    | 0.668                                    |
| 1000 * 1000    | 0.937  | 0.969                                    | 0.960                                    |
| <b>32 CPUs</b> |  |  |  |
| 100 * 100      | 0.108  | 0.101                                    | 0.107                                    |
| 500 * 500      | 0.656  | 0.656                                    | 0.644                                    |
| 1000 * 1000    | 0.944  | 0.921                                    | 0.926                                    |
| <b>64 CPUs</b> |  |  |  |
| 100 * 100      | 0.067  | 0.062                                    | 0.059                                    |
| 500 * 500      | 0.581  | 0.574                                    | 0.555                                    |
| 1000 * 1000    | 0.828  | 0.785                                    | 0.764                                    |
| <b>100 CPU</b> |  |  |  |
| 100 * 100      | 0.054  | 0.062                                    | 0.061                                    |
| 500 * 500      | 0.685  | 0.735                                    | 0.790                                    |
| 1000 * 1000    | 0.957  | 0.867                                    | 0.925                                    |

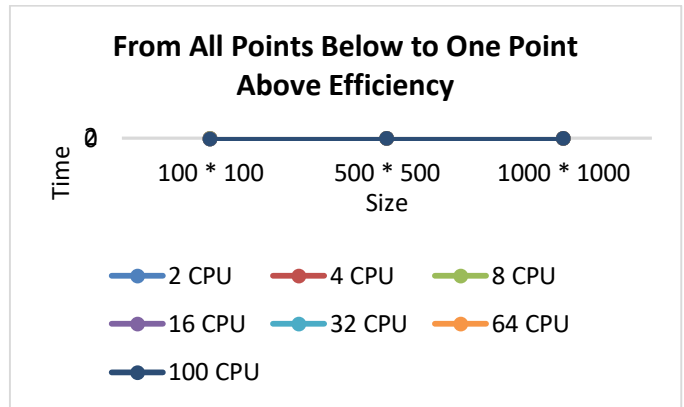


Fig. 11. From All Points below to One Point above Efficiency Plotting.

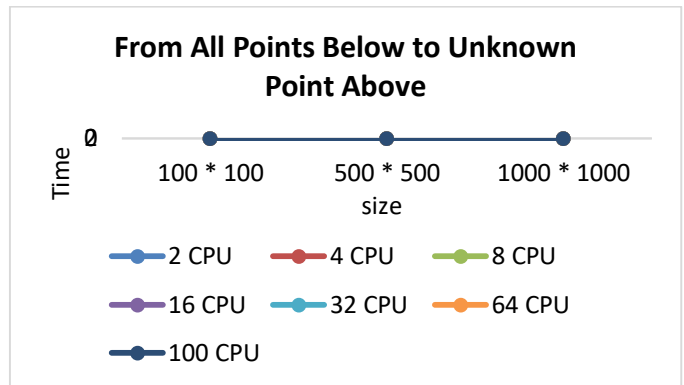


Fig. 12. From All Points below to unknown Point above Efficiency Plotting.

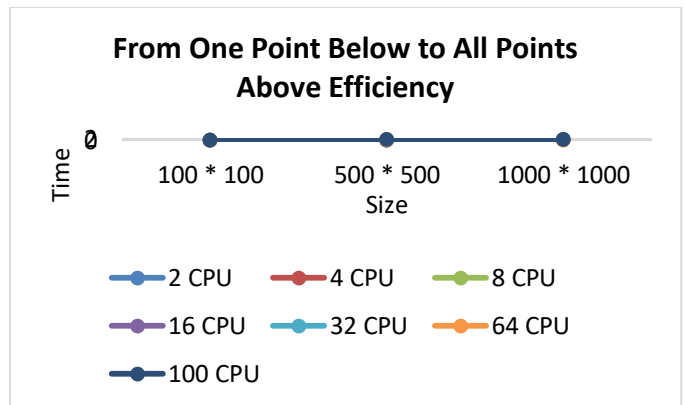


Fig. 13. From One Point below to All Points above Efficiency Plotting.

## VI. CONCLUSION

Heuristic search is a search process that uses domain knowledge in heuristic rules or procedures to direct the progress of a search algorithm. Hill climbing is a heuristic search technique for solving certain mathematical optimization problems in the field of artificial intelligence. In this technique, starting with a suboptimal solution is compared to starting from the base of the hill, and improving the solution is compared to walking up the hill. The optimal solution of the hill climbing technique can be achieved in polynomial time and is an NP-complete problem in which the numbers of local maxima can lead to an exponential increase in computational time. To address these problems, the proposed hill climbing algorithm

based on the local optimal solution is applied to the message passing interface, which is a library of routines that can be used to create parallel programs by using commonly available operating system services to create parallel processes and exchange information among these processes. Experimental results show that parallel hill climbing outperforms sequential methods.

This research uses different matrix sizes that contain points that need to go from the current point to the next point in a certain matrix. The cost of moving from the current point to the next point is calculated using the sum of the obstacles between the current point and all the 10 other points. The path is decided depending on the minimum value among the total obstacle weights in each path. The proposed algorithm is tested in sequential and parallel forms coded by MPI. The results are compared in terms of efficiency and speedup ratio.

The comparison between MPI results and sequential methods indicates that MPI is always faster and more efficient than sequential methods for different matrix sizes. Fig. 7, 8, and 9 show the comparison results for sizes 100×100, 500×500, and 1000×1000, respectively. The MPI outperforms the sequential methods; thus, the research goal is achieved.

#### REFERENCES

- [1] Snir, M., Otto, S., Huss-Lederman, S., Dongarra, J., & Walker, D. (1998). MPI--the Complete Reference: The MPI core (Vol. 1). MIT press.
- [2] Chira, C., Horvath, D., & Dumitrescu, D. (2011). Hill-Climbing search and diversification within an evolutionary approach to protein structure prediction. *BioData mining*, 4(1), 23.
- [3] Selman, B., & Gomes, C. P. (2006). Hill-climbing search. *Encyclopedia of Cognitive Science*, 81, 82.
- [4] Cook, C. M., Rosenfeld, A., & Aronson, A. R. (1976). Grammatical inference by hill climbing. *Information Sciences: an International Journal*, 10(2), 59-80.
- [5] Apter, M. J. (1970). *The Computer Simulation of behaviour*. London: Hutchinson. Allgemeinverständliche, inzwischen etwas veraltete Darstellung der Simulationsmethodik mit Diskussion von Anwendungen aus Bereichen des Lernens, des Problemlösens, des Mustererkennens, der Sprache und der Persönlichkeitstheorie bis hin zum Problem des Bewußtseins.
- [6] Masadeh, R., Mahafzah, B. A., & Sharieh, A. (2019). Sea Lion optimization algorithm. *International Journal of Advanced Computer Science and Applications*, 10(5), 388-395.
- [7] Goswami, S., Das, A. K., Guha, P., Tarafdar, A., Chakraborty, S., Chakrabarti, A., & Chakraborty, B. (2017). An approach of feature selection using graph-theoretic heuristic and hill climbing. *Pattern Analysis and Applications*, 1-17.
- [8] Hansen, E. A., & Zhou, R. (2007). Anytime heuristic search. *Journal of Artificial Intelligence Research*, 28, 267-297.
- [9] Masadeh, R., Sharieh, A., & Sliet, A. (2017). Grey wolf optimization applied to the maximum flow problem. *International Journal of Advanced and Applied Sciences*, 4(7), 95-100.
- [10] Fox, Geoffrey C., Steve W. Otto, and Anthony JG Hey. "Matrix algorithms on a hypercube I: Matrix multiplication." *Parallel computing* 4.1 (1987): 17-31.
- [11] Masadeh, R., Alzaqebah, A., Smadi, B., Masadeh, E. (2020). Parallel Whale Optimization Algorithm for Maximum Flow Problem. *Modern Applied Science*, 14(3), 30-44.
- [12] Gropp, W. D., Gropp, W., Lusk, E., & Skjellum, A. (1999). *Using MPI: portable parallel programming with the message-passing interface* (Vol. 1). MIT press.
- [13] Rashid, M. H., & Tao, L. (2017, October). Parallel Combinatorial Optimization Heuristics with GPUs. In *Computer Science and Intelligent Controls (ISCSIC), 2017 International Symposium on* (pp. 118-123). IEEE.
- [14] Jiang, W., Wang, P., Wang, J., & Wang, L. (2017). Optimal power allocation for parallel grid-connected inverters based on lagrangian function method. *Chinese Journal of Electrical Engineering*, 3(3), 68-76.
- [15] Kim, J., Kim, G., & Park, Y. I. (2018). Component Sizing of Parallel Hybrid Electric Vehicle Using Optimal Search Algorithm. *International Journal of Automotive Technology*, 19(4), 743-749.
- [16] Jiang, W., Wang, P., Wang, J., & Wang, L. (2017). Optimal power allocation for parallel grid-connected inverters based on lagrangian function method. *Chinese Journal of Electrical Engineering*, 3(3), 68-76.
- [17] Chen, B., Mo, W., Chattopadhyay, I., & Lipson, H. (2018). Autostacker: an Automatic Evolutionary Hierarchical Machine Learning System.
- [18] Mincu, R. S., & Popa, A. (2018, July). Heuristic Algorithms for the Min-Max Edge 2-Coloring Problem. In *International Computing and Combinatorics Conference* (pp. 662-674). Springer, Cham.
- [19] Harman, M., & McMinn, P. (2007, July). A theoretical & empirical analysis of evolutionary testing and hill climbing for structural test data generation. In *Proceedings of the 2007 international symposium on Software testing and analysis* (pp. 73-83). ACM.
- [20] Gámez, J. A., Mateo, J. L., & Puerta, J. M. (2011). Learning Bayesian networks by hill climbing: efficient methods based on progressive restriction of the neighborhood. *Data Mining and Knowledge Discovery*, 22(1-2), 106-148.
- [21] Khari, M., & Kumar, P. (2017). Empirical Evaluation of Hill Climbing Algorithm. *International Journal of Applied Metaheuristic Computing (IJAMC)*, 8(4), 27-40.
- [22] Chan, W. K. V., D'Ambrogio, A., Zacharewicz, G., Mustafee, N., Wainer, G., & Page, E. A Global and Local Search Approach to Quay Crane Scheduling Problem.
- [23] Nicolau, M., & McDermod, J. (2017). Late-Acceptance and Step-Counting Hill-Climbing GP for Anomaly Detection.

# High-Speed and Secure Elliptic Curve Cryptosystem for Multimedia Applications

Mohammad Alkhatib

College of Computer and Information Sciences  
Al Imam Mohammad Ibn Saud Islamic University (IMSIU)  
Riyadh, Saudi Arabia

**Abstract**—The last few years witnessed a rapid increase in the use of multimedia applications, which led to an explosion in the amount of data sent over communication networks. Therefore, it has become necessary to find an effective security solution that preserves the confidentiality of such enormous amount of data sent through unsecure network channels and, at the same time, meets the performance requirements for applications that process the data. This research introduces a high-speed and secure elliptic curve cryptosystem (ECC) appropriate for multimedia security. The proposed ECC improves the performance of data encryption process by accelerating the scalar multiplication operation, while strengthening the immunity of the cryptosystem against side channel attacks. The speed of the encryption process has been increased via the parallel implementation of ECC computations in both the upper scalar multiplication level and the lower point operations level. To accomplish this, modified version of the Right to Left binary algorithm as well as eight parallel multipliers (PM) were used to allow parallel implementation for point doubling and addition. Moreover, projective coordinates systems were used to remove the time-consuming inversion operation. The current 8-PM Montgomery ECC achieves higher performance level compared to previous ECC implementations, and can reduce the risk of side channel attacks. In addition, current research work provides performance and resources-consumption analysis for Weierstrass and Montgomery elliptic curve representations over prime field. However, the proposed ECC implementation consumes more resources. Presented ECCs were implemented using VHDL, and synthesized using the Xilinx tool with target FPGA.

**Keywords**—*Elliptic curves cryptosystem; performance; binary method; projective coordinates; security applications*

## LIST OF NOTATIONS

|          |  |
|----------|--|
| ECC      | Elliptic Curve Cryptosystem                        |
| RLA      | Right to Lift Algorithm                            |
| SPA      | Simple Power Attack                                |
| STA      | Simple Time Attack                                 |
| SM       | Sequential Multiplication                          |
| SA       | Sequential Addition                                |
| PM       | Parallel Multiplier                                |
| TSM      | Time-consumption for one sequential multiplication |
| $T_M$    | Time-consumption for one multiplication operation  |
| $T_{KP}$ | Time-consumption for scalar multiplication         |
| GF       | Galious Field                                      |
| NAF      | Non-Adjacent-Form                                  |

## I. INTRODUCTION

Elliptic Curve Crypto-system (ECC) is a type of public key cryptosystems that depend on the discrete logarithm problem for elliptic curves. It has been introduced by Miller and Koblitz in 1985 [1,2]. Since that date, it has been widely used in many security applications due to its reliability and efficiency. By using much shorter key length, ECC can provide equivalent security level to that obtained by other asymmetric ciphers such with consuming less time and resources, which made it very efficient for multimedia applications that need to provide the security services for huge amount of data in the shortest possible period of time and, of course, with the least amount of resources consumed.

A variety of ECC representations over  $GF(p)$  and  $GF(2^n)$  were presented and used for different elliptic curves applications. First, ECC represents the plaintext as a point on an elliptic curve. Then, it encrypts the plaintext by performing a number of arithmetic operations over finite fields. ECC computations can be categorized into upper and lower layers. The upper layer's computations are mainly point doubling and point addition operations, which are performed by the scalar multiplication operation. It is worth mentioning here that the scalar multiplication is the key operation in ECC encryption process. On the other hand, lower level of computations includes addition, multiplication, and inversion operations. The latter is the most time-consuming operation [3].

Previous research works focused on improving the performance and security of ECC encryption. These works studied several possible techniques to accelerate the encryption process by speeding up scalar multiplication operation as well as increasing the cryptosystem immunity against side channel attacks such as simple time (STA) and simple power (SPA) attacks. The major performance improvement techniques include the use of projective coordinates to avoid the costly inversion operation and the parallel implementation of ECC arithmetic computations, especially in the lower level [1-5].

The current research utilizes both the use of projective coordinates and the inherited parallelism in ECC computations, and perform these computations in parallel for both upper and lower computational layers. This is achieved by using parallel hardware components, which are multipliers and adders. In addition to performing the lower level computations in parallel, this study implements the upper layer's operations in parallel to achieve higher speed for encryption process. In particular, proposed ECC performs the two main operations (point



doubling and point addition) in scalar multiplication in parallel. This plays crucial role in increasing the speed of scalar multiplication and thus ECC encryption.

In order to enable parallel implementation of point doubling and point addition operations, the current ECC uses modified version of the Right to Left binary algorithm (RLA), which is widely used to perform scalar multiplication. The modified RLA has the ability to apply both operations in parallel once required. This is obtained by assigning an independent variable to save the point operations result of certain iteration of the RLA and use it for point addition in the next iteration. To accomplish optimum performance, proposed ECC uses eight parallel multipliers (8-PM).

The parallel ECC implementation can be realized by using parallel hardware components for hardware implementations, and the known multithreading technique for software implementations [6-8]. This research uses hardware implementations because they are faster and more secure than software implementation for ECC.

Several projective coordinates systems are studied and implemented in this research to achieve greater speed for ECC encryption. Moreover, the main ECC representations over  $GF(p)$  are implemented to study their characteristics in terms of performance and security. It should be mentioned that proposed ECC implementation strengthen the security against SPA and STA.

The core contribution of this research is represented by developing a novel and fast ECC using modified version of the RLA. The proposed ECCs utilize the maximum parallelization levels for both Montgomery and Weierstrass curves to achieve optimum performance.

Results showed that the use of proposed RLA to implement ECC operations improved the speed of the encryption considerably. Such ECC designs and implementations are highly recommended for securing applications that need a high-speed ECC, such as multimedia applications.

The remaining parts of this article are the background and related works, ECCs computations and architectures, Results and analysis, and conclusion.

## II. BACKGROUND AND RELATED WORKS

ECC is a public key cryptography that can be used to provide different security services including confidentiality of data transmitted over communication networks. ECC supersedes other public key cryptosystems because it can provide equivalent security level with using much shorter key size, which represents considerable improvement in terms of performance and resources consumption [1,6].

Scalar multiplication is the main operation in ECC encryption, and it consists of two operations; point doubling and point addition. Several algorithms were used to perform scalar multiplication. It can be noticed that scalar multiplication algorithms use iterative approach to perform point operations, which leads in the end to converting the plaintext to the ciphertext [1, 7-9].

ECC computations performed during scalar multiplication are called upper level computations. The Binary method, the Non Adjacent Form (NAF), and the Montgomery ladder are the main algorithms used to apply ECC scalar multiplication. These algorithms vary in performance and security. The RLA is a form of the binary algorithm which is intensively used for ECC encryption due to its security advantages and the ability to withstand against side channel attacks [10-12].

The RLA, scans the binary bits of the key and always performs point doubling operation regardless the value of the key bit. If the value of the bit is one, the point addition operation is performed as well. Therefore, the RLA assumes that, on average, the point addition operation is executed half times of point doubling during the entire scalar multiplication. However, previous ECC that uses standard RLA implements point operations sequentially, which increases the critical path delay of scalar multiplication [13, 14].

Point doubling is the dominating operation in scalar multiplication. In this operation, the elliptic curve point is added to itself, while in point addition, two different points are added. Computations performed within each point doubling and addition operation are called lower level of computations, which represent arithmetic operations over finite fields. Those operations are mainly modular multiplication, addition, and inversion operations [2-6].

The inversion is the most time-consuming operation in elliptic curve cryptography. Previous researches suggested to use projective coordinates instead of affine coordinates to eliminate the inversion operation by converting it to consecutive multiplications. This played important role in reducing the time delay of scalar multiplication operation. Several projective coordinates systems were presented including homogenous, López-Dahab, and Jacobean coordinates systems [5, 6, 14-18].

Since its first introduction in 1985, different elliptic curves forms over  $GF(p)$  were presented [1-3, 16-20]. Some curves such Montgomery, and Tripling Oriented curves have less computational complexity in comparison with other curves such as Weierstrass and Binary Edwards. Thus, particular types of elliptic curves can reduce the time delay for the scalar multiplication operations. This made the curves with lower computational complexity more suitable for security applications that require high-speed cryptosystem [21-23].

The majority of researches in this field focused on hardware implementations of ECC since they are faster and more secure than software implementations. Researchers studied the effect of using the different projective coordinates systems with the main forms of elliptic curves. The majority of previous research works studied the performance and resources-consumption of Weierstrass curve over  $GF(p)$  [7-9].

Researchers found that the homogenous projection system  $(X/Z, Y/Z)$  accomplished the highest performance levels when implemented with many  $GF(p)$  elliptic curve representations; including [8, 11, 12-21]. On the other side Jacobean and López-Dahab coordinates showed better performance when used with other types of elliptic curves. The performance or speed of ECC is usually estimated by the number of sequential

multiplication (SM) and addition (SA) levels consumed by point doubling and point addition operations [7, 9-10].

In addition to using projective coordinates, researchers used specific hardware components such as Montgomery multiplier to increase the speed of ECC operations [8].

However, the majority of previous researches studied the use of usual serial implementation of ECC computations, in which only one multiplication or addition operation is performed in every level of computations. Although, serial ECC implementation consumes the least possible resources, it cannot satisfy performance requirements for applications that need to provide the confidentiality of many data streams simultaneously as in multimedia applications [20-23]. A high-performance cryptosystem is in demand to satisfy the requirements of many emerging cloud services technologies such as the frame work presented in [31], which aims to provide secure environment for cloud infrastructures that allows to serve clients in efficient and secure way.

Another study developed hardware implementations for Weierstrass ECC over  $GF(2^n)$ . The presented fast ECC is suitable for smart card implementations [25]. In [26], researchers introduced Weierstrass ECC, which uses the López-Dahab projective coordinates to eliminate the costly inversion operation.

In [24], authors developed Weierstrass ECC design that can support both types of finite fields. However, the results obtained from proposed ECC showed that it needs extra memory resources.

Recent research studies proposed parallel implementations of ECC computations to increase the speed of encryption process. Actually, this technique improved ECC performance, but with consuming more resources [10, 13-23]. It should be highlighted that the inherited parallelism in elliptic curve computations make it possible to perform lower level operations in parallel manner, which shortens the time delay of the scalar multiplication [19].

In addition to improving the performance, parallel ECC implementation can strengthen the security of the cryptosystem against the STA and SPA. Furthermore, using parallel ECC designs considerably enhance the area $\times$ time-consumption<sup>2</sup> (AT<sup>2</sup>) factor [20-23].

Authors in [10] developed parallel hardware designs for Weierstrass ECC over  $GF(p)$ . Homogenous projective coordinates were used to remove inversion operation. The ECC implementation with four parallel multipliers (4-PM) consumed the least time for the encryption process. However, with consuming more resources compared to usual serial ECC implementation. Similar parallel ECC implementation over  $GF(2^n)$  was presented in [14]. This study found that using Jacobean projection with Weierstrass curve obtained shorter time-delay compared to other projective coordinates systems.

Another ECC implementation that uses Weierstrass curve over  $GF(2^n)$  was proposed in [15]. Although it provided considerable trade-off between performance and power consumption, the presented ECC implementation has low

system utilization level, which is necessary for efficient ECC [10].

Many ECC design introduced in previous researches [10, 12, 14, 15] worked on the standard Weierstrass elliptic curve representation, while there are many elliptic curve representations that have not been sufficiently studied. Furthermore, parallel ECC implementations in the aforementioned studies consume additional resources and area. Therefore, they are not appropriate for applications with limited resources.

Other research works provided different ECC design choices by using variable degree of parallelism to mitigate the problem of resources consumption. In other words, researchers studied and implemented all possible ECC design schemes that provided significant trade-off between performance and resources consumption. These research efforts aim to introduce a variety of ECC designs that can satisfy the requirements of several security applications in terms of performance and resources-consumption [17-23].

A number of parallel ECC design schemas over  $GF(p)$  were introduced in [17]. Proposed designs used variable degrees of parallelism for ECC computations. Experimental results showed that the 4-PM design achieved the shortest time delay for ECC point addition, which is estimated by four SMs. Other researchers [19] studied the use of variable degrees of parallel designs to increase the speed of point doubling operation. Experiments of both researches found that Weierstrass ECC accomplishes the highest performance level when implemented using four PMs and homogenous coordinates system. Other presented design schemes in [17, 19] provide important trade-off between area and performance, which could be useful for a variety of elliptic curves applications.

Few research works investigated the parallel implementation of ECC forms. In [22] authors analyzed the performance and resources-consumption levels of Montgomery elliptic curve over  $GF(p)$  when implemented using different projective coordinates and parallel hardware designs. A set of design choices for Montgomery ECC were presented. Experimental results illustrated that the 2-PM design accomplished the best trade-off between area and performance, where the 4-PM ECC obtained the highest speed for encryption process.

Similar research works studied the characteristics of Tripling Oriented ECC when implemented using different degrees of parallelism [23].

Authors in [20] and [21] used similar methodology to study the performance and cost features of Binary Edward and Edward elliptic curves respectively. Experiments showed that the greatest speed level of Edward ECC can be reached using the 5-PM design. On the other side, Binary Edward ECC accomplished the shortest time delay when implemented using the 7-PM design. Homogenous projective coordinates system represented best choice for both forms since it involves less arithmetic computations than other projections.

The vast majority of the previous studies used the known Binary method to apply the ECC scaler multiplication. It should be mentioned here that recent few researches investigated the use of other algorithms to process ECC point operations.

Authors in [29] used the NAF algorithm and parallel hardware designs the perform ECC upper level of computations. The results of that research proved that the use of NAF algorithm to perform Montgomery ECC operations can reduce the time delay of the encryption process considerable in comparison with usual Binary method. These outcomes are supported by the fact that NAF algorithm requires performing less number of point addition operations during the scaler multiplication.

Another ECC implementation that uses NAF algorithm was reported in [30]. Researchers studied possible improvement for both Edward and Binary Edward ECCs using parallel hardware designs as well as the NAF algorithm for scaler multiplication. It was found that NAF algorithm can enhance the performance of the encryption process compared to parallel ECC implementations using Binary method. However, the Montgomery ECC presented in [29] remains ahead of the later ECCs developed by [30] in terms of performance.

In [32], researchers presented secure ECC suitable for compact devices by using a modified addition formula for usual RLA. Although proposed method needs less space complexity, the time delay of the encryption process was increased because of using affine coordinates. This makes it unable to satisfy the performance requirements for many web applications. Another promising research work tried to accelerate ECC computations using a modified approach for the Montgomery ladder algorithm [33]. However, the proposed algorithm did not utilize inherited parallelism in the different levels of computations in the encryption process, which is essential to develop high speed crypto processor. Moreover, ECCs presented in [32-33] are vulnerable to side channel attacks.

The research works toward improving the speed and security of ECC needs to give more attention to the development of fast and secure algorithms that benefit from the fact that encryption's computations can be implemented in parallel in both the upper and lower levels of operations.

The current research work proposes a modified version of the RLA that enables parallel implementation of the upper level of computations for ECC represented by point addition and point doubling operations. In addition, proposed ECC implementations utilize the inherited parallelism of the lower level of computations for ECC represented by arithmetic computations over  $GF(p)$ . This research analyzes the performance and resources consumption level of the major types of elliptic curves, which are Weierstrass, Edward, and Montgomery curves. The use of different projective coordinates systems is also investigated to find the most suitable ECC design for applications that need high-performance cryptosystem such as multimedia applications.

The next section presents the research methodology followed in this study to obtain the research outcomes.

### III. RESEARCH METHODOLOGY AND METHODS

The research methodology used to conduct this study is depicted in Fig. 1.

It can be noticed from Fig. 1 that the research methodology of the current study is divided into three phases, as follows:

- Preparation of elliptic curve computational schemes,
- Improving the speed of ECC's lower level of computations, and
- Improving the speed of ECC's scaler multiplication (upper level of computations).

In the first stage, the main research efforts focused on studying the key ECC representations over  $GF(p)$ . According to security, and computational complexity features, three forms of elliptic curves were chosen for implementation in the current research. These forms are Weierstrass, Edward, and Montgomery curves.

In addition to selecting the elliptic curves forms, the use of different projective coordinates systems was investigated. In particular, the computations of point doubling and point addition operations for each ECC form were performed using different projection. The use of projective coordinate avoids the need for the time-consuming inversion operation. Inversion is converted to a number of multiplications. This contributes in reducing the time delay of the lower level of computations. The projective coordinates implemented in this study are homogenous, López-Dahab, and Jacobean systems.

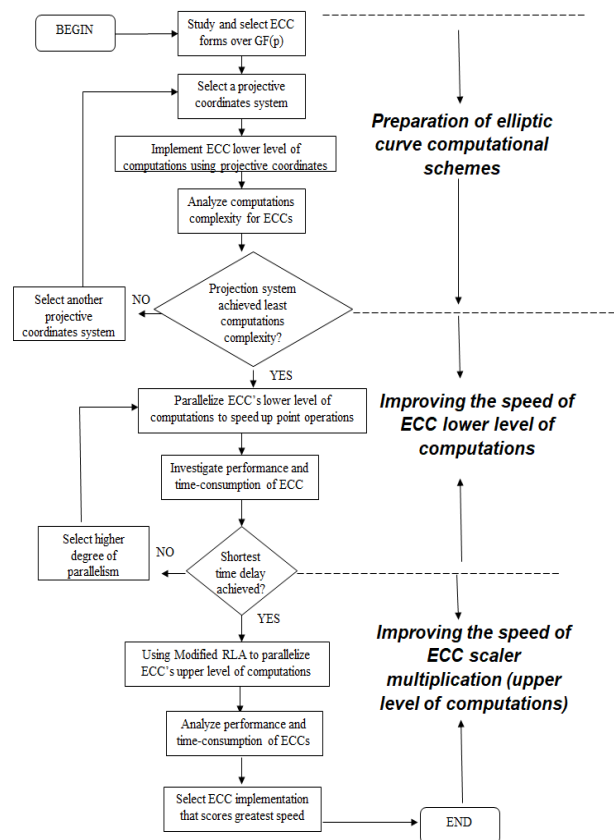


Fig. 1. Research Methodology.

The last step in the first stage involves analyzing the performance or time consumption level for the ECC computations. Such analysis aims to find the most efficient projective ECC that has the least computational complexity. This is essential for the process of developing a high-speed ECC for different security applications.

At the second stage of this work, researchers performed ECC's lower level of computations using different degrees of parallelism in order to reduce the time delay for finite fields arithmetic included in each point doubling and addition. The ECC design or computational scheme that gives the shortest possible time delay for each ECC form is then implemented using parallel hardware designs. Performance and resources-consumption features for ECCs are thoroughly analyzed and studied.

It should be mentioned that the performance of proposed ECCs is estimated using the number of sequential multiplication levels consumed during point doubling or point addition operation. On the other side, resources consumption is estimated by the number of parallel hardware components needed for each design.

In the final stage of this research, a modified version of RLA is used to apply scalar multiplication. This modification allows to perform point addition in parallel to point doubling via saving the results of the previous RLA's iteration in a separate variable, and then it is added to the current elliptic curve point, while performing the point doubling computations.

In other words, researchers parallelized the ECC's upper level of computations to accomplish the maximum speed for encryption process. Such ECC implementation increases the performance but needs more resources to apply ECC computations.

Proposed ECCs will be first studied theoretically by observing the parallel computational schemes and estimating the time and resources consumption levels. Then, presented ECCs are implemented using the Xilinx tool to find out, more accurately, the time and resources consumed for encryption process.

In addition to performance, other factors that are important to determine the efficiency and appropriateness of ECC are investigated. Those factors include the area  $\times$  time (AT), area  $\times$  time<sup>2</sup> (AT<sup>2</sup>), and hardware/system utilization.

The next section of this research illustrates the computations and hardware designs for proposed ECCs.

#### IV. COMPUTATIONAL SCHEMES AND DESIGNS FOR HIGH-SPEED ECCS

This section presents the equations of the main ECC forms studied in this research as well as the points computations involved in ECC encryption process.

In addition, this section shows computational schemes for each curve and the parallel hardware designs needed to perform ECC computations.

The modified version of RLA used in scalar multiplication is illustrated as well.

#### A. ECC Points Computations

The two main operations in ECC scalar multiplication are point doubling and point additions. The key difference in the computations of the two operations is in finding the slope (m), which is the derivation of elliptic curve equation.

For point addition the calculation of m is similar for all ECC forms over GF(p), while in point doubling the slope differs from one form to another. Equations required to compute the slope for point doubling and point addition operations are presented in (1) and (2), respectively. For more information about calculating the slope for ECC point operations, the reader may refer to [27-29].

In point doubling, the slope (m) =

$$\frac{dy}{dx} = \frac{3[x^2+2a(x+1)]}{(2y)} \quad (1)$$

In point addition,

$$\text{the slope } (m) = (y_2 - y_1)/(x_2 - x_1) \quad (2)$$

It can be noticed that elliptic curve points are represented by x and y coordinates. In point doubling (3), the elliptic curve point is added to itself where in point addition (4), to different points are added. The result is a third point P<sub>3</sub>(x<sub>3</sub>, y<sub>3</sub>) on an elliptic curve, which can be computed as follows:

$$x_3 = m^2 - x_1 - x_2 \quad (3)$$

$$y_3 = m(x_1 - x_3) - y_1 \quad (4)$$

This research studies two major forms of elliptic curves, which are Weierstrass, and Montgomery curves represented in equations (5), and (6), respectively. These curves are very important because the Weierstrass curve is the most used form for ECC cryptographic operations, while the Montgomery curve is distinguished from its counterparts in terms of the degree of complexity for ECC computation, which may lead to the development of high-speed cryptosystem.

$$E: Y^2 = X^3 + aX + b \quad (5)$$

where a and b belong to GF (p) and 4a<sup>2</sup> + 27 and b  $\neq$  0.

$$by^2 = x^3 + ax^2 + x \quad (6)$$

where a, b belong to GF (p), and with b\*(a<sup>2</sup> - 4)  $\neq$  0.

According to previous studies, homogenous projective coordinates (X/Z, Y/Z) achieves the least computational complexity for points operations in the above mentioned elliptic curve forms represented by equations (5), and (6) [25-30].

Using projective coordinates, the elliptic curve points are represented by three coordinates, which are x, y, and z.

In Montgomery ECC, the result of point doubling can be calculated as follows:

$$\text{The slope } M = \frac{3X^2 + 2aXZ + Z^2}{2bYZ}$$

$$X_3 = (2bYZ)[(3X^2 + 2aXZ + Z^2)^2 - 8b^2Y^2ZX]$$

$$Y_3 = \{(3X^2 + 2aXZ + Z^2)[12b^2Y^2ZX - (3X^2 + 2aXZ + Z^2)^2] - 8b^3Y^4Z^2\}$$

$$Z_3 = (2bYZ)^3$$

The result of point doubling operation for Weierstrass is another point represented by the three coordinates  $x_3$ ,  $y_3$ , and  $z_3$ , as follows:

The slope  $M = \frac{3X^2 + aZ^2}{2YZ}$

$$X_3 = 2YZ * [(3X^2 + aZ^2)^2 - 8XZY^2]$$

$$Y_3 = (3X^2 + aZ^2) * [12XZY^2 - (3X^2 + aZ^2)^2] - 8Y^4Z^2$$

$$Z_3 = 8Y^3Z^3$$

It can be noticed that the slope (M) depends on the elliptic curve equation, while in point addition the calculation of the slope does not relate to the elliptic curve equation. Therefore, point addition computation is similar for all curves using the homogenous projection, and can be computed as follows:

$$M = \frac{Y_1Z_2 - Y_2Z_1}{X_1Z_2 - X_2Z_1}$$

$$X_3 = (X_1Z_2 - X_2Z_1) * [Z_1Z_2 * (Y_1Z_2 - Y_2Z_1)^2 - (X_1Z_2 + X_2Z_1) * (X_1Z_2 - X_2Z_1)^2]$$

$$Y_3 = \left\{ \left[ \begin{array}{l} (Y_1Z_2 - Y_2Z_1) * \\ X_1Z_2(X_1Z_2 - X_2Z_1)^2 - \\ (Z_1Z_2(Y_1Z_2 - Y_2Z_1)^2 - (X_1Z_2 + X_2Z_1) * (X_1Z_2 - X_2Z_1)^2) \\ Y_1Z_2 * (X_1Z_2 - X_2Z_1)^3 \end{array} \right] - \right\}$$

$$Z_3 = Z_1Z_2 * (X_1Z_2 - X_2Z_1)^3$$

In order to find the result of point doubling and point addition for each curve, a number of arithmetic operations, such as modular multiplication and addition over GF(p), should be performed. The required computations for point doubling and addition are called computational scheme for certain ECC and they play important role in determining the performance of the cryptosystem.

The next section presents the computational schemes for proposed Weierstrass and Montgomery ECCs as well as the hardware designs needed to implement them.

### B. ECC Computational Schemes and Hardware Designs

The computational scheme for ECC is the series of modular multiplication, addition, and subtraction operations required to perform point doubling and addition operations and hence calculating the values of  $x_3$ ,  $y_3$ , and  $z_3$  presented in the previous section.

Usually, ECC computations are implemented sequentially using hardware components, which are multipliers (M) and adders (A). In this research, all possible parallelization levels for performing ECC computations were investigated to find out the parallelization level that obtains the shortest time delay.

Table I presents the estimated performance levels for both Weierstrass and Montgomery ECCs when implemented using different parallelization levels as well as the required hardware devices for each implementation.

As can be noticed from Table I, the performance of ECC implementation is estimated by the number of sequential multiplication (SM) and sequential Addition (SA) levels required to perform point doubling operation. It is worth mentioning that the time consumed by SAs is neglected compared to SMs. So, the later will be the main factor used to estimate the time-consumption of ECC.

It is obvious that the 4-PM implementation for Montgomery ECC achieves the shortest time delay estimated by three SMs and three SAs. Similar performance was reported by Weierstrass ECC implementation but with using higher degree of parallelism, and hence more resources. In particular, the 5-PM Weierstrass ECC implementation can get comparable performance level to that obtained by Montgomery curve using 4 PMs.

Theoretical results showed that the use of more than 4 PMs for Montgomery ECC has no impact on the performance level. The same applies for the 5-PM ECC implementation for Weierstrass curve. In other words, the performance level is saturated at this degree of parallelism and cannot be improved further by adding more parallel Ms and As.

Experiments showed that the use of parallel hardware components is controlled by the inherited parallelism in ECC computations for both curves. For example, no additional parallel computations for Montgomery ECC point doubling can be performed in the first SM level, which involves five parallel multiplication operations. In order to start the second level's computations, the crypto processor needs to obtain the results of the first level, and so on.

Although other ECC implementations presented in Table I requires less resources because of the use of less number of hardware components, it seems that they consume greater time to perform point doubling. This research focuses on ECC implementations that score the highest possible performance level.

Similarly, Table II presents the theoretical results for studying the performance and resources consumption for ECC point addition. It is worth remembering here that the computations of point addition are similar for all ECC forms.

TABLE I. COMPUTATIONS REQUIRED TO PERFORM ECC POINT DOUBLING WITH DIFFERENT PARALLELIZATION LEVELS USING PROJECTIVE COORDINATES (X/Z, Y/Z)

| Elliptic Curves Form                                     | Hardware Design | Parallel Hardware Units | Sequential multiplication and addition levels |
|--|-----------------|-------------------------|---|
| ECC Point Addition for Montgomery and Weierstrass Curves | Serial Design   | 1M, 1A                  | 16 SM, 6 SA                                   |
|  | 2-PM            | 2M, 2A                  | 8 SM, 5 SA                                    |
|  | 3-PM            | 3M, 2A                  | 6 SM, 5 SA                                    |
|  | 4-PM            | 4M, 3A                  | 4 SM, 4 SA                                    |
|  | 5-PM            | 5M, 3A                  | 4 SM, 4 SA                                    |

TABLE II. COMPUTATIONS REQUIRED TO PERFORM ECC POINT ADDITION WITH DIFFERENT PARALLELIZATION LEVELS USING PROJECTION (X/Z, Y/Z)

| Elliptic Curves Form | Hardware Design | Parallel Hardware Units | Sequential multiplication and addition levels |
|----------------------|-----------------|-------------------------|---|
| Weierstrass ECC      | Serial Design   | 1M, 1A                  | 12 SM, 4 SA                                   |
|                      | 2-PM            | 2M, 2A                  | 6 SM, 3 SA                                    |
|                      | 3-PM            | 3M, 2A                  | 4 SM, 3 SA                                    |
|                      | 4-PM            | 4M, 2A                  | 4 SM, 3 SA                                    |
|                      | 5-PM            | 5M, 2A                  | 3 SM, 3 SA                                    |
| Montgomery ECC       | Serial Design   | 1M, 1A                  | 12 SM, 4 SA                                   |
|                      | 2-PM            | 2M, 2A                  | 6 SM, 3 SA                                    |
|                      | 3-PM            | 3M, 2A                  | 4 SM, 3 SA                                    |
|                      | 4-PM            | 4M, 2A                  | 3 SM, 3 SA                                    |
|                      | 5-PM            | 5M, 2A                  | 3 SM, 3 SA                                    |

It can be noticed from Table II that the 4-PM ECC implementation for point addition can achieve the shortest time delay estimated by four SM and four SA levels. Using higher degrees of parallelism has no positive impact on the performance as can be seen from the table.

Other presented ECC implementations represent a considerable trade-off between performance and time-consumption. This could benefit the purpose of developing efficient ECC for different security applications in accordance with required speed and the limitations on available resources. This research concerns the high-speed ECC implementations appropriate for multimedia security applications.

Fig. 2 and 3 present the parallel hardware designs of the computational schemes for high-speed ECCs using Montgomery and Weierstrass curves over GF(p), respectively.

Note that only the hardware designs that provide the shortest time delay are introduced, which the 4-PM and 5-PM designs for Montgomery and Weierstrass ECCs, respectively.

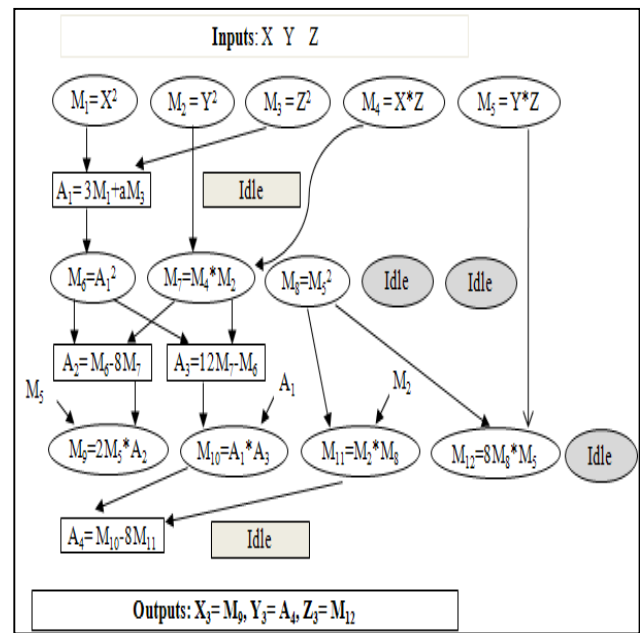


Fig. 3. Hardware Design for Weierstrass ECC Computations.

It should be highlighted here that the time consumed by one SM level is equivalent to the time consumption of one multiplication operation regardless how many multiplications are performed in that level.

For example, the first level SM level in Weierstrass ECC involves five multiplications. In despite of that, it consumes similar time to that needed for one multiplication operation. This represents a great benefit in terms of performance that can be achieved using parallel hardware implementation for ECC.

Fig. 4 shows the parallel hardware design for the computational scheme of ECC point addition operation. This design applies for both elliptic curves studied in this research.

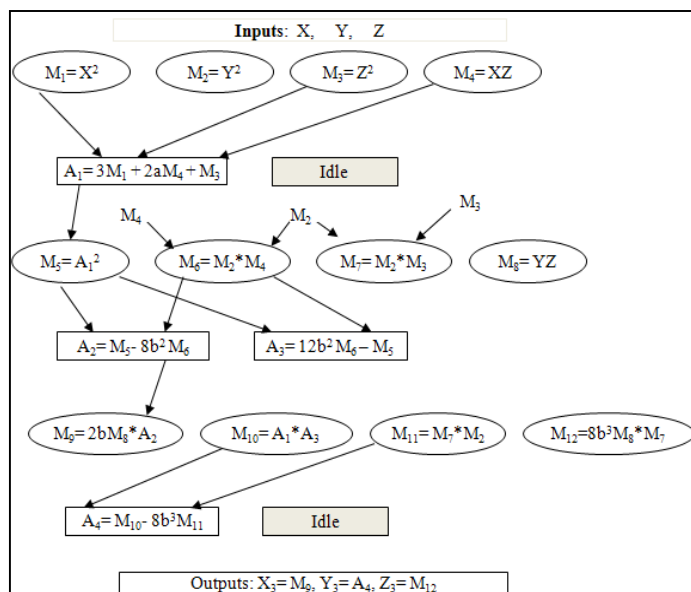


Fig. 2. Hardware Design for Montgomery ECC Computations.

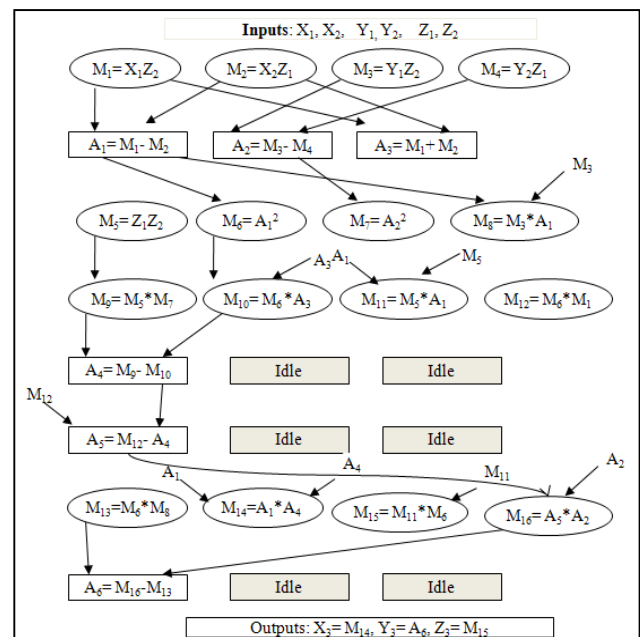


Fig. 4. Hardware Design for Point Addition Computations.

The results of point addition can be obtained after performing four SM and four SA levels.

As can be seen from Fig. 2, 3, and 4, the parallel implementation of the lower level of computations for ECC can improve the performance greatly. In addition, the next section presents the modified version of RLA, which allows the parallel implementation of upper level of computations for ECC. This contributes in accelerating the scalar multiplication even further.

### C. Modified Implementation of RLA

This section introduces an implementation of a modified version of known RLA for scalar multiplication, which is the main operation in ECC encryption process. The current proposed implementation allows to perform the upper level computations, represented by point doubling and point addition operations in parallel once required during scalar multiplication. The modified version of RLA is depicted in Fig. 5.

As can be noticed from Fig. 5, the inputs to the RLA are as follows:

- P which is a point on an elliptic curve E. The point P represents the original plaintext.
- K is the scalar represented by a series of binary bits (from 0 to n-1). N is the key length.

Outputs of RLA is another point (Q), which represents the encrypted message. Q is the result of scalar multiplication operation ([k]P).

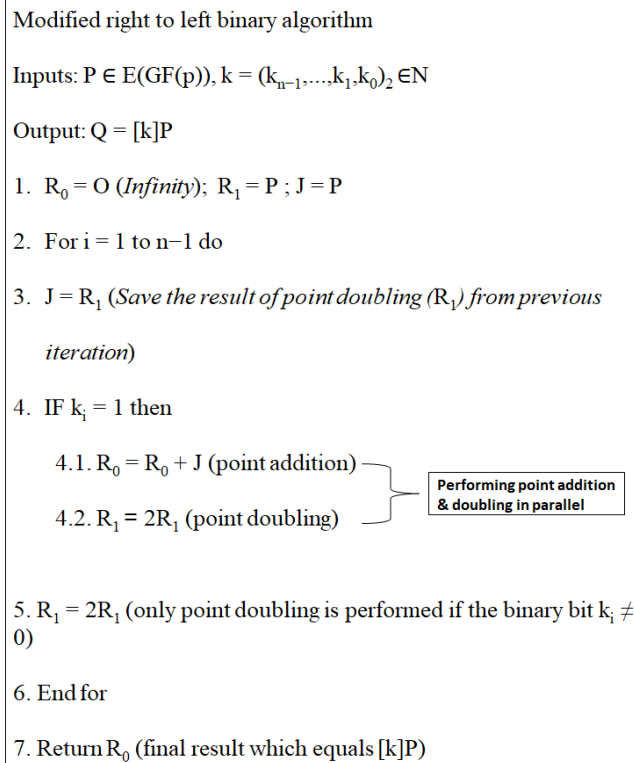


Fig. 5. Modified Version of RLA for Scalar Multiplication.

At the beginning of the algorithm, three elliptic curve points are defined, as follows:

- $R_0 = O$ , which is used to save the result of point addition during the scalar multiplication.
- $R_1 = P$ , which is used to save the result of point doubling during the scalar multiplication.
- $J = P$ , which is additional elliptic curve point used for temporary saving the value of  $R_1$  of previous iteration in scalar multiplication.

The RLA depends on the iterative approach and starts by scanning the binary bits of k from right to left.

For each bit in k, if the value of  $k_i$  equals one then both point addition and doubling operations are performed. Otherwise, only the point doubling operation is performed.

Note that the variable point J is used in each iteration to save the value of  $R_1$  obtained from the previous iteration of scalar multiplication. Then, J is used to perform the point addition operation in the current iteration. Thus, the point addition does not require the current value of  $R_1$  which is being calculated by point doubling operation. In other words, both operations become independent and can be processed in parallel manner.

The use of proposed modified version of RLA allows to perform the upper level computations of ECC in parallel to reduce the time delay of scalar multiplication and hence improve the performance of the cryptosystem.

In order to implement point doubling and point addition operations in parallel, extra hardware resources are needed. In particular, eight PMs are required for Montgomery curve. The 8-PM ECC design allows to perform the four parallel multiplications included in each point doubling and point addition as illustrated in Fig. 2 and 4. If the value of  $k_i$  is zero, only the point doubling is performed which makes four out of eight multipliers idle as this iteration. On the other hand, the Weierstrass ECC needs nine PMs to implement point operations in parallel. Five PMs for point doubling and four PMs for point addition as clarified in Fig. 3 and 4. The 8-PM and 9-PM designs accomplish the highest performance level for Montgomery and Weierstrass ECCs, respectively.

It should be mentioned here that the time complexity of proposed RLA is estimated by  $O(n)$ , where n is the size of inputs key or the number of binary bits of the key. This means it takes linear time complexity, which can satisfy required performance level for many applications.

Another advantage of using the proposed ECC is that it can increase the immunity of the cryptosystem against side channel attacks such the known STA. Each iteration of the modified RLA consumes approximately similar time regardless the value of the binary bit of the scalar k. This is because both operations are conducted in parallel if the value is one. So, the attacker cannot reveal the value of the binary bit  $k_i$  in each iteration by analyzing the time consumed. In this way the secret key, represented by the binary bits' vector of the scalar k, is kept confidential and cannot be exposed to the attacker by using the SPA.

On the other side, usual ECC implementations that use serial design implements point doubling and addition operations sequentially. This enables the attacker to trace the time consumed by certain iteration of the scalar multiplication operation and thus determining whether one or two operations have been performed by that iteration. Thus, it is possible to infer the values of the key bits by observing and analyzing the time consumption of RLA's iterations.

As mentioned previously, the current research work aims to develop secure and high-performance ECC that can encrypt/decrypt huge amount of data within the least possible time.

The next subsection describes the hardware implementation environment for proposed ECCs.

#### D. Implementation Environment

This research work seeks to develop high-speed ECC using hardware implementation. Thus, experiments conducted in this study investigates possible hardware design choices and analyzes the time delay and resources consumption for proposed hardware implementations. For this purpose, this study uses the popular hardware description language VHDL to implement proposed designs. The code for each proposed ECC is written in VHDL and then simulated using the ModelSim tool for validation purposes.

Authors used the standard carry save multiplier and carry save adder to perform field arithmetic operations for elliptic curve.

In order to simulate and synthesize suggested ECC implementations and obtain performance and resources-consumption results, researchers also use the Xilinx tool with the target FPGA (Field Programmable Gate Array) chip family chosen to be virtex5 (XC5VLX30).

The next section presents and analyzes the implementation results for proposed Montgomery and Weierstrass ECCs.

### V. RESULTS AND ANALYSIS

This section discusses the performance and resources consumption results for proposed ECCs in the current research.

Researchers observed the architectures of the proposed ECC implementations for Weierstrass and Montgomery forms. It can be noticed from previously presented hardware designs in Fig. 1 and 3 that such design requires eight parallel multipliers to implement the Montgomery ECC point doubling and point addition operations in parallel manner. In particular, four parallel multipliers are needed for each operation.

The Weierstrass ECC designs presented in Fig. 2 and 3 needs nine parallel multipliers; four for point addition and five for point doubling.

The required hardware resources, the estimated time consumption for point doubling and point addition as well as the overall time delay for the scalar multiplication are shown in Table III. It also shows the AT and AT<sup>2</sup> cost factors.

It should be mentioned here that time consumption for proposed ECCs is estimated in terms of the number of SM levels consumed by each point operation. The overall time consumption is calculated using the following equation:

$$\text{Overall time} = (\text{No. SMs for point doubling}) + (0.5 \times (\text{TIME of POINT ADDITION})) \quad (7)$$

The RLA assumes that point addition happens in half the number of the binary bits for the scalar (k). So, the time of point addition is multiplied by 0.5 in the above equation.

Remember that point doubling and addition are performed in parallel and the point addition consumes one SM level more than point doubling. Therefore, if the binary bit of k equals one, the time consumed is equivalent to the time of point doubling plus one SM level. The time delay of one SM is about one third the time allocated for point doubling operation as can be seen from Fig. 2 and 3. It can be estimated by (0.33 × TIME OF POINT DOUBLING). Therefore, equation (7) can be written as follows:

$$\text{Overall time} = (\text{No. SMs for point doubling}) + (0.5 \times (0.33 \times (\text{No. SMs for point doubling}))) \quad (8)$$

Table III shows performance and area consumption features for the proposed Montgomery and Weierstrass ECCs. This research uses the modified version of RLA to apply point doubling and point addition operations in parallel as depicted in Fig. 4. The 8-PM ECC for Montgomery curve scores the shortest time delay, estimated by 3.5 multiplication cycles. It also achieves the best AT<sup>2</sup> results compared to other ECC implementations presented in Table III. It is worth mentioning that the AT<sup>2</sup> factor focuses more on the time delay and it should be considered very important for developing high-speed ECC. Results showed that proposed 9-PM Weierstrass ECC achieves similar performance level but with consuming more resources.

It can be noted from Table III that proposed Montgomery and Weierstrass ECCs outperform their counterparts that use normal RLA in terms of performance. However, they consume more resources to enable the parallel implementation of the upper level of computations for the scalar multiplication operation.

Table III presents estimated results obtained by studying proposed ECC designs and analyzing their time delay and resources consumption features.

This research provides hardware implementation for proposed ECCs using the modified RLA. Implementation results show actual performance and resources consumption levels for Montgomery and Weierstrass ECCs when implemented using 8-PM and 9-PM designs. Obtained results are then compared with the corresponding ECC implementations that use usual RLA.

The performance is the most important factor to be considered when developing a high-speed ECC. It is assessed using the time required to perform the scalar multiplication operation (T<sub>KP</sub>), which can be computed through the following four steps of calculations.



TABLE III. COMPARISON BETWEEN PROPOSED PARALLEL ECC IMPLEMENTATIONS USING MODIFIED RLA AND ECCs THAT USE USUAL RLA

| Scalar Multiplication Algorithm         | Elliptic Curve Representation over GF (p) | ECC design | Consumed Resources | Time Delay     |                |              | Cost Factors |                 |
|---|---|------------|--------------------|----------------|----------------|--------------|--------------|-----------------|
|   |   |            |                    | Point Doubling | Point Addition | Overall Time | AT           | AT <sup>2</sup> |
| Proposed Parallel RLA                   | Montgomery ECC                            | 8-PM       | 8 M, 5 A           | 3              | 4              | 3.5          | 28           | 98              |
|   | Weierstrass ECC                           | 9-PM       | 9 M, 5 A           | 3              | 4              | 3.5          | 31.5         | 110.25          |
| Usual RLA implemented in [10-13, 17-21] | Montgomery ECC                            | 4-PM       | 4 M, 2 A           | 3              | 4              | 6            | 24           | 144             |
|   | Weierstrass ECC                           | 5-PM       | 5 M, 2 A           | 3              | 4              | 6            | 30           | 180             |

In the first step, we need to calculate the time consumption of one multiplication ( $T_M$ ) operation by using the following equation:

$$T_M = (\text{cycles/bit}) \times m \times \text{clockperiod} \quad (9)$$

where  $m$  is the key size. The current research experiments implement proposed ECCs using the key sizes of 256 bits, 512 bits, and 1024 bits.

For example, by using equation (9), the  $T_M$  can be computed as follows:

$$T_M = 1 \times 256 \times 9.079 = 2324.224 \text{ nano sec (n sec)}$$

Remember that in our proposed parallel ECC implementations the time consumption of one multiplication operation equals the time consumed by an entire SM level ( $T_{SM}$ ).

**Secondly**, compute the time consumption for each point doubling and point addition operation, which represent the main building blocks of the scalar multiplication.

This research concerns the ECC design that accomplish the heist performance level with consuming the least resources. Therefore, the 8-PM Montgomery ECC is implemented. It requires 3 SMs for point doubling and 4 SMs for point addition. The time consumed by one point addition ( $T_{ADD}$ ) and one point doubling ( $T_{DBLE}$ ) can be computed using the following equation:

$$T_{ADD} = 4 * T_{SM} = 9296.869 \text{ n sec}, T_{DBLE} = 3 * T_{SM} = 6972.672 \text{ n sec} \quad (10)$$

In the third step, we compute the time of one inversion operation ( $T_{INV}$ ), which is required to convert the projective coordinates back to the affine coordinates. The time of one inversion is equivalent to the time of three SMs [20-24]. Hence, the  $T_{INV}$  can be computed as follows:

$$T_{INV} = 3 * T_{SM} = 6972.672 \text{ n sec}$$

Finally, in the last step the total time consumption of the scalar multiplication ( $T_{KP}$ ) is computed via the following equation:

$$T_{KP} = ((256 \times 0.5 \times T_{ADD}) + (256 \times T_{DBLE}) + T_{INV}) \times 10^{-6} \quad (11)$$

Remember that using proposed parallel implementation for point doubling and point addition the  $T_{ADD}$  can be estimated one SM ( $T_{SM}$ ), which is equivalent to  $(0.33 \times T_{DBLE})$ , approximately. Therefore, equation 11 can be rewritten as follows:

$$T_{KP} = ((256 \times 0.5 \times (T_{SM})) + (256 \times T_{DBLE}) + T_{INV}) \times 10^{-6} \quad (12)$$

The result of the equation (12) is multiplied by  $10^{-6}$  to convert the final value from nano to mille seconds.

The performance or time consumption of the proposed 8-PM Montgomery ECC can be calculated using equation 12 as follows:

$$T_{KP} = ((256 \times 0.5 \times 2324.224) + (256 \times 6972.672) + 6972.672) \times 10^{-6} = 2.089 \text{ m sec.}$$

The performance of proposed Montgomery ECC equals 2.089 m sec. It overcomes the corresponding ECC implemented using normal RLA as well as the fastest known implementation using NAF algorithm introduced in [30]. ECCs are implemented using key sizes 256 bits 512 bits, 1024 bits, and 2048 bits as can be noticed from the table.

Table IV shows a comparison between proposed Montgomery ECC and other ECC implementations using NAF algorithm and usual RLA presented in [30] and [22] respectively. It is obvious that proposed 8-PM ECC accomplishes the best performance results compared to other ECC implementations using the different key sizes. For security reasons, it is recommended to use key sizes of 2048 bits and above for encryption and decryption processes. Current implementation shows the time consumption of ECC encryption using different key sizes. It was noticed from the experiments that the time consumptions of encryption and decryption operations are comparable when using the same key size. The Montgomery ECC using NAF algorithm obtains higher performance in comparison with the corresponding ECC implementation that uses normal RLA.

It should be mentioned that Montgomery elliptic curve has less computational complexity compared to the majority of other elliptic curves forms. Therefore, it is highly recommended to use it when developing high-speed cryptosystem.

The 8-PM ECC introduced in this research improves the cryptosystem immunity against STA. Using such cryptosystem's architecture makes it difficult to identify the binary bits of the decryption key via tracing and analyzing the time consumed by each RLA's iteration. The parallel implementation of upper level's operations prevents the attacker from distinguishing the points operation performed in every iteration, and thus the key bit cannot be revealed.

For example, assume that the time-consumption of point doubling is equivalent to TD while the time of point addition is estimated by TA. In usual RLA, if the current key bit equals one the time consumed by the cryptosystem for the processing of the finite field arithmetic computations can be estimated by TD+TA, while it can be estimated by TD if the bit value is zero. The proposed ECC algorithm consumes approximately similar time regardless the value of the key bit. This prevents the attacker from tracing the time-consumption of ECC to find out the sequence of bits in the key.

However, it seems that proposed ECC consumes more area and resources than other ECC implementations presented previously. In despite of that, the additional cost can be afforded for the sake of developing high-speed and secure ECC that fits certain security applications.

Although the proposed 9-PM Weierstrass ECC achieves comparable performance results to that achieved by the 8-PM Montgomery ECC, it needs more resources.

Fig. 6 shows a comparison between the performance levels of ECC implementations presented in Table IV. It can be clearly noted that proposed 8-PM Montgomery ECC using modified RLA outperforms its counterparts for all key lengths ranging from 256 bits to 2048 bits. Therefore, it represents considerable choice for the development of high-performance cryptosystem.

Table V presents time-consumption comparison between proposed Montgomery ECC and the main ECC implementations found in previous research works. As can be noticed, the proposed high-speed ECC achieves better performance results for the encryption/decryption process than results reported in previous literature.

This highlight the importance of the parallel implementation of point doubling and point addition once needed in the modified RLA iterations.

The next section of this article presents the key conclusions and future research works.

TABLE IV. COMPARISON BETWEEN PROPOSED 8-PM ECC USING MODIFIED RLA AND OTHER MONTGOMERY ECCS USING NORMAL RLA AND NAF ALGORITHM

| ECC Implementation                                     | 256 Key bits       | 512 Key bits       | 1024 Key bits      | 2048 Key bits       |
|--|--------------------|--------------------|--------------------|---------------------|
| <b>Proposed 8-PM Montgomery ECC using modified RLA</b> | <b>2.089 m sec</b> | <b>4.171 m sec</b> | <b>8.337 m sec</b> | <b>16.674 m sec</b> |
| The 4-PM Montgomery ECC using NAF algorithm [30]       | 2.506 m sec        | 5.012 m sec        | 10.024 m sec       | 20.048 m sec        |
| The 4-PM Montgomery ECC using normal RLA [22]          | 2.979 m sec        | 5.957 m sec        | 11.916 m sec       | 23.832 m sec        |

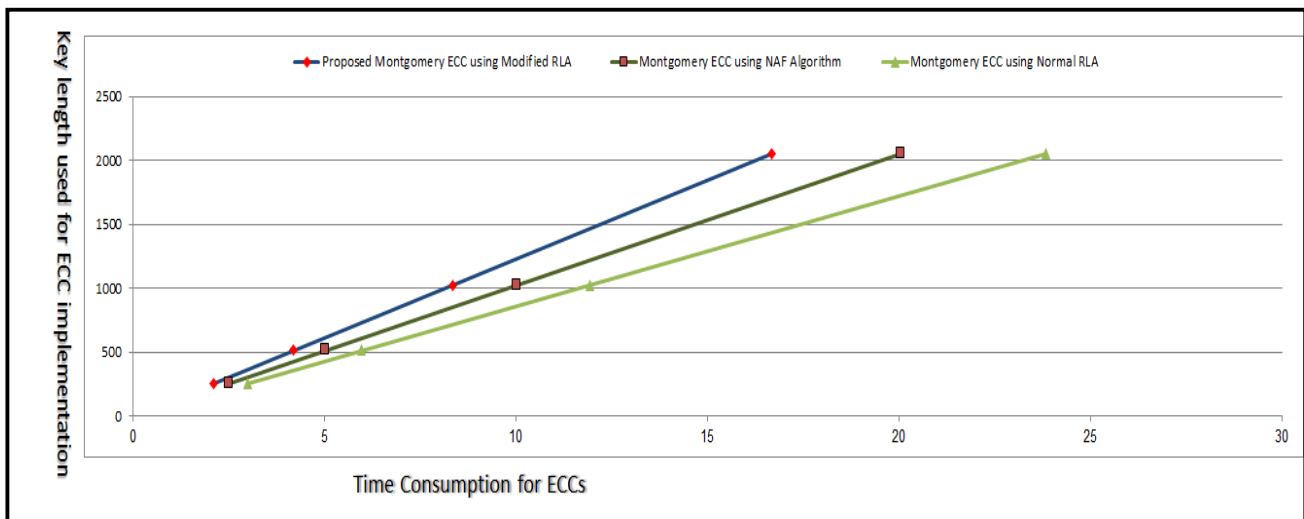


Fig. 6. Comparison between Performance Levels for Presented Montgomery ECC using different Key Lengths.

TABLE V. PERFORMANCE COMPARISON WITH PREVIOUS ECC IMPLEMENTATIONS

| No | Name, Ref.   | Finite Field, Key size        | KP time (m sec) | Device          |
|----|--|-------------------------------|-----------------|-----------------|
| 1  | Lo'ai et al. [12] (Edwards)                            | GF (p), 512-bit               | 33.301          | XC5VLX30        |
| 2  | Kerins et al. [25]                                     | GF (2 <sup>n</sup> ), 233-bit | 13.2            | XCV2000E        |
| 3  | Nele et al. [26]                                       | GF (2 <sup>n</sup> ), 160-bit | 3.8             | XCV800-4        |
| 4  | Kocabas et al. [28] (Binary Edwards)                   | GF (2 <sup>n</sup> ), 163-bit | 149.50          | ASIC            |
| 5  | Alkhatib [18] (Binary Edward curve)                    | GF (p), 256-bit               | 3.874           | XC5VLX30        |
| 6  | Alkhatib [21] (Edward curve)                           | GF (p), 256-bit               | 4.469           | XC5VLX30        |
| 7  | Montgomery ECC implementation using NAF [30]           | GF (p), 256-bit               | 2.506           | XC5VLX30        |
| 8  | Montgomery ECC implementation using Normal RLA [22]    | GF (p), 256-bit               | 2.979           | XC5VLX30        |
| 9  | <b>Proposed 8-PM Montgomery ECC using modified RLA</b> | <b>GF (p), 256-bit</b>        | <b>2.089</b>    | <b>XC5VLX30</b> |

## VI. CONCLUSION AND FUTURE WORK

This article introduced high-speed Montgomery ECC appropriate for security applications that requires fast encryption/decryption process.

Experimental results illustrated that proposed ECC surpasses previously developed ECCs in terms of performance. Much shorter time delay was reported by proposed ECC using different key lengths. The least time-consumption reported in this study is 2.089 m sec. This research developed high-performance cryptosystems for both Montgomery and Weierstrass curves. The later consumes more resources and hence the priority is given for Montgomery curve implementation.

Furthermore, proposed ECC improves the cryptosystem security against side channel attacks such as STA. It hinders the attacker's mission to analyze the time consumed by each iteration of the RLA. This was achieved by applying point doubling and point addition in parallel manner to hide the variation of time consumed by the RLA iterations.

A number of factors contributed in improving the performance and security level for proposed ECCs. Projective coordinates were used to avoid the time-consuming inversion operation. Finite field computations in each point operation were implemented using parallel hardware design. This increase the speed of performing the lower level of computations in ECC.

This research also proposed modified version of RLA to allow the parallel implementation of the upper level of computations represented by scalar multiplication operations. In order to achieve this, 8-PM and 9-PM designs were used for Montgomery and Weierstrass ECCs, respectively.

Results and comparisons illustrated that proposed high-speed cryptosystem accomplished better performance level compared to Montgomery ECCs using the usual RLA and NAF algorithm. It exceeds the performance of the major ECC implementations proposed in previous studies.

Such high-speed cryptosystem is considered the optimum choice for security applications that need fast encryption/decryption process to provide confidentiality for huge amount of data with consuming the least possible time. In multimedia applications, the cryptosystem needs to perform

encryption/decryption very fast to handle the enormous amount of data being transmitted inbound or outbound.

However, proposed ECCs consumes more area and resources in comparison with previous research works

In future, researchers may investigate the parallel implementation of ECC's upper computational level using NAF algorithm. Studying the performance and resources consumption levels for other forms of elliptic curves when implemented using proposed crypto processor architecture is another significant research direction.

## REFERENCES

- [1] Wade Trappe, Lawrence. c, Introduction to Cryptography with Coding Theory. Washington, Pearson Prentice Hall, 2002.
- [2] N. Koblitz, "Elliptic curve cryptosystem", Mathematics of Computation, Vol. 48, pp. 203-209, 1987.
- [3] V. Miller, "Uses of elliptic curves in cryptography", Lecture Notes in Computer Science, Vol. 218, pp. 417-426, 1986.
- [4] Blake, Seroussi, and Smart. Elliptic Curves in Cryptography. Cambridge University Press: New York, 1999.
- [5] Tanja Lange, "A note on L'opez-Dahab coordinates", Faculty of Mathematics, Technical University of Denmark, 2006.
- [6] David, Nigel, and Jacques, "Projective coordinates Leak", Applied research and security center, France.
- [7] GuericMeurice de Dormale, Jean-Jacques Quisquater. High-speed hardware implementations of Elliptic Curve Cryptography: A survey. Journal of Systems Architecture 53 (2007) 72-84, by Elsevier.
- [8] A. Satoh, K. Takano, "A scalable dual-field elliptic curve cryptographic processor," IEEE Transactions Computers 52 (4) (2003) 449-460.
- [9] G. Orlando, and C. Paar, "A scalable GF (p) elliptic curve processor architecture for programmable hardware," Cryptographic Hardware and Embedded Systems - CHES 2001, Paris, France, May 14-15, 2001.
- [10] Adnan Gutub and Mohammad K. Ibrahim, "High Radix Parallel Architecture For GF(p) Elliptic Curve Processor", IEEE Conference on Acoustics, Speech, and Signal Processing, ICASSP 2003, Pages: 625-628, Hong Kong, April 6-10.
- [11] Adnan Gutub. Efficient Utilization of Scalable Multipliers in Parallel to Compute GF (p) Elliptic Curve Cryptographic Operations. Kuwait Journal of Science & Engineering (KJSE) 2007, 34(2): 165-182.
- [12] L. Tawalbeh and Q. Abu Al-Haija. Speeding up Elliptic Curve Cryptography Computations by Adopting Edwards Curves over GF (P). International Journal of Security (IJS) 2009, CSC Journals, Malaysia, Vol.3, Issue.4, IJS-19.
- [13] Q. Abu Al-Haija and Mohammad Al-Khatib. Parallel Hardware Algorithms & Designs for Elliptic Curves Cryptography to Improve Point Operations Computations Using New Projective Coordinates. Journal of Information Assurance and Security (JIAS) 2010, By Dynamic Publishers Inc., USA, Vol.4, No.1, Paper 6: (588-594).

- [14] Adnan Abdul-Aziz Gutub and Mohammad K. Ibrahim, "High performance elliptic curve GF ( $2^k$ ) crypto-processor architecture for multimedia", IEEE International Conference on Multimedia & Expo, ICME 2003, pages 81- 84, Baltimore, Maryland, USA, July 6-9, 2003.
- [15] Adnan Gutub, "High Speed Low Power GF ( $2^k$ ) Elliptic Curve Cryptography Processor Architecture", IEEE 10th Annual Technical Exchange Meeting, KFUPM, Dhahran, Saudi Arabia, March 23-24, 2003.
- [16] L. Tawalbeh and Q. Abu Al-Haija. Enhanced FPGA Implementations for Doubling Oriented and Jacobi-Quartics Elliptic Curves Cryptography. Journal of Information Assurance and Security (JIAS), 2010, Volume 6, pp. 167-175.
- [17] Mohammad Al-khatib, Qacem, and AzmiJaafar. Hardware Architecture & Designs for Projective Elliptic Curves Point Addition Operation using Variable Levels of Parallelism. International Review on Computers and Software 2011 Vol. 6 N. 2, pp. 237-243.
- [18] Mohammad Al-Khatib, Q. Abu Al-Haija, and Ramlan Mahmud. Performance Evaluation of Projective Binary Edwards Elliptic Curve Computations with Parallel Architectures. Journal of Information Assurance and Security (JIAS) 2011, By Dynamic Publishers Inc., USA, Vol.6, No.1, Paper1: (001-009).
- [19] Mohammad Al-khatib, AzmiJaafar, and Q. Sbu Al-Haija. Choices on Designing GF (p) Elliptic Curve Coprocessor Benefiting from Mapping Homogeneous Curves in Parallel Multiplications. International Journal on computer science and engineering 2011, Vol.3, No.2, Paper 2: (467-480).
- [20] Mohammad Alkhatib, Azmi B. Jaafar, ZuriatiZukarnain, and Mohammad Rushdan. On the Design of Projective Binary Edwards Elliptic Curves over GF (p) Benefiting from Mapping Elliptic Curves Computations to Variable Degree of Parallel Design. International Journal on computer science and engineering 2011, Vol.3, No.4, Paper 44: (1697-1712).
- [21] Mohammad Alkhatib, Azmi B. Jaafar, ZuriatiZukarnain, and Mohammad Rushdan, "Trade-off between Area and Speed for Projective Edwards Elliptic Curves Crypto-system over GF (p) using Parallel Hardware Designs and Architectures", International Review on Computers and Software, July 2011 Vol. 6 N. 4, pp. 163-173.
- [22] Mohammad Al-khatib, AzmiJaafar, Zuriati Ahmad Zukarnain, and MohamadRushdanMd Said, "Hardware Designs and Architectures for Projective Montgomery ECC over GF (p) benefiting from mapping elliptic curve computations to different degrees of parallelism", International Review on Computers and Software, Vol. 6, N. 6, November 2011.
- [23] Mohammad Al-khatib, and Adel Al-Salem. Efficient hardware implementations for Tripling Oriented Elliptic Curve Crypto-system. International Review on Computers and Software. 2014, Vol. 9, N. 4.
- [24] A. Satoh, K. Takano, "A scalable dual-field elliptic curve cryptographic processor," IEEE Transactions Computers 52 (4) (2003) 449–460.
- [25] T. Kerins, E. M. Popovici and W. P. Marnane. "An FPGA Implementation of a Flexible, Secure Elliptic Curve Cryptography Processor", International Workshop on Applied Reconfigurable Computing-ARC 2005, IADIS press, pp.22-30.
- [26] Nele Mentens, Siddika Berna Ors, Bart Preneel, "An FPGA Implementation of an Elliptic Curve Processor over GF ( $2^m$ )", In *Proceedings of the 2004 ACM Great lakes Symposium on VLSI, GLSVLSI 2004: VLSI in the Nanometer Era.* pp. 454-457.
- [27] Turki F. Al-Somani. Performance Evaluation of Elliptic Curve Projective Coordinates with Parallel GF (p) Field Operations and Side-Channel Atomicity. Journal of Computers 2010. JCP.5.1.99-109.
- [28] Kocabas, U., J. Fan, and I. Verbauwhede, "Implementation of Binary Edwards curves for very-constrained devices," In 21st IEEE International Conference on Application-specific Systems Architectures and Processors, ASAP 2010, Rennes, France, pages 185 –191.
- [29] Mohammad Alkhatib, "Cost-Effective Implementations for Weirstrass and Montgomery Elliptic Curve Crypto-systems". International Journal of Computer Science and Information Security (IJCSIS) 2016, Vol. 14 N. 9, pp. 98-109.
- [30] Mohammad Alkhatib, "Improved ECC Performance Using NAF Algorithm for Binary Edward and Edward Elliptic Curves", IJCSNS International Journal of Computer Science and Network Security, 2019, Vol. 19 No. 6 pp. 98-109.
- [31] I Khan, H Rehman, Mohammad Al-Khatib, Z Anwar, M Alam, "A thin client friendly trusted execution framework for infrastructure-as-a-service clouds", Future Generation Computer Systems 89, 239-248, 2017.
- [32] Jin Y., Miyaji A. (2019) Secure and Compact Elliptic Curve Cryptosystems. In: Jang-Jaccard J., Guo F. (eds) Information Security and Privacy. ACISP 2019. Lecture Notes in Computer Science, vol 11547. Springer, Cham.
- [33] Susella, R., Montrasio, S.: A compact and exception-free ladder for all short Weierstrass elliptic curves. In: Lemke-Rust, K., Tunstall, M. (eds.) CARDIS 2016. LNCS, vol. 10146, pp. 156–173. Springer, Cham (2017).

#### AUTHOR'S PROFILE

<sup>1</sup>Imam University, faculty of Computer and Information Sciences, Department of Computer Science.



**Mohammad Al-khatibe** is an assistant Professor in faculty of Computer and Information Sciences at Imam University. He received the bachelor degree in Computer Science. His Master degree was obtained from Depaul University in the field of Information Systems. He also achieved PhD in Computer Science (Security in Computing) from University PUTRA Malaysia (UPM). His research interest includes: information security, cryptography, and Elliptic Curve algorithm.

# A Survey on Privacy Vulnerabilities in Permissionless Blockchains

Aisha Zahid Junejo<sup>1</sup>, Manzoor Ahmed Hashmani<sup>2</sup>

Computer and Information Sciences Department  
Universiti Teknologi PETRONAS  
Sri Iskandar, Malaysia

Abdullah Abdulrehman Alabdulatif<sup>3</sup>

Computer Department, College of Sciences and Arts  
Qassim University, P.O. Box 53, Al-Rass  
Saudi Arabia

**Abstract**—Blockchain decentralization not only ensures transparency of transactions to eliminate need of trusting third party, but also makes the transactions of the network to be publicly accessible to all the participating peers in the network. As a result, data anonymity and confidentiality are compromised making several business enterprises and industrialists hesitant to adopt the technology. Although research community has proposed various privacy-preserving solutions for blockchain, however, they still lack in efficiency resulting in distrust of industries in opting for the technology. This study is conducted for contributing to the existing body of knowledge corresponding to privacy in blockchains. The fundamental goal of this study is to delve into privacy vulnerabilities of the blockchain network in a permissionless setting by identifying non-trivial roots of factors causing privacy breach in blockchain and presenting limitation of existing privacy preserving mechanisms. Studies with superficial comparison of privacy preserving techniques are available in literature but a detailed and in-depth analysis of their limitations and causes of privacy breach in blockchain is yet not done. Therefore, in this paper we first present comprehensive analysis of various privacy breaching factors of the blockchain networks. Next, we discuss existing cryptographic and non-cryptographic solutions in literature. We found out that these existing privacy preserving mechanisms have their own set of limitations and hence are inefficient at current point of time. The existing privacy preserving mechanisms need further consideration of the research community before they're widely adopted and benchmarked. Therefore, in the end, we identified some future directions that need to be addressed to model an efficient privacy preserving mechanism for wider adoption of the blockchain technology.

**Keywords**—Blockchains; privacy vulnerabilities; cryptographic primitives; anonymity; confidentiality

## I. INTRODUCTION

The Blockchain technology is one of the most promising technological trends in the world today. It is a horizontal innovation that has the potential to impact every area of human endeavor [1]. The first application of Blockchains, widely known as Bitcoin, was introduced around a decade ago in October 2008 by S. Nakamoto [2]. Succeeding it, various other cryptocurrencies have been introduced [3] [4] [5] [6]. Initially introduced for the financial transactions of the cryptocurrency, the blockchain technology gradually spread to other sectors as well due to its inherent features. Over the years, the technology has been profusely researched and experimented to bring its benefits to other application areas. The technology has

eliminated the need of trusting third parties (i.e., banks) for authorization and record keeping of several transactions by providing transparency [7] and tamper resistance [8]. Transparency in Blockchain networks ensure the availability of the transactions to each node in a distributed network, whereas tamper-resistance makes each recorded transaction to be unmodifiable [9] or removable. Over the years, the technology has been profusely researched and experimented to bring its benefits to other application areas [10]. It is because of the decentralized, immutable and transparent nature of blockchain, that its applications have also been witnessed in non-financial areas like education, internet of things IoT, healthcare, big data, cloud computing, supply chain management, cyber security and so on. The blockchain ledger is written on a base and shared among the participating nodes for verification. This enables even the mutually distrusting nodes verify the data through consensus to achieve consistency and maintain the integrity of the blockchain network. Therefore, despite the fact that blockchain provides greater efficiency, reduced capital costs and greater data protection, it is still vulnerable to privacy issues. The data on the blockchain must be public because different nodes need to calculate and verify the same data so it must be accessible across the network. The transparency and credibility of the data is increased due to public availability of the data, however, it introduces the risk of privacy too as business enterprises and industrial organizations are not willing to make any business details public for adversaries to infer the personal information and extort the clients [11]. It is possible to set access control on the network using permissioned blockchains [12], however, the use of this type of blockchain makes the system more centralized and nullifies the purpose of using decentralized system, altogether. With the recent advancements in blockchain research and the eagerness of industries towards blockchain adoption makes privacy one of the key issues that need to be solved. The research in this paper has been carried out to highlight the issue of privacy in blockchain and the reasons behind it. This will help future researchers to solve the existing issues to get a better privacy protection in blockchain networks for a much wider adoption of this breakthrough technology.

### A. Gap Analysis and Contribution

According to the best of our knowledge, various studies [13] [14] [15] have highlighted the importance of privacy preservation in blockchain networks. Although these studies have contrasted existing mechanisms of ensuring privacy, however, they lack comprehensive insight towards possible

factors resulting in privacy disclosure. The study in the paper, therefore, presents comprehensive discussion on root causes of privacy breach in a blockchain network. Based on existing body of knowledge in the domain, we have managed to deduce some meaningful insights that will help research community to design more private blockchain networks. This research study is a multifold: i) describes blockchain technology and its benefits over traditional transaction systems, ii) elaborates the concept and need of privacy in relation to blockchain networks, iii) discusses privacy threats to blockchain and deduces the causes of privacy breach with respect to these threats, iv) discusses existing privacy solution and their limitations, v) suggests future directions to overcome privacy vulnerabilities in blockchain.

**B. Organization of the Paper**

The organization of the paper is as follows: Section II gives an overview of blockchain and its working mechanism followed by Section III that describes the issue of privacy in various settings of blockchain networks. Section IV discusses various factors causing privacy breach in blockchains. Further, Section V elucidates the existing privacy preserving mechanisms in blockchains and their limitations. Discussion and proposed future directions are presented in Section VI and Section VII concludes the study.

**II. INTRODUCTION TO BLOCKCHAIN TECHNOLOGY**

In 1991, S. Haber and W. S. Stornetta introduced the concept of a cryptographically secured network of blocks [16]. This concept was adopted by Nick Szabo as he worked upon and introduced decentralized digital currency called Bitgold. A decade later, in 2008, the concept was brought into practical implementation by S. Nakamoto [17] in the form of a cryptocurrency that is widely known is Bitcoin. It was since 2008, that the blockchain has been used to implement different cryptocurrencies. Additionally, due to the decentralized, immutable and transparent nature of blockchain, its applications have also been witnessed in non-financial areas like education [18] [19], internet of things IoT [20] [21], healthcare [8] [22] [23], big data, cloud computing, supply chain management [24] [25], cyber security and so on.

Since blockchain networks are distributed, hence the record of transactions is not stored on a single centralized server instead in a case a transaction occurs in the blockchain, it is distributed among all participating nodes where each node maintains a copy of the ledger [26]. This means that there exists thousands and millions of copies of the same blockchain where each node has access to the transaction details. Spreading the information across the network to multiple computers makes the information difficult to be manipulated hence providing transaction record integrity. Fig. 1 depicts the working mechanism of a blockchain network. A user A initiates the transaction that meant for a user B. This transaction is stored on a block and hence the block is created. Once the block has been created it is broadcasted to all participating nodes, also referred as peers, for verification of the transaction. If the transaction is validated by majority of the network, the newly created block is added to the existing chain and a copy of the updated ledger is maintained at each peer for record keeping. This completes a typical blockchain

transaction from user A to user B. The authenticity of transactions in a blockchain network is validated via asymmetric cryptography, also widely known as public key cryptography. It is one of the core components of blockchain technology [27]. More information on the types of cryptography can be found in [28] and is not discussed in detail as it is beyond the scope of this paper.

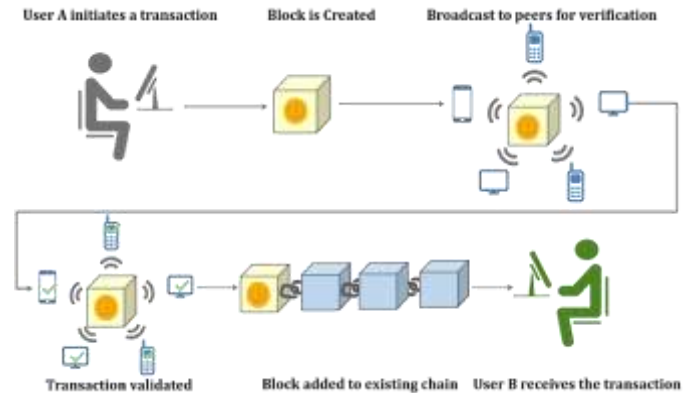


Fig. 1. Blockchain Working Mechanism.

**III. PRIVACY VULNERABILITY IN BLOCKCHAIN**

The blockchain networks are fundamentally transparent and distributed in nature, due to which they are widely being adopted and experimented. However, this means that all the data on a blockchain network is readily available for anyone on the network to view, causing privacy breach.

Blockchain networks can broadly be classified into two categories i.e. permissioned and permissionless blockchains. In a permissionless blockchain, a user requires no permission to enter the network. These kind of blockchains are open for anyone to join and participate. These systems have gained the attention of research community due to their decentralized consensus system [29]. On the other hand, special permissions are required in order to join a permissioned blockchain network. In a permissioned blockchain, the owner has the authority to decide who can join and become a part of the network. This means the blockchain owner has the ability and control to dictate the structure of the network, issue updates of the software, and control whatever operation and process occurs on that blockchain network.

|         |  | PERMISSIONLESS   | PERMISSIONED   |
|---------|--|--|--|
| PUBLIC  |  | <ul style="list-style-type: none"> <li>● Anyone</li> <li>● Anyone</li> <li>● Anyone</li> <li>● Anyone</li> </ul>                                     | <ul style="list-style-type: none"> <li>● Anyone</li> <li>● Anyone</li> <li>● Authorized User</li> <li>● Authorized User</li> </ul>                     |
|         |  | Requires Identity Privacy and Data Privacy   | Requires Identity Privacy and Data Privacy   |
| PRIVATE |  | <ul style="list-style-type: none"> <li>● Authorized User</li> <li>● Authorized User</li> <li>● Authorized User</li> <li>● Authorized User</li> </ul> | <ul style="list-style-type: none"> <li>● Authorized User</li> <li>● Authorized User</li> <li>● Network Operator</li> <li>● Network Operator</li> </ul> |
|         |  | Requires Data Privacy Only   | Requires Data Privacy Only   |

Fig. 2. Permissioned vs. Permissionless.



Private and public blockchains can have either permissioned or permissionless setting. This is illustrated in Fig. 2. Public and Permissionless allow anyone to join, read, write and commit to the transactions in the network. This means, all our data, be it personal or not, will be accessible by anyone in the network. This is where the issue of privacy arises. Moreover, in public and permissioned blockchains anyone can join and read the transactions, however only authorized users can write or commit. This improves trust in the blockchain but still doesn't guarantee the privacy of our assets. Similar is the case in Private and Permissionless blockchains. Lastly, in private and permissioned blockchains, although all users are known to the authorities, but this still doesn't guarantee the privacy of the data being transacted. So whatever type of blockchain it is, it does require privacy guarantee.

#### IV. CAUSES OF PRIVACY BREACH IN BLOCKCHAIN

Blockchains provide efficiency, reduced costs, transparency and trust but is still prone to privacy breach. For wider adoption, the privacy of blockchain networks must be strengthened. This section covers several causes resulting in privacy disclosure in blockchain networks.

##### A. Anonymization Inefficiency

In blockchain networks, anonymization refers to hiding the identity of the user. Anonymity is achieved when:

- Public address of the user cannot be mapped to his real identity.
- Blockchain transactions do not contain any personal identifiable information (PII).

Despite of blockchain claims of anonymity, it does not provide enough privacy. Several techniques are available in literature through which the anonymity of a blockchain network can be broken to identify the actual participants involved in a certain transaction. The phenomenon of disclosing user anonymity is known as deanonymization. In deanonymization, analysis of the network and network listening can help identify the blockchain user by unmasking him [13]. Further elaboration on deanonymizing blockchain users is presented in following subsections. Note that since cryptocurrencies are the first and widest applications of blockchain networks, hence the discussion carried out in following few sections will mainly focus cryptocurrencies to understand privacy mechanism and vulnerable areas of the technology. The same idea can further be applied to different applications.

1) *Deanonymizing via network analysis*: Each successful transaction in blockchain is added to transaction network where every node represents a transaction, and every (directed) edge represents a flow of data from an output of one transaction to an input of another. Analyzing the network relationships can be used to deanonymize a user's identity, thereby compromising the privacy. Since blockchain is a P2P network, hence IP address of nodes can be leaked [13] while transaction broadcasting.

2) *Deanonymizing via address clustering*: It is possible for transaction contents, transactions relationship with other transactions and the way transaction is broadcasted, to unintentionally leak information about the parties involved in the transaction to interested third parties. It is in fact noticed that various interested third parties systematically gather this kind of information to analyze various user patterns for multiple reasons including market research, competitor analysis, compliance and law enforcement. This analysis can (though not easily) be carried out using address clustering. The idea is to partition the set of addresses involved in a transaction to as many numbers of subsets as possible. Each subset, known as address cluster, most likely corresponds to the same entity. By combining address clusters with address tagging and graph analysis [30], the activity in blockchain can be effectively analyzed.

3) *Deanonymizing via transaction fingerprinting*: Another threat to anonymity is transaction fingerprinting. Androulaki investigated Bitcoin privacy provisions in a university setting. A simulator to mimic Bitcoin system was used and the results depicted that about 40% of the users' identities can be recovered despite of using Bitcoin's privacy measures [31].

Table I shows various deanonymization attacks on blockchain based cryptocurrencies.

TABLE I. DEANONYMIZATION ATTACKS ON CRYPTOCURRENCIES

| S.No | Paper Title   | Privacy Threat             | Success Rate | Test Case                    |
|------|---|----------------------------|--------------|------------------------------|
| 1    | An analysis of anonymity in Bitcoin using P2P network traffic [33]                            | Network Analysis           | >90%         | Bitcoin                      |
| 2    | Deanonymization of clients in Bitcoin P2P network [34]  | Network Analysis           | 11% - 60%    | Bitcoin,                     |
| 3    | Deanonymization and linkability of cryptocurrency transactions based on network analysis [35] | Network Analysis           |              | Bitcoin, Zcash, Dash, Monero |
| 4    | Data-Driven De-Anonymization in Bitcoin [36]  | Address Clustering         | 68.59%       | Bitcoin                      |
| 5    | Evaluating User Privacy in Bitcoin [31]   | Transaction Fingerprinting | 40%          | Bitcoin                      |

##### B. Transaction Pattern Linkability

Transaction information following through the public network can be used to reach out to statistical distributions on Cryptocurrencies revealing some new regulation within blockchain applications [13].

1) *Threat of transaction graph analysis*: M. Moser et al. [32] developed a framework based on transaction graph analysis to deanonymize the identities of users from publicly available transaction information in Bitcoin. Monero was taken as test case in the study and was empirically evaluated. Mix-ins used in Monero resulted in about 62% of the

transactions being unshielded to chain reaction i.e. deducing the actual input by elimination method. Moreover, The sampling of mix-ins in Monero is done in such a way that it gets easier to distinguish them from the real coins using their age distribution; in short, the real input is usually the “newest” input.

The authors estimated this phenomenon to guess the real input with around 80% accuracy. Further, each transaction in cryptocurrencies have some number of inputs and outputs that consume and create new coins respectively to conserve the total balance. Each input spends the new coins created in prior transaction and hence a transaction graph is formed. The public nature of blockchain data poses a potential privacy hazard to users. Since each transaction is publicly broadcast and widely replicated, any potentially identifying information can be determined for even years after a transaction is committed. The study depicted that a huge amount of data in Monero is traceable.

In another study [37] the authors focused on the typical behavior of users, the way they acquire spend their bitcoins, the balance of bitcoins they keep in their accounts, the way they move bitcoins between their various accounts in order to better protect their privacy. In addition, the research study isolated all the large transactions in the system, and discovered close relation of all these transactions to a single large transaction that took place in November 2010, even though the associated users apparently tried to hide this fact with many strange looking long chains and fork-merge structures in the transaction graph. Similarly, another study was carried out to test transaction linkability with the test case being Monero, again. In this study, three attack routines were developed to test against Monero’s privacy guarantee. The results of the study depicted in 88% of the cases it was easy to determine the origin of funds transferred.

2) *Web payment*: When a user makes a payment through web or online wallets, the consumer identity is prone to be linked to his real identity via browser cookies. When the user pays with a cryptocurrency, the service provider can link the real identity to the token history in the blockchain which also states that the attack is resilient against mixing mechanisms like CoinJoin [14].

In [38], two attacks are presented. The first attack shows that web trackers can extract substantial amount of information for advertising and analytics purposes when the user makes purchases on shopping websites. This information is enough to identify the blockchain transaction uniquely for linking it with the web cookies of the user to further reveal user’s identity. The second attack depicts that by linking even two purchases of the same user, the web tracker can identify his cluster of addresses even if anonymity techniques of blockchain such as CoinJoin are deployed. Moreover, it is possible to apply the attacks to past purchases as well. Thus, in the study, it is shown that third party web trackers have the ability of deanonymizing the cryptocurrency users.

A summary of studies carried out under this kind of privacy threats is given in Table II.

### C. Crisis of Private Key Theft

Private keys in a blockchain network are very critical to ensure the security and privacy of the user because these keys are used for signing each transaction in the network. Participant’s assets are controlled through private key in the blockchain systems. Hence, it is very important that proper key management systems [39] are enforced. If compromised, it can not only lead to privacy leakage but may also result in identity theft.

Although, private key allows a user to have sovereignty over his assets, however it comes under the responsibility of securing and managing one’s own private keys. Currently, there are no efficient mechanisms for recovery of the keys in a case of loss. Table III summarizes some of private key theft incidences compromising the security and privacy in blockchain systems.

TABLE II. TRANSACTION PATTERN LINKABILITY ATTACKS ON CRYPTOCURRENCIES

| S.No | Paper Title   | Privacy Threat             | Success Rate   | Test Case      |
|------|---|----------------------------|--|----------------|
| 1    | A Traceability Analysis of Monero’s Blockchain [40]   | Transaction Graph Analysis | 88%  | Monero, RingCT |
| 2    | Quantitative Analysis of the Full Bitcoin Transaction Graph [37]                              | Transaction Graph Analysis | 62%  | Monero         |
| 3    | When the cookie meets the blockchain: Privacy risks of web payments via cryptocurrencies [38] | Web Payment                | Attack I: 90%<br>Attack II: T = 2 – 89%<br>T = 3 – 99% | Bitcoin        |

TABLE III. PRIVATE KEY THEFT INCIDENTS

| S.No | Incident   | Amount                     | Year | Victim      |
|------|--|----------------------------|------|-------------|
| 1    | A hacker took possession of the administrative account was hacked and private keys were stolen. BTC price was changed to 1 cent and bought BTC from Mt. Gox users. | 2643 BTC                   | 2011 | Mt. Gox     |
| 2    | Attacker got access to bitcoinica database, obtained private information of users for theft.   | 38,000 BTC                 | 2012 | Bitcoinic a |
| 3    | Unencrypted Private Keys stored online for backup were stolen  | 24,000 BTC                 | 2012 | Bitfloor    |
| 4    | The attacker under the nickname Lucky7Coin inserted the Trojan code into the code of Cryptsy—a cryptocurrency exchange. A hacker got access to BTC and LTC keys.   | 13,000 BTC<br>300,000 LTC. | 2014 | Cryptsy     |
| 5    | Hackers infected the internal network of the exchange with a virus that was transmitted through email, and it allowed them to steal private keys.                  | 523M NEM                   | 2018 | Coinchec k  |
| 6    | Phishing and malware tactics were used to steal user 2FA codes and API keys alongwith customers’ private details.  | 7,000 BTC                  | 2019 | Binance     |



## V. INEFFICIENCY OF EXISTING PRIVACY-PRESERVING FRAMEWORKS

Blockchain technology has two categories when it comes to preserving privacy. The first category involves protecting the identity of the user by assigning him complete anonymity while making transactions. The second category involves protecting the transaction data from unauthorized entities and hackers thus maintaining data confidentiality. The classification of various privacy preserving techniques surveyed in the literatures are depicted in Fig. 3 and detailed in the subsequent section. The classification is done based on which technique contributes towards achieving what kind of privacy in blockchains.

The privacy preserving frameworks, reviewed in literature can broadly be classified into two categories, i.e.:

- **Mixing Methods:** Mixing methods or services are used to retain the transaction data privacy of the blockchain networks.
- **Cryptographic Primitives:** Cryptographic primitives are mathematical functions that are used in cryptography to verify data authenticity.

### A. Privacy Vulnerability in Mixing Services

Link between sender and receiver in a blockchain network can be known by analyzing the publicly available content. Introduction to mixers provides a solution to the stated problem. The concept of mixing service was first presented in [41] by Chaum. It allows users to hide who a participant communicates with as well as the content of the communication.

In Fig. 4, the basic architecture of a mixer is depicted. There are two types of mixing services, i.e., centralized mixing and decentralized mixing. Both concepts are elaborated:

1) **Centralized mixing:** Multiple mixing websites are available for use. These offer mixing of the transactions anonymously on exchange of mixing fees. The websites swap the transactions among various users so that the relationship between incoming and outgoing transactions can be hidden. Centralized mixing suffers from various limitations (discussed in section 3.3) including the mixing server being prone to denial of service (DOS) attacks as the server remains a single point of failure. Resultantly, it becomes an obstruction of the distributed blockchain network

2) **Decentralized mixing:** Decentralized mixing overcomes the limitations of centralized mixing which makes it vulnerable to DOS attack. A decentralized mixing pattern is proposed to enable a set of mutually untrusted peers to publish their messages simultaneously and anonymously without the need of a third-party anonymity proxy. Moreover, decentralized mixing eliminates the need of paying mixing fees. CoinJoin [42] and MultiParty [43] Computation are only two methods in literature that has successfully implemented decentralized mixing services.

3) **Critical analysis of mixing services:** Although mixing services can provide a substantial amount of identity privacy,

however, it has its own set of concerns which shall be taken into account before opting out for such a privacy preserving mechanism. These issues are discussed below:

a) **Waiting delay:** In order to use mixing services, user must wait for other participants to swap their transactions in order to hide and relationships between a transaction inputs and outputs. This incurs high waiting delay for a transaction to be completed.

b) **Third party involvement:** Since mixing servers are usually websites or other third-party software, hence they're not an appropriate solution to the privacy vulnerability of blockchain networks.

c) **Malicious mixing services:** Although mixing services hide the relationship between a user's transaction's input and output from an adversary, however, the server itself knows about all the input-output pairs and hence, the privacy in this scenario solely relies on how honest the intermediary is and becomes prone to breaches.

d) **Mixing fees:** Mixing services usually incur cost of hiding the identities of the users via mixing.

### B. Privacy Vulnerability in Cryptographic Primitives

There are two categories of cryptographic algorithms when it comes to blockchain networks. The first ones are primary, which are important for data transaction and communication in blockchain networks, the second ones are optional which are used for preserving and enhancing user and transaction data privacy [44] in blockchain networks.

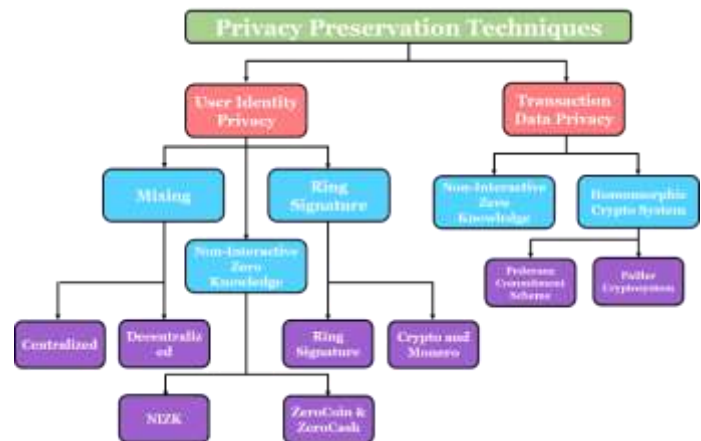


Fig. 3. Classification of Privacy Preserving Techniques.

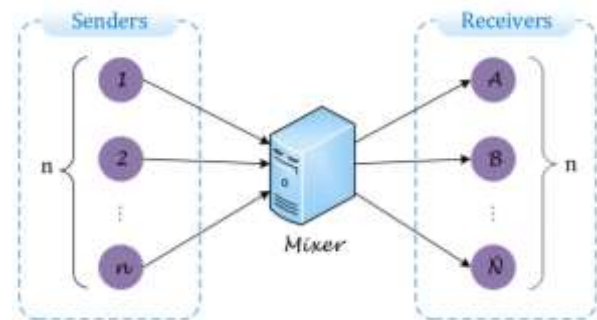


Fig. 4. Mixing Service Architecture [57].

In permissionless blockchain networks, any peer is able to join the network as participant at any point in time. No centralized authority manages or supervises that who joins the network or who should be banned from the network in permissionless scenario. This results in the content of the blockchain to be readable by any peer in the network. However, using optional cryptographic primitives, a permissionless blockchain network can be designed in such a way that privacy of the network is enhanced and each peer gets only relevant information [44]. Currently, the most widely used technologies to achieve blockchain privacy are ring signatures and zero-knowledge proofs.

1) *Ring signature*: In cryptography various kinds signatures, such as blind signature, ring signature, group signature and DC-nets, from which only ring signature and its variants are used to achieve anonymity in blockchains [44].

Ring signature was introduced in 2001 by Rivest et al. [45]. The concept behind ring signature is that a user chooses a set of participants to create a ring, including himself. Each participant in the ring has a public key. The user initiating the ring signs the message with his/her private key and public keys of all participants. Verifying node knows that one of the members signed the message but can't tell who actually signed it. Hence, anonymity is achieved.

The working mechanism of ring signature is illustrated in Fig. 5. The signature is analogous to the signature for a cheque in joint bank account where all participants sign the transaction with their public keys along with the originator's private key. After each participant of the ring has signed the transaction, it goes further for validation and verification.

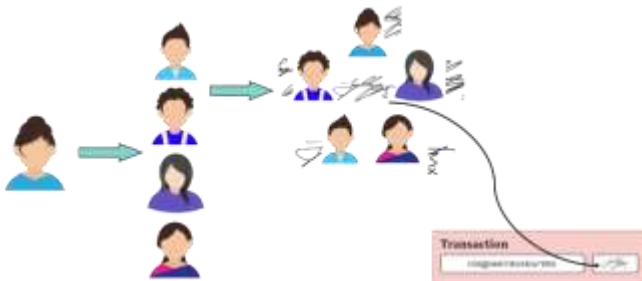


Fig. 5. Transaction Signing in Ring Signature.

Two basic advantages provided by ring signatures include unforgeability and anonymity [46]. Anonymity can further be sub divided into two properties i.e. unlinkability and untraceability [44]. Unlinkability refers to the verifier not being able to decide the link between two transactions whereas untraceability refers to the signer not being identified. These properties have led to development of several ring based privacy preserving protocols [47] [48] [49] [50] [51] which are widely used blockchain networks.

A signature scheme known as linkable spontaneous anonymous group (LSAG) was proposed in 2004 [47]. It is a variant of linkable ring signature in which groups are formed spontaneously without any group manager. The concept of ring signature was extended in [48] into traceable ring signature where an issue related tag was added to the signature. This idea was further adopted in [49] for the design on Ring-Coin with

improved efficiency. In this case, anyone in the ring, pretending to be another person to sign the same message, would face the risk of revealing his/her identity immediately. This idea was further adopted for preventing double-spending attack in blockchain and became the basis of CryptoNote [50] with a slight modification.

Furthermore, a concept of confidential transaction, using homomorphic commitment protocol, was proposed [51] for hiding transaction amounts. Later, three techniques i.e. ring signature, confidential transaction and multilayered linkable spontaneous anonymous group signature (MLSAG) [52] were combined to form Ring Confidential Transactions (RingCT), with its implementation being in Monero. Besides these, one-time signature [53], borrowmean signature and multisignatures are also used for preserving privacy in blockchain networks [44].

2) *Critical analysis of ring signature*: Monero [50], based on ring signature is considered to be the most efficiently privacy preserving cryptocurrency, however, Monero (due to vulnerabilities in the architecture of ring signature) also faces privacy issues. Some issues with ring signature include:

a) *Large ring size*: The size of the ring is directly proportional to the number of participants involved in the ring; this increases the ring size. To keep the ring size limited, usually the no. of participants that can take a part in ring formation is limited. This reduces the anonymity set size, hence increasing the risk of deanonymization.

b) *Lack of scalability*: Transaction size in ring signature is large – almost thousands of bytes per transaction. This will require more storage space to keep the records of the entire blockchain, hence compromising the scalability of blockchains.

c) *Transaction timing attack*: When a user creates the ring for his transaction, he usually collects other transactions of the same denomination available in the blockchain. Since each transaction in blockchain is time stamped, hence the newest created transaction in the anonymity set is considered to be the one to be redeemed. A study [40] depicts that 98% of the transactions are prone to time attack for traceability.

3) *Zero-knowledge proof*: Zero-knowledge protocols, introduced in 1980s [54], are one of the most widely used cryptographic techniques to enable the transfer of assets across a distributed, peer-to-peer blockchain network with improved privacy. The goal of zero-knowledge proofs is to prove the validity of a transaction with zero knowledge provided to the verifier about the transaction. The concept involves the certifier to formulate a formal proof to prove that a certain assertion is true without the need of providing any additional and useful information to the verifier [15]. A variant of ZKP, known as Non-Interactive Zero-knowledge Proof (NIZK proof), is widely used in blockchains as it eliminates the need of to and fro communication between the prover and the verifier and instead, requires only one time message to be sent from prover to the verifier. It is important to remark that not all ZKP schemes are non-interactive. Most of the ZKP protocols available in literature are interactive. Usually, in

ZKP scenario, the prover is required to answer various challenges sent by verifier, resulting in multiple rounds of communication. However, for blockchains and other distributed ledger technologies (DLT), it is desirable to avoid the communication because either (i) validating nodes can't properly agree on how to choose those challenges, since in many constructions we have to choose them randomly, while the verification algorithm must be deterministic in order to reach consensus; or (ii) because it would make the communication complexity of the system very poor. This property makes it suitable for anonymous and distributed verification of messages in blockchains.

The concept first appeared in [54] and is accepted for creating privacy preserving protocols in blockchain networks. NIZK proofs must meet the following three properties:

- Completeness: Everything that is true has a proof.
- Soundness: Everything that can be proved is true.
- Zero knowledge: Only the proven statement is revealed.

Zerocoin, introduced in [55] uses NIZK proof cryptography for providing anonymity by preventing transaction graph analysis i.e. by breaking the trace of coins. However, it fails to provide complete anonymity due to following reasons:

- Fixed denomination coins are used.
- Before payment is made, anonymous coins need to be converted into non-anonymous ones.
- The amount of transactions, or other metadata is not hidden.

To overcome the limitations of Zerocoin, zerocash was introduced [56]. Identity and transaction privacy were simultaneously provided in Zerocash to overcome the limitations of Zerocoin. It uses anonymous coins to provide privacy in blockchains. Further, size of transaction and time of verification of transactions were also significantly reduced. Zerocash uses ZK-SNARKS. However, the NIZK protocol incurs high computation overheads, especially in the proof generation phase of zk-SNARKs protocol used in Zcash.

4) *Critical analysis of zero knowledge proof:* Despite of providing both identity privacy and data privacy, ZKPs still have not perfected at preserving privacy in blockchain networks. A few issues with ZKPs include:

a) *Trusted Setup Problem:* The working of ZKPs involve a parameter generator that can issue prover and verifier keys to verify a transaction. This is where vulnerability to privacy breach arises as it is very significant to consider who to trust for parameter generation and how to ensure no record keeping at the generator. If compromised, this may result in forgery of the data.

b) *High Computation Overhead:* Theoretically, ZKPs achieve the highest level of anonymity and transaction privacy protection for the blockchain but at the expense of high computational costs it requires when it generates the transaction proofs.

c) *Prone to deanonymization:* A study [57] empirically shows that 98% transactions in Zcash are linkable.

## VI. DISCUSSION FOR WAY FORWARD

Maintaining privacy in blockchain based networks is very significant for its wide acceptance and adoption as shown in the literature. Besides the actual data, metadata also flows through the blockchain network. This metadata can be used to infer additional information about the users participating in the transaction. Additional information inferred may include the identity of the user and this identity unmasking can further reveal all the transactions related to the user. In other words, even with the most powerful privacy preserving mechanisms, this metadata still flows through the network. This is one of the biggest challenges for any privacy protecting approach that might be used in public permissionless blockchain networks. Adding mix-ins to transactions do not have an impactful effect either. Temporal analysis makes it evident that timing plays a major role in analysis of user identity thereby nullifying the effect of mix-ins. Analyzing transaction relationships, patterns, time and links, it becomes easier to trackback the headnode and determine the identity of a person. Once the identity of an individual is leaked, all the corresponding transaction information of the individual also gets prone to leakage.

In certain organizations, it is not desirable to make the confidential data publicly available, for instance patient records in healthcare, sensor data in IoT devices, private goods' information on supply chain management systems, business transactions in financial sector and so on. Hence, keeping privacy intact when blockchains are deployed for those applications, has a great significance. If privacy is not ensured, the integration of blockchain in such application areas may not progress and soon come to a halt. Setting access control is possible by permissioned blockchain, however, using those kinds of blockchains nullifies the purpose of using a decentralized system altogether.

Privacy in a blockchain network can be preserved in various ways but the most prominent one in literature is preserving privacy through the use of efficient cryptographic primitives. A brief summary of type of privacy offered and limitations of existing privacy protecting mechanisms' implementations is presented in Table IV.

It can clearly be seen from the table that existing approaches have a number of limitations and thus need further research for reduction of the privacy risk in blockchain systems. Hence, a few research directions are presented that can be investigated further.

### A. Transparency vs. Privacy

Blockchain is transparent by virtue of its design. Transparency, however, can be a double-edged sword when it comes to blockchain transactions. On one hand, blockchain is trusted for its transparency whereas on the other hand, this results in serious privacy concerns for a variety of potential application domains. The desire of stronger privacy in some applications leads to limited usage of the technology. Hence, the biggest challenge to achieve privacy in blockchain systems is finding the correct balance between the degree of transparency and the degree of privacy leveraged.

TABLE IV. SUMMARY OF EXISTING PRIVACY PRESERVATION TECHNIQUES

| S.no | Privacy Preservation Method                        | Type of Privacy |      | Fundamental Framework | Limitations  |
|------|--|-----------------|------|-----------------------|--|
|      |  | Identity        | Data |                       |  |
| 1    | Mixing websites, (Cryptomix, Bitmix.Biz, SmartMix) | ✓               | ✗    | Centralized Mixing    | Long waiting delay<br>High Mixing Fees (4-5% of the transaction for these particular types)<br>Prone to DDoS and Sybil attacks   |
| 2    | Centralized Tumblers [58]                          | ✓               | ✗    | Centralized Mixing    | Long waiting delay<br>High Mixing Fees<br>Prone to money laundry attacks<br>Depends on the trusted party<br>Cannot guarantee safety from theft                                   |
| 3    | CoinSwap [59]                                      | ✓               | ✗    | Centralized Mixing    | Long waiting delay<br>High Mixing Fees<br>Prone to DDoS and Sybil attacks<br>No proof that the mixer is not storing transaction record   |
| 4    | CoinJoin [42]                                      | ✓               | ✗    | Decentralized Mixing  | Long waiting delay<br>Prone to DDoS and Sybil attacks<br>Lacks internal unlinkability  |
| 5    | CoinShuffle [60]                                   | ✓               | ✗    | Decentralized Mixing  | High communication and computation overhead<br>Can be frustrated by dishonest participants<br>Prone to Sybil attack  |
| 6    | CoinParty [43]                                     | ✓               | ✗    | Decentralized Mixing  | 2/3 users are honest (in theory)<br>Lesser theft prevention  |
| 7    | RingCT [53]  | ✓               | ✗    | Ring Signature        | Large Transaction Size<br>Increasing no. of participants increases ring size   |
| 8    | CryptoNote [50]                                    | ✓               | ✗    | Ring Signature        | Limited Ring Size<br>Lacks Scalability due to larger transaction size<br>Smaller anonymity set   |
| 9    | Zerocoin [55]                                      | ✓               | ✓    | Zero-knowledge Proofs | Requires larger proof size (Computationally complex)<br>Leakage of trusted setup parameters can lead to forgery of coins<br>Requires fix denominations<br>Requires trusted setup |
| 10   | Zerocash [56]                                      | ✓               | ✓    | Zero-knowledge Proofs | Computationally intensive<br>Leakage of trusted setup parameters can lead to forgery of coins<br>Requires trusted setup  |

### B. Scalability

Some privacy preserving techniques provide a sufficient amount of privacy for a wide variety of applications. In addition, advanced versions of some of the existing techniques i.e. ring signatures and zero-knowledge proofs (ZKP) provide both user privacy and data content privacy. However, this privacy protection is done at the cost of scalability of the network. Scalability, itself, is one of the major concerns in the technology of blockchain these days, hence, industrialists do not opt for the privacy solutions that further increase the issue. The need of scalable solutions make it another significant challenge in terms of privacy protection of user and user assets. Therefore, researchers should delve further into the cryptography of these techniques to find out the loopholes in existing techniques. The identified loopholes will further help the researchers to model scalable privacy preserving mechanisms.

### C. Private Key Management Systems

Loss or theft of private is another major issue that may result in privacy breach of the user and loss of user assets associated with the key. Proper private key management

systems should, therefore, be incorporated. Moreover, mechanisms to recover or report the lost keys should be brought into practical implementation.

## VII. CONCLUSION

Invention of blockchain eliminated the need of trusting a third party for record keeping and transaction verification. Blockchains promote transparency by introducing publicly verifiable transactions. However, this transparency has led the blockchain community to an emerging issue of privacy. Privacy in blockchain refers to safeguarding the identity of the user involved in a transaction and protecting the secrecy of transaction data. Although researchers and industrialists have proposed some privacy preserving mechanisms over the years, however, these mechanisms are still prone to privacy breaches and do not provide complete privacy. For instance, mixing services and ring signatures can provide user identity privacy only and does not provide transaction data privacy. Similarly, homomorphic cryptosystems aim at providing transaction data privacy but does not provide user identity privacy. Moreover, although ZKPs provide both kinds of privacy in blockchains but it does so at the cost of system performance. Poor

performance of the techniques restricts universal adoption of blockchain technology. Hence, the need for a more efficient privacy preservation framework that doesn't only retain user identity and transaction data privacy, but also ensures the performance of the system doesn't lag arises. For development of an effective solution to problem of privacy in blockchain, understanding the root cause of the issue is important. Therefore, in this study we have highlighted some privacy breaching causes by the virtue of blockchain design. These causes include (i) additional information flowing through the network that aids in deanonymizing a blockchain user; (ii) linking the time and pattern of transactions; and (iii) absence of effective private key management systems in the case of private key thefts. In order to be completely benefitted by the variety of features that blockchain has to offer, it is essential that the privacy in blockchain systems shall be strengthened.

#### ACKNOWLEDGMENTS

The authors would like to extend their gratitude to Universiti Teknologi PETRONAS for provision of necessary equipment and resources to carry out the research.

#### REFERENCES

- [1] Junejo A.Z., Memon M.M., Junejo M.A., Talpur S., Memon R.M. (2020) Blockchains Technology Analysis: Applications, Current Trends and Future Directions—An Overview. In: Peng SL., Son L., Suseendran G., Balaganesh D. (eds) Intelligent Computing and Innovation on Data Science. Lecture Notes in Networks and Systems, vol 118. Springer, Singapore. [http://doi-org-443.webvpn.fjmu.edu.cn/10.1007/978-981-15-3284-9\\_47](http://doi-org-443.webvpn.fjmu.edu.cn/10.1007/978-981-15-3284-9_47).
- [2] Satoshi Nakamoto, "Bitcoin: A peer-to-Peer Electronic Cash," 2009.
- [3] C. Lee, "Litecoin," 2011. Whitepaper.
- [4] "Peercoin—secure & sustainable cryptocurrency," August 2012. Available: <https://peercoin.net/whitepaper>.
- [5] S. King, "Primecoin," 7 July 2013.
- [6] I. Grigg, "EOS: An Introduction,"
- [7] Guy Zysking, Oz Nathan, Alex 'Sandy' Pentland, "Decentralizing Privacy: Using Blockchain to Protect Personal Data.," in Security and Privacy Workshops, San Jose, 2015.
- [8] Asad Ali Siyal, Aisha Zahid Junejo, Muhammad Zawish, Kainat Ahmed, Aiman Khalil, Geroglia Sorsou, "Applications of Blockchain Technology in Medicine and Healthcare: Challenges and Future Perspectives.," MDPI Cryptography, vol. 3, no. 3, 2019.
- [9] D. Yaga, Peter Mell, N. Roby, K Scarfone, "Blockchain Technology Overview," NIST, 2018.
- [10] V. K. Supriya Thakur, "Blockchain and Its Applications – A Detailed Survey," International Journal of Computer Applications, vol. 180, no. 3, 2017.
- [11] Yuchong Cui, Bing Pan, Yanbin Sun, "A Survey of Privacy-Preserving Techniques for Blockchain," in Artificial Intelligence and Security, New York, USA, 2019.
- [12] T. K. Sharma, "Permissioned And Permissionless Blockchains: A Comprehensive Guide," Blockchain Council, 13 November 2019. [Online]. Available: <https://www.blockchain-council.org/blockchain/permissioned-and-permissionless-blockchains-a-comprehensive-guide/>. [Accessed 23 March 2020].
- [13] Qi Feng, Debiao He, Sherali Zeadally, Muhammad Khurram Khan, Neeraj Kumar, "A survey on privacy protection in blockchain system," Journal of Network and Computer Applications, vol. 126, pp. 45-58, 2019.
- [14] Jorge Bernal Bernabe ; Jose Luis Canovas ; Jose L. Hernandez-Ramos ; Rafael Torres Moreno ; Antonio Skarmeta , "Privacy-Preserving Solutions for Blockchain.," IEEE Access, vol. 7, pp. 164908 - 164940, 2019.
- [15] R. Zhang, R. Xue, and L. Liu, "Security and privacy on blockchain.," CoRR (to appear in ACM Computing Survey), vol. abs/1903.07602, 2019.
- [16] H. S. a. S. W., "How to Time-Stamp a Digital Document.," Journal of Cryptology , pp. 99-112, 1991.
- [17] G. Wood, Ethereum: A secure decentralised generalised transaction ledger, 2014.
- [18] J. D. M Sharples, "The Blockchain and Kudos: A Distributed System for Educational Record, Reputation and Reward.," in Adaptive and adaptable learning, 2016.
- [19] D. Skiba, "The potential of Blockchain in education and health care.," Nursing Education Perspectives, vol. 38, no. 4, pp. 220-221, 2017.
- [20] M. S. Ali, M. Vecchio, M. Pincheira, K. Dolui, F. Antonelli, and M. H. Rehmani, "Applications of blockchains in the Internet of Things: A comprehensive survey," IEEE Communication Surveys and Tutorials, vol. 21, no. 2, pp. 1676-1717, 2019.
- [21] M. A. P. I. 2.-0. O. I. 2.-0. p. 1.-6. V. M. H. Miraz, "Applications of Blockchain Technology beyond cryptocurrency.," Annals of Emerging Technologies in Computing (AETIC), vol. 2, pp. 1 - 6, 2018.
- [22] C. Esposito, A. De Santis, G. Tortora, H. Chang and K. R. Choo, "Blockchain: A Panacea for Healthcare Cloud-Based Data Security and Privacy," IEEE Cloud Computing, vol. 5, no. 1, pp. 31-37, 2018.
- [23] Ekblaw, A.; Azaria, A.; Halamka, J.D.; Lippman, A., "A Case Study for Blockchain in Healthcare: "MedRec" prototype for electronic health records and medical research data," in IEEE Open and Big Data Conference, Vienna, 2016.
- [24] S. & M. R. Abeyratne, "Blockchain ready manufacturing supply chain using distributed ledger," International Journal of Research in Engineering and Technology, vol. 5, no. 9, pp. 1-10, 2016.
- [25] J. F Calzadilla, A. Villa, "Systematic Literature Review of the use of Blockchain in Supply Chain.," in 12th European Research Seminar (ERS) On Logistics and SCM, Barcelona, 2017.
- [26] D. & F. T. Bhowmik, "The multimedia blockchain: A distributed and tamper-proof media transaction framework.," in 22nd International Conference on Digital Signal Processing (DSP) , 2017.
- [27] Zibin Zheng ; Shaoan Xie ; Hongning Dai ; Xiangping Chen ; Huaimin Wang, "An Overview of Blockchain Technology: Architecture, Consensus, and Future Trends.," in IEEE International Congress on Big Data (BigData Congress), 2017.
- [28] M. S. Henriques and N. K. Vernekar, "Using symmetric and asymmetric cryptography to secure communication between devices in IoT", Proc. Int. Conf. IoT Appl. (ICIOT), pp. 1-4, May 2017.
- [29] T. N. a. H. Hartenstein, "Network Layer Aspects of Permissionless Blockchains.," IEEE Communications Surveys & Tutorials, vol. 21, no. 1, pp. 838-857, 2019.
- [30] C. F. Martin Harrigan, "The Unreasonable Effectiveness of Address Clustering," in IEEE Conferences on Ubiquitous Intelligence & Computing, Advanced and Trusted Computing, Scalable Computing and Communications, Cloud and Big Data Computing, Internet of People, and Smart World Congress, 2016.
- [31] Elli Androulaki, Ghassan O. Karame, Marc Roeschlin, Tobias Scherer, Srdjan Capkun, "Evaluating User Privacy in Bitcoin," in Financial Cryptography and Data Security, Japan, Springer, 2013, pp. 34-51.
- [32] Malte Möser\*, Kyle Soska, Ethan Heilman, Kevin Lee, Henry Heffan, Shashvat Srivastava, Kyle Hogan, Jason Hennessey, Andrew Miller, Arvind Narayanan, and Nicolas Christin, "An Empirical Analysis of Traceability in the Monero Blockchain," Proceedings on Privacy Enhancing Technologies, vol. 2018, no. 3, p. 2018, 143–163.
- [33] P. Koshy, D. Koshy, P McDaniel, "An analysis of anonymity in Bitcoin using P2P network traffic," in International Conference on Financial Cryptography and Data Security., 2014.
- [34] A. Biryukov, D. Khovratovich, and I. Pustogarov, "Deanonymisation of clients in Bitcoin P2P network," in ACM Conference on Computer and Communications Security, 2014.
- [35] A. Biryukov, Tikhomirov, Sergei, "Deanonymization and linkability of cryptocurrency transactions based on network analysis," Proceedings of 2019 IEEE European Symposium on Security and Privacy (EuroS&P), 2019.

- [36] D. Nick, *Data-Driven De-Anonymization in Bitcoin*, Zurich: Swiss Federal Institute of Technology, 2015.
- [37] Dorit Ron, Adi Shamir, "Quantitative Analysis of the Full Bitcoin Transaction Graph," in *Financial Cryptography and Data Security*, SpringerLink, 2013, pp. 6-24.
- [38] S. Goldfeder, H. Kalodner, D. Reisman, and A. Narayanan, "When the cookie meets the blockchain: Privacy risks of Web payments via cryptocurrencies," *Proceedings of Privacy Enhancing Technol*, vol. 2018, no. 4, pp. 179-199, 2018.
- [39] Mercer, N. T. Courtois and R., "Stealth address and key management techniques in blockchain systems," in *International Conference on Information Systems Security and Privacy*, Porto, 2017.
- [40] Amrit Kumar, Clement Fischer, Shruti Tapole, Prateek Saxena, "A Traceability Analysis of Monero's Blockchain," in *European Symposium on Research in COmputer Security*, Oslo, Norway, 2017.
- [41] Chaum, D.L., "Untraceable electronic mail, return addresses, and digital pseudonyms," *Communication ACM*, vol. 24, no. 2, p. 84-90., 1981.
- [42] G. Maxwell, "CoinJoin: Bitcoin Privacy for the Real World.," *Bitcoin Forum*, 2013. [Online]. Available: <https://bitcointalk.org/index.php?topic=279249>.
- [43] Jan Henrik Ziegeldorf, Fred Grossmann, Martin Henze, Nicolas Inden, and Klaus Wehrle., "Coinparty: Secure multi-party mixing of bitcoins.," in *Proceedings of the 5th ACM Conference on Data and Application Security and Privacy*, 2015.
- [44] Licheng Wang, Xiaoying Shen, Jing Li, Jun Shao, Yixian Yang, "Cryptographic primitives in blockchain," *Journal of Network and Computer Applications*, vol. 127, pp. 43-58, 2019.
- [45] R. Rivest, A. Shamir, and Y. Tauman., "How to leak a secret," *Asiacrypt 2001*, vol. 2248, p. 552-565, 2001.
- [46] Yifan Wu., "An E-voting System based on Blockchain and Ring Signature," *University of Birmingham*, 2017.
- [47] Liu, J.K., Wei, V.K., Wong, D.S., "Linkable spontaneous anonymous group signature for ad hoc groups (extended abstract). In: *Information Security and Privacy*," in 9th Australasian Conference, ACISP, Sydney, Australia, 2004.
- [48] Fujisaki, E., Suzuki, K., "Traceable ring signature. In: *Public Key Cryptography*," in 10th International Conference on Practice and Theory in Public-key Cryptography, Beijing, China, 2007.
- [49] Back, A., "Bitcoins with Homomorphic Value (Validatable but Encrypted).," 2015. [Online]. Available: <https://bitsharestalk.org/index.php/topic,16797.msg214814.html#msg214814>.
- [50] N. van Saberhagen, "Cryptonote V 2.0.," 2013. . [Online]. Available: <https://cryptonote.org/whitepaper.pdf>. [Accessed 11 December 2019].
- [51] Maxwell, G., "Confidential transactions.," 2017. [Online]. Available: [https://people.xiph.org/greg/confidential\\_values.txt](https://people.xiph.org/greg/confidential_values.txt). [Accessed 11 December 2019].
- [52] Danny Yang Jack Gavigan Zooko Wilcox'O'Hearn, "Survey of confidentiality and privacy preserving technologies for blockchains," [Online]. Available: [https://z.cash/static/R3\\_Confidentiality\\_and\\_Privacy\\_Report.pdf](https://z.cash/static/R3_Confidentiality_and_Privacy_Report.pdf).
- [53] Noether, S., "Ring Signature Confidential Transactions for Monero," *IACR Cryptology*, p. 1098, 2015.
- [54] S Goldwasser, S Micali, and C Rackoff., "The Knowledge Complexity of Interactive Proof-systems.," 1985.
- [55] Miers, I, Garman, C., Green, M., Rubin, A.D., , "Zerocoin: anonymous distributed e-cash from bitcoin," in *IEEE Symposium on Security and Privacy*, Berkeley, CA, USA, 2013.
- [56] Ben-Sasson, E., Chiesa, A., Garman, C., Green, M., Miers, I., Tromer, E., Virza, M., "Zerocash: decentralized anonymous payments from bitcoin," in *IEEE Symposium on Security and Privacy*, Berkeley, CA, USA, 2014.
- [57] George Kappos, Haaron Yousaf, Mary Maller, Sarah Meiklejohn, "An Empirical Analysis of Anonymity in Zcash," in *27th USENIX Security Symposium*, USA, 2018.
- [58] Seres, István András et al., "MixEth: efficient, trustless coin mixing service for Ethereum," *IACR Cryptology ePrint Archive*, 2019.
- [59] G. 2. Maxwell, "Coinswap: Transaction Graph Disjoint Trustless Trading.," October 2013.
- [60] Ruffing, T., Moreno-Sanchez, P., Kate, A., "Coinshuffle: practical decentralized coin mixing for bitcoin.," in *European Symposium on Research in Computer Security*, 2014.

# Efficient Method for Three Loop MMSE-SIC based Iterative MIMO Systems

Zuhaibuddin Bhutto<sup>1\*</sup>, Saleem Ahmed<sup>2</sup>, Syed Muhammad Shehram Shah<sup>3</sup>  
Azhar Iqbal<sup>4</sup>, Faraz Mehmood<sup>5</sup>, Imdadullah Thaheem<sup>6</sup>, Ayaz Hussain<sup>7</sup>

Computer System Engineering Department, Balochistan University of Engineering & Technology, Pakistan<sup>1</sup>

Computer System Engineering Department, Dawood University of Engineering & Technology, Pakistan<sup>2</sup>

Software Engineering Department, Mehran University of Engineering & Technology, Pakistan<sup>3</sup>

Department of Basic Sciences and Mathematics, Dawood University of Engineering & Technology, Pakistan<sup>4,5</sup>

Energy Systems Engineering Department, Balochistan University of Engineering & Technology, Pakistan<sup>6</sup>

College of Information and Communication Engineering, SungkyunKwan University, Suwon, Republic of Korea<sup>7</sup>

**Abstract**—Iterative decoding is one of the promising methods to improve the performance of MIMO systems. In iterative processing channel decoder and MIMO detector share the information in order to enhance the overall system performance. However, iterative processing requires a lot of computations therefore it is considered as a computationally complex approach due to complex detection schemes involving iterative processing. There are several promising detection methods that require further improvements and they can be candidates in order to practically implement iterative processing. In this paper, the propose method to improve the efficiency of three loop-based minimum mean squared errors with soft interference cancellation (MMSE-SIC) method by reducing its complexity with a single inverse operation. Simulation results are given in order to provide detail analysis of the proposed MMSE-SIC based approach for iterative detection and decoding (IDD).

**Keywords**—MIMO; Iterative Detection and Decoding (IDD); sphere decoding; Minimum Mean Squared Errors with Soft Interference Cancellation (MMSE-SIC)

## I. INTRODUCTION

To acquire a higher data rate, MIMO techniques are widely used in most current wireless communication systems. The channel coding or forward error correction (FEC) scheme is an important part of MIMO communication systems if one targets high QoS for mobile users. It is essential to exploit high-performance FEC methods to achieve the performance gains in MIMO based communication systems. The FEC methods like turbo codes and LDPC codes [1] promises to come close to the Shannon capacity limit. The harsh channel conditions demand to use FEC schemes with iterative decoding to achieve the performance goals. Turbo codes are one of the coding schemes that are based on the concept of iterative decoding [1].

The iterative decoding approach can be employed as the outer loop, which is a connection between the MIMO detector and FEC decoder. In such scenario, the iterative loop between the MIMO detector and FEC decoder utilizes extrinsic LLRs iteratively [2]-[4]; therefore, it is known as joint iterative detection and decoding (JIDD). One major implementation difficulty of the JIDD based MIMO systems is the signal

detection issue at the receiving side and its computational complexity which makes it impractical.

Previous research has employed reduced search MAP methods and equalization method for SISO detectors. The sphere decoding and tree search-based methods have been very promising approaches for JIDD systems [4]-[6]. These approaches target to reduce the MAP search space by finding the likely candidates. However, tree search methods are still computationally complex. In [7][8], a minimum mean squared error with soft interference cancellation (MMSE-SIC) detector is derived from the MMSE detector with the interference cancelation as the pre-process by considering a priori information from the channel decoder. The MMSE-SIC detection method does not provide promising performance improvement in the JIDD system, but they benefit from low complexity. There have been considerable performance enhancements in the MMSE-SIC method. In [9], the author used two approaches to enhance the performance of the MMSE-SIC method. First, a posteriori information is used to enhance the performance of the MMSE-SIC method and complexity is reduced by hard decision threshold (HDT) method. However, still, there is a performance gap between MMSE-SIC and other existing methods like tree search-based methods. In [10] author employed a three-loop approach to further enhance the performance of the MMSE-SIC method. The third loop enhances the performance of the method with the expenses of additional processing. The main complexity of the MMSE-SIC method lies in its filtering matrix inversion method which is addressed in [11] by proposing a single matrix inversion method.

In this paper, the proposed method is used to reduce the complexity of the three-loop method which involves repeated filtering matrix inversion for each transmitted layer. The purpose to use a single inversion of filtering matrix for all layers in the three-loop method. Simulation results show that the proposed method is more suitable with negligible performance degradation with only single inversion instead of layer-based inversion of the filtering matrix.

The rest of the research paper is outlined as follows. The three-loop method is explained with its system diagram in Section II. Then Section III is about the proposed single

\*Corresponding Author

inversion approach for the three-loop MMSE-SIC method. Section IV is about simulation results. Finally, the paper is concluded in the conclusion section.

## II. CONVENTIONAL DETECTION METHODS FOR JIDD SYSTEMS

The system diagram of the conventional JIDD based MIMO system is shown in Fig. 1. The MIMO system of size is considered. The system diagram shows that information bits,  $\mathbf{u}$  are sent to FEC encoder to produce the coded output represented as  $\mathbf{c}$ . The size of each codeword length is of  $n$ . Then bit interleaving is done with codewords having a size of  $M \times K$  resulting in interleaved bits  $\mathbf{x}$ .

Then the bits are divided into MIMO frames of size  $M \times K$  bits. Each MIMO frame containing transmitting bits can be represented as:

$$\mathbf{s} = [s_{1,1}, \dots, s_{1,K}, s_{2,1}, \dots, s_{M,K}] \quad (1)$$

where  $s_{m,k}$  is the  $k$ th bit which is mapped onto the  $m$ th transmitting symbol.

The transmitted information vector is represented as,  $\mathbf{x} = [x_1, x_2, \dots, x_M]^T$ , where symbols are from constellation,  $\Omega$ ,  $\mathbf{x} \in \Omega^M$ . The received information vector is represented as,  $\mathbf{y} = [y_1, y_2, \dots, y_M]^T$ , it can be represented with an  $N \times M$  complex channel matrix, as follows:

$$\mathbf{y} = \mathbf{H}\mathbf{x} + \mathbf{n} \quad (2)$$

where  $\mathbf{n}$  is an  $N \times 1$  complex Gaussian noise vector. The entries of channel matrix  $\mathbf{H}$  are assumed to be known at the receiver.

There are several existing detection methods which are employed in JIDD system. The optimum soft MAP demodulator calculates the exact LLR values in JIDD systems. However it is impractical to use it due to its computational complexity. The LSD is also one of the methods to generate soft information based on SD algorithm. Another approach which searches for the ML solution and its counter hypothesis for soft output is STS [5]. Below are the details of existing detection methods employed in JIDD systems.

### A. MAP based Detection for JIDD System

Fig. 1 illustrates the overall MIMO detector system divided into transmitter and receiver. As depicted in Fig. 1, at receiver, first the soft bit information (SBI),  $L$ , is estimated from  $\mathbf{y}$  for all transmitted bits. Using the MAP detection process,  $L(x_{m,k})$ , is calculated as:

$$L(x_{m,k}) = \ln \left( \frac{P(x_{m,k} = +1 | \mathbf{y})}{P(x_{m,k} = -1 | \mathbf{y})} \right) \quad (3)$$

where,  $k, m$  are the  $k$ th bit of the  $m$ th symbol for which SBI information is calculated.

By applying Bayes' theorem, and use of interleaving operation which makes all information bits in a symbol vector statistically independent and under the max-log approximation, the  $L(x_{m,k})$  can be represented as [4]:

$$L(x_{m,k}) = \ln \left( \frac{\sum_{s \in \mathcal{X}_{m,k}^+} P(\mathbf{y} | x) \cdot P(x)}{\sum_{s \in \mathcal{X}_{m,k}^-} P(\mathbf{y} | x) \cdot P(x)} \right) \approx \max_{x \in \mathcal{X}_{m,k}^+} d_x - \max_{x \in \mathcal{X}_{m,k}^-} d_x \quad (4)$$

where  $d_x$  can be found as follows:

$$d_x = -\frac{1}{N_0} \|\mathbf{y} - \mathbf{H}\mathbf{x}\|^2 + \frac{1}{2} \sum_{m,k} x_{m,k} L_a^d(x_{m,k}). \quad (5)$$

The performance of the JIDD system employing MAP detection method is illustrated in Fig. 2. The performance is generated for  $2 \times 2$  and  $4 \times 4$  MIMO detector with 16 QAM modulation scheme. The performance is compared for different outer iterations which are between decoder and MIMO detector. The code rate of 1/3 is set for the encoding and decoding purposes with the constraint length of 3 is set for each RSC component. The turbo decoder was employed with 8 inner iterations. As it can be seen that as outer iterations increases the performance of  $4 \times 4$  system approaches  $2 \times 2$  MIMO system.

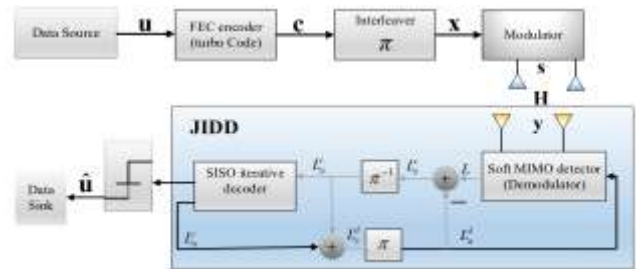


Fig. 1. Transmitter and Receiver Blocks for MIMO System with JIDD.

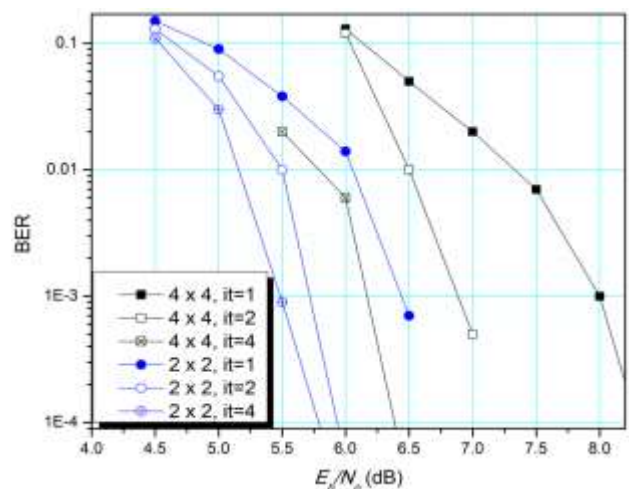


Fig. 2. BER Performance of JIDD based MIMO System Employing MAP Detector for the  $2 \times 2$  and  $4 \times 4$  System.



**B. Sphere Decoding Based Detection Methods**

The SD scheme aims to search the transmitted information by searching the ML metric with minimum distance from received signal, which is considered as the final solution metric [4][12-14]. SD is a suboptimal solution to find the best solution by considering small space vector within a set radius instead of searching whole space of possibilities in a full space. The LSD is the method which is based on SD algorithm. In LSD, the search is not limited for one best solution which is the case in SD method, the LSD builds a list of  $N_c$  best solutions by searching through the space of radius  $r$ . The candidates whose radius falls inside  $r$  whose are used to build a subset list  $L_{LSD} \subset \chi$  of size  $N_c$  [6]. Even the list is full, LSD method continues its search for better candidates until it reaches the end of the search space. The LSD method performance varies for different sizes of  $N_c$ , therefore, the increased size of candidate the performance of LSD methods enhances. However, a large size of  $N_c$  will result in higher candidate search complexity due to increase in search space.

Fig. 3 shows the BER performance [14-16] of Conventional JIDD detection methods over  $2 \times 2$  MIMO system with a 16-QAM scheme. The turbo decoder is set with 8 inner iterations and there are 3 outer iterations performed between decoder and MIMO detector. The full search MAP can produce optimal performance compared to LSD and STS methods. Compared to the conventional LSD with 40 candidates, the STS method produces better performance. However performance of LSD depends on the list size which can be improved with increase in list size with the cost of computational complexity.

**C. Three Loop MMSE-SIC Method**

As depicted in Fig. 4, loops 1 and 2 are those which the JIDD system conventionally employs. In order to reduce the error propagation (EP) which occurs in MMSE-SIC based JIDD systems due to utilization of soft information of other layers, the additional Loop 3 is used and it can enhance the error-rate performance by reflecting soft information within the MIMO detector; so it uses 1) a priori soft information feedback by the FEC decoder,  $L_a^d$ ; 2) a posteriori soft information which is generated from the first  $L^0$  and second  $L^1$  inner MMSE-SIC iterations. During the execution of Loop 3, the MIMO detector estimates  $L^0$  and it is stored so that it can be utilized in the second inner MMSE-SIC iteration jointly with  $L^1$ . While SIC-MMSE detector performing the first inner iteration, the  $i^{th}$  streams soft symbol is calculated by utilizing  $L_a^d$  and  $L^0$  [10]. After completing the first iteration, generated soft information is feedback within the MMSE-SIC detector, both soft information extracted in the first MMSE-SIC iteration, and the a posteriori soft

information  $L^1$  found at the second MMSE-SIC process, are incorporated. This process continues with the detector until there is no further change in performance and finally generated a posteriori information in the last iteration is de-interleaved and forwarded to the channel decoder [10].

In Fig. 5, the impact of Loop 3 which is inside MMSE-SIC detector is elaborated. The third loop utilizes more reliable information during different steps involved in MMSE-SIC method. Fig. 2 illustrates the impact of Loop 3 on BER performance of  $4 \times 4$  MIMO system. The number shown in parentheses represents the total layers utilizing information of previous detection iteration. In Fig. 5, the “3 loop method (1)” where (1) indicates information from previous iteration of Loop 3 is utilized only at first layer during current detection process. Similarly, (2) indicates that information is utilized in 2 current layers and so on. Therefore, as more layers use the information of previous loop from Loop 3, the BER performance is improved because information of previous layers is more reliable and it can reduce the error propagation resulting in performance enhancement.

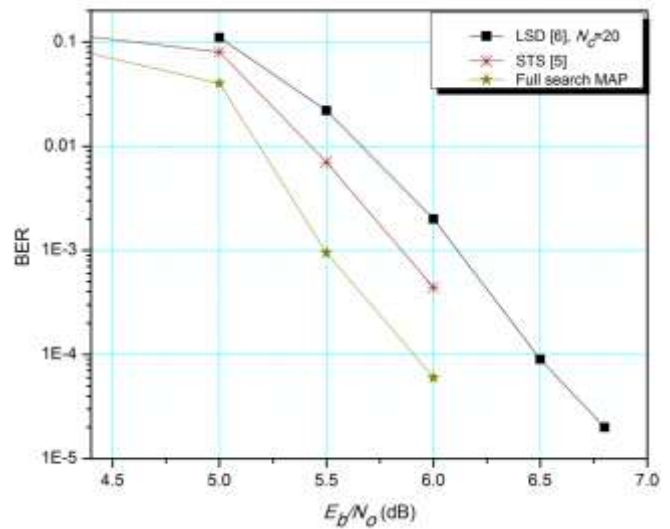


Fig. 3. BER Performance of Various Nonlinear JIDD Detection Methods Over  $2 \times 2$  MIMO System with 16-QAM Modulation.

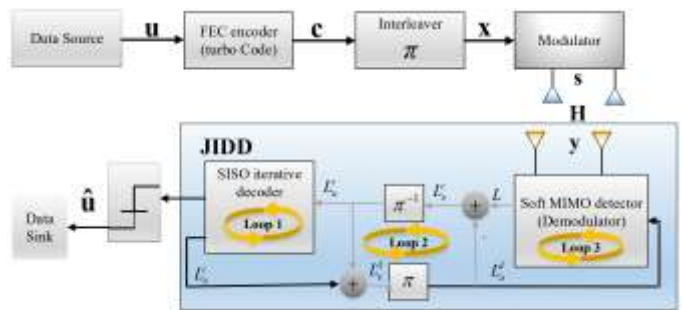


Fig. 4. MMSE-SIC based JIDD System having Three Iterative Loops.

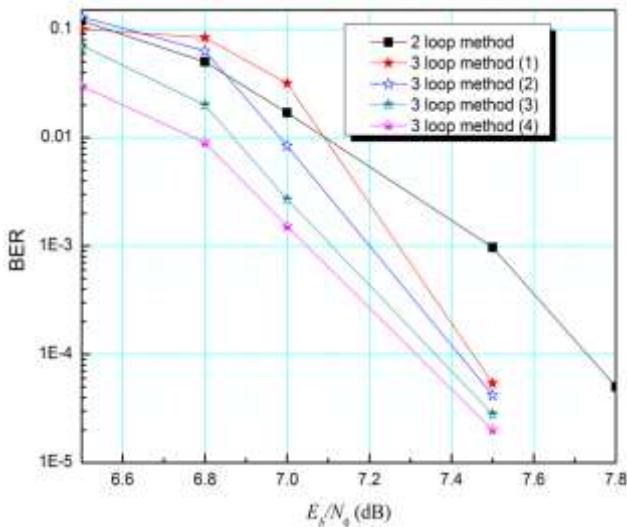


Fig. 5. BER Performance of Three Loop Method for  $4 \times 4$  MIMO System.

### III. SINGLE INVERSION BASED THREE LOOP METHOD

Three loop approach for MMSE-SIC scheme can improve the performance but the cost of computational complexity. Each iteration in the third loop involves  $M$  inverse operations of the filtering matrix. The single inversion method [11] can be employed together with three loops MMSE-SIC method to improve the overall efficiency of the system in terms of computational complexity with a negligible impact on performance. In this paper, the proposed method is used to reduce the complexity of the three-loop MMSE-SIC method by performing the filtering operation with a single inversion for each loop three iterations. By using the single filtering matrix inversion in each iteration of loop 3, the large amount of complexity is reduced while keeping the advantage of performance enhancement brought by the loop 3 approach. During the  $l$ th iteration inside Loop 3, the soft symbols  $\tilde{x}_j^l$  for each later are estimated by utilizing different soft information in each of its iterations.

For the first iteration when  $l = 0$ , the soft symbols are calculated using, which is soft information (a priori) provided by the channel decoder. During the first iteration in loop 3, the soft symbols for  $j = 1; \dots; M$ ; are found as:

$$\tilde{x}_j^0 = E[x_j] = \sum_{q \in \Omega} q \prod_{k=1}^K \frac{1}{2} \left( 1 + \tilde{s}_{j,k} \tanh(L_a^d(s_{j,k})) \right) \quad (6)$$

where  $\tilde{s}_{j,k}$  is put as  $\pm 1$ , which depend on the complex symbol  $q$  taken from the complex space  $\Omega$ . The error,  $e_j = x_j - \tilde{x}_j$ , which can impact the overall performance of the system due to its impact in error propagation is the error calculated between the transmitted stream  $x_j$  and its soft symbol  $\tilde{x}_j$ . The variance of each soft symbol depicts its reliability, which is found using [7][8]:

$$V[j,0] = E\left[|e_j|^2\right] = \sum_{q \in \Omega} |q|^2 \prod_{k=1}^K \frac{1}{2} \left( 1 + \tilde{s}_{j,k} \tanh(L_a^d(s_{j,k})) \right) - |\tilde{x}_j^0|^2. \quad (7)$$

In loop 3, only a priori information is used for soft symbol and their variance calculation. In the successive iterations,  $l > 0$ , soft symbols and their variance are found by using both a posteriori information of previous loop 3 iterations. Therefore, for successive iteration, found as:

$$\tilde{x}_j^l = \sum_{q \in \Omega} q \prod_{k=1}^K \frac{1}{2} \left( 1 + \tilde{s}_{j,k} \tanh(L^{l-1}(s_{j,k}) - L_a^d(s_{j,k})) \right) \quad (8)$$

and their variance is calculated using:

$$V[j,l] = E\left[|e_j^l|^2\right] = \sum_{q \in \Omega} |q|^2 \prod_{k=1}^K \frac{1}{2} \left( 1 + \tilde{s}_{j,k} \tanh(L^{l-1}(s_{j,k}) - L_a^d(s_{j,k})) \right) - |\tilde{x}_j^l|^2. \quad (9)$$

After calculating the soft symbols and their variance we estimate Gram matrix and matched filter output. The interference cancelation for the layer is done using:

$$\tilde{\mathbf{y}}_i^l = \mathbf{y}^{MF} - \sum_{j,j \neq i} \mathbf{g}_j \tilde{x}_j^l = \mathbf{h}_i x_i + \tilde{\mathbf{n}}_i^l, \quad (10)$$

where  $\tilde{\mathbf{n}}_i^l = \sum_{j,j \neq i} \mathbf{h}_j e_j^l + \mathbf{n}$ , and  $\mathbf{g}_j$  is the  $j^{\text{th}}$  column of  $\mathbf{G}$ , and  $\tilde{x}_i^l = [\tilde{x}_1^l, \dots, \tilde{x}_{i-1}^l, 0, \tilde{x}_{i+1}^l, \dots, \tilde{x}_M^l]^T$ .

The next step is filtering the MMSE matrix which involves  $M$  inverse operation (one for each layer). In order to perform single inversion for each iteration in the loop 3, we use the layer independent variance matrix,

$$\Sigma^l = \text{diag}\{V[1,l], \dots, V[M,l]\}. \quad (11)$$

Then the MMSE filtering matrix,  $\mathbf{W}^l$ , is calculated only once in each iteration of loop 3 by using:

$$\mathbf{W}^l = (\mathbf{G} \Sigma^l + \sigma^2 \mathbf{I}_N)^{-1}. \quad (12)$$

It is noticeable that previously three-loop approach for MMSE-SIC used (13) during the filtering process which is the layer dependent variance matrix and is found as:

$$\Sigma_i^l = \text{diag}\{V[1,l], \dots, V[i-1,l], 1, V[i+1,l], \dots, V[M,l]\}. \quad (13)$$

The symbol,  $\hat{x}_i^l$  is calculated by using:

$$\hat{x}_i^l = (\mathbf{w}_i^l) \tilde{\mathbf{y}}_i^l. \quad (14)$$

The symbol estimation in equation (14) can be expended as:

$$\hat{x}_i^l = \mu_i^l x_i + \eta_i^l \quad (15)$$

where

$$\mu_i^l = (\mathbf{w}_i^l) \mathbf{g}_i, \quad (16)$$

and denotes Gaussian random variable having a variance  $(\tilde{\sigma}_i^l)^2$ , which is obtained by using:

$$(\tilde{\sigma}_i^l)^2 = \mu_i^l - |\mu_i^l|^2. \quad (17)$$

The further complexity can be reduced by using a single distance calculation method in order to find the soft information for each bit. The hard decision threshold (HDT) scheme is an efficient single distance calculation method which can calculate of the  $k$ th bit using [9]:

$$L^l(s_{i,k}) \approx -2\rho_i^l \tilde{b}_k^l, \quad (18)$$

where  $\tilde{b}_k^l$  depicts the distance between the estimated symbol, and the HDT line for the estimating bit, and is the signal-to-interference-plus-noise ratio (SINR) for the layer.

#### IV. SIMULATION RESULTS

Performance of the proposed method is compared for different MIMO systems over the 16-QAM modulation scheme. The Rayleigh fading channel model is deployed to evaluate the performance of the proposed method. The code rate of the turbo code is 1/3. The constraint length employed was three for each recursive systematic convolutional (RSC) code. Eight iterations are deployed in the turbo decoder. The exchange of information between turbo decoder and MIMO detector was performed 4 times.

Fig. 6 illustrates the performance comparison of a single inversion based three-loop MMSE-SIC method with other conventional methods for a  $2 \times 2$  MIMO system with a 16-QAM modulation technique. The proposed method has outperformed conventional MMSE-SIC method while it has similar performance as conventional loop 3 method while much less complex because of single inversion in each iteration while conventional loop 3 requires  $M$  inversion operations in each iteration.

Fig. 7 illustrates the performance comparison of a single inversion based three-loop MMSE-SIC method with other conventional methods for a  $4 \times 4$  MIMO system with a 16-QAM modulation technique. As shown in Fig. 7, the proposed method has outperformed the conventional MMSE-SIC method while it has similar performance as the conventional loop 3 method while much less complex because of single inversion in each iteration while conventional loop 3 requires  $M$  inversion operations in each iteration.

##### A. Discussion on Complexity

The complexity of the conventional MMSE-SIC method is dominated by its filtering operation which requires  $M$  inverse operations in each iteration. The number of matrix inverse operations is directly proportional to the number of antennas. As the antenna size increases so do the number of matrix inversion operations. Therefore, the proposed method has no impact on the increase in the antenna size because it is layer independent which means it requires only a single inversion operation. Referring to Fig. 7, the conventional loop 3 method in [10] requires eight matrix inversion operation for 2 iterations while the proposed scheme needs only two matrix inversion operation (one for each iteration). Additionally, instead of Max-Log MAP approximation for SBI estimation, we employ a single distance calculation approach which further reduces the complexity of the proposed method.

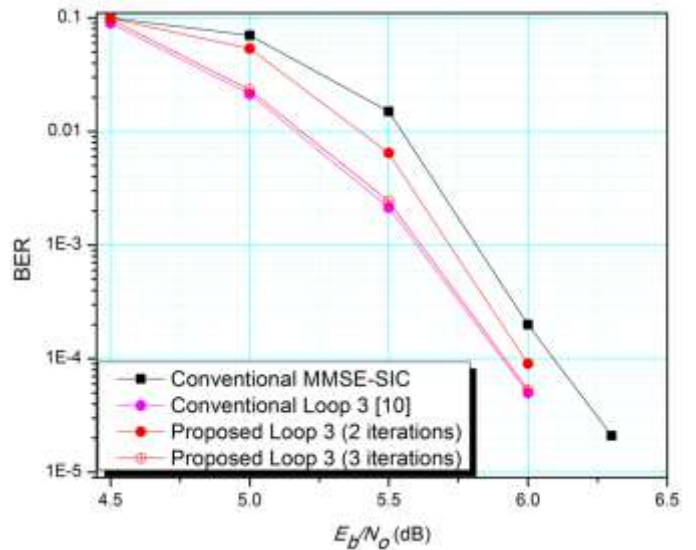


Fig. 6. BER Performance of Proposed Single Inversion based Three-Loop MMSE-SIC Method for a MIMO System with 16-QAM Modulation.

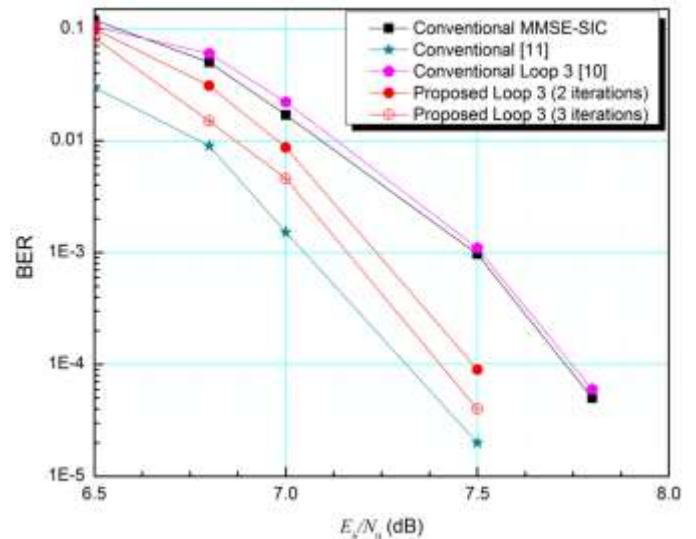


Fig. 7. BER Performance of Proposed Single Inversion based Three-Loop MMSE-SIC Method for a  $4 \times 4$  MIMO System with 16-QAM Modulation.

#### V. FINDINGS AND CONCLUSION

Due to the large computational complexity of multiple matrix inversion in three loop MMSE-SIC method, the proposed single inversion method for three loop-based MMSE-SIC method. The proposed method improves the efficiency of the MMSE-SIC method by lowering the computational complexity. The simulation results are given which shows that the proposed method can improve the performance in each Loop 3 iteration with much lower complexity.

#### ACKNOWLEDGMENT

This work was supported by the Balochistan University of Engineering and Technology, Khuzdar Pakistan, Research Fund.

REFERENCES

- [1] C. Berrou, A. Glavieux, and P. Thitimajshima, "Near Shannon limit error-correcting coding and decoding turbo codes", in Proc. of the IEEE Intern. Conf. on Communications (ICC), Geneva, USA, pp. 1064-1070, May 1993.
- [2] N. Khan, S. Ahmed, D. M. Saqib, "Study of MIMO Detection schemes for Emerging Wireless Communications", in International Journal of Computer Science and Network Security, vol.18 no.3, Mar. 2018.
- [3] M. Zhang, S. Ahmed and S. Kim, "Iterative MMSE based Soft MIMO Detection with Parallel Interference Cancellation" in IET Communications, vol. 11, no. 11, pp. 1775-1781, Sept. 2017.
- [4] B. M. Hockwald, and S. Ten Brink, "Achieving near-capacity on a multiple-antenna channel", IEEE Transaction Communication, vol. 51, no. 3, pp. 389-399, Apr. 2003.
- [5] C. Studer, H. Blcskei, "Soft-input soft-output single tree-search sphere decoding", IEEE Transactions on Information Theory, vol. 56, no. 10 pp. 4827-4842, Oct. 2010.
- [6] S. Ahmed, S. Kim, "Efficient list-sphere detection scheme for joint iterative multiple-input multiple-output detection", IET Communication, vol. 8, no. 18, pp. 3341-3348, Dec. 2014.
- [7] X. Wang, and H. V. Poor, "Iterative (turbo) soft-interference cancellation and decoding for coded CDMA", IEEE Transaction on Communication, vol. 47, no. 7, pp. 1046-1061, Jul. 1999.
- [8] M. T. Tuchler, A. C. Singer, and R. Koetter, "Minimum mean squared error equalization using a priori information", IEEE Transaction on Signal Processing, vol. 50, no. 3, pp. 673-683, Mar. 2002.
- [9] S. Ahmed, S. Kim, "Efficient soft bit estimation for joint iterative multiple-input multiple-output detection", IET Communications, vol. 9, no. 17, pp. 2107-2113, Nov. 2015.
- [10] S. Ahmed and S. Kim, "Efficient SIC-MMSE MIMO detection with three iterative loops", International Journal of Electronics and Communications, vol. 72, pp. 65-71, Feb. 2017.
- [11] C. Studer, S. Fateh, D. Seethaler, "ASIC implementation of soft input soft-output MIMO detection using MMSE parallel interference cancellation", IEEE Journal of Solid-State Circuits, vol. 46, no. 7, pp. 1754-1765, Jul. 2011.
- [12] H. Vikalo and B. Hassibi, "On the expected complexity of sphere decoding", in the Thirty-Fifth Asilomar Conf. on Signals, Systems and Computers, vol. 2, pp. 1051-1055, Nov. 2001, Pacific Grove, CA, USA.
- [13] O. M. Damen, H. E. Gamal, and G. Caire, "On maximum-likelihood detection and the search for the closest lattice point", IEEE Transactions on Information Theory, vol. 49, no. 10, pp. 2389-2402, Oct. 2003.
- [14] Z. Bhutto, *et al.*, "Efficient Method for SBI Estimation in Iterative Coded MIMO Systems", International Journal of Computer Science and Network Security (IJCSNS), vol. 19, no. 6, pp. 135-140, June 2019.
- [15] A. Hussain, Z. Bhutto, "BER Performance of Opportunistic Relaying with Direct Link using Antenna Selection", International Conference on Information Technology (AIT), Bangkok, Thailand, June 2012.
- [16] A. Hussain, Z. Ahmed, and Z. Bhutto, "BER Performance of Opportunistic Relaying Scheme using Antenna Selection", 6th Joint IFIP Wireless and Mobile Networking Conference (WMNC 2013), Dubai, UAE, April 2013.

# An Empirical Comparison of Fake News Detection using different Machine Learning Algorithms

Abdulaziz Albahr<sup>1</sup>, Marwan Albahar<sup>2</sup>

College of Applied Medical Sciences, King Saud bin Abdulaziz University for Health Sciences, Alahsa, Saudi Arabia<sup>1</sup>  
King Abdullah International Medical Research Center, Alahsa, Saudi Arabia<sup>1</sup>  
Umm Al Qura University<sup>2</sup>

**Abstract**—Relying on social networks to follow the news has its pros and cons. Social media websites indeed allow the spread of information among people quickly. However, such websites might be leveraged to circulate low-quality news full of misinformation, i.e., "fake news." The wide distribution of fake news has a considerable negative impact on individuals and society as a whole. Thus, detecting fake news published on the various social media websites has lately become an evolving research area that is drawing great attention. Detecting the widespread fake news over the numerous social media platforms presents new challenges that make the currently deployed algorithms ineffective or not applicable anymore. Basically, fake news is deliberately written on the first place to mislead readers to accept false information as being true, which makes it difficult to detect based on news content solely; consequently, auxiliary information, like user social engagements on social media websites, need to be taken into account to help make a better detection. Using such auxiliary information is challenging because users' social engagements with fake news produce noisy, unstructured, and incomplete Big-Data. Due to the fact that fake news detection on social media is fundamental, this research aims at examining four well-known machine learning algorithms, namely the random forest, the Naïve Bayes, the neural network, and the decision trees, distinctively to validate the efficiency of the classification performance on detecting fake news. We conducted an experiment on a widely used public dataset i.e. LIAR, and the results show that the Naïve Bayes classifier defeats the other algorithms remarkably on this dataset.

**Keywords**—Fake news; classification; machine learning; performance comparison

## I. INTRODUCTION

Many people follow the news through different social media platforms because of their ease of access. For instance, about two-thirds of the Americans follow the news through social media websites [1][2]. Newman et al. [3] reported the increased usage of various digital platforms in Great Britain as the main source of the news feed. Because of circulating the breaking news swiftly, social media platforms are significantly better than traditional media [4]. However, not all posted news items are true. There are many economic, social, and political reasons behind people's manipulation of data and information changing. Therefore, these manipulated data leads to creating news items that are neither totally true nor totally false [5]. This, in turn, leads to misleading information on social media networks that causes several predicaments in society. Such misinformation (also known as "Fake News") has a broad spectrum of types and forms. For example, rumors, fake

advertisements, satires, and false political reports are different types of fake news [1]. The spread of fake news becomes more viral than the true news items [6] urged many researchers to concentrate on innovating efficient automated solutions for detecting fake news [7]. Google has announced a new service named "Google News Initiative" aimed at tracking and eliminating fake news [8]. This project will assist users in distinguishing fake news and reports [9]. In fact, the task of detecting fake news is challenging. A fake news detection model aims at identifying purposely misleading news relying on investigating the previously reviewed fake and real news. This brings us to shed light on the availability of large-scale top-quality training data as one of the cornerstones. The fake news detection framework's task can be considered a simple binary classification or a fine-tuned classification in a challenging setting [10]. After 2017, various fake news datasets were introduced. Researchers sought to improve the deployed models' performance using these different datasets such as (ISOT, Kaggle, and LIAR datasets), which are well-known publicly available datasets [11].

In the current research paper, we compare different machine learning classifiers' performance for detecting fake news. The key contributions of this research paper are as follows:

- A detailed performance analysis of four machine learning algorithms using different NLP techniques for detecting fake news.
- Different machine learning-based models are implemented to detect and classify fake news. Each model's performance is measured to categorize various news items correctly, which revealed each model's ability to improve its accuracy of detecting fake news.

This paper is arranged as follows: Section II presents the related works. The objective of this study is clearly highlighted in Section III. In Section IV, we review different classifiers. Section V will explain the data collection process and provide an analysis of the dataset. In Section VI, we present the experimental setup and the evaluation metrics. In Section VII, the examined models' methodology containing the data preparation and handling the missing data problem is discussed in detail. The experimental results of the implemented models are discussed in Section VIII. A discussion of the obtained results and a conclusion of this study are shown in Section XI and X.

## II. RELATED WORKS

As people tend to consume more news on social media, fake news on social media has emerged as a critical problem that has a negative impact on society and government [25]. An early study on detecting fake news concentrated on detecting rumors on twitter, and these studies were conducted by social scientists [12]. Later, researchers have focused on understanding the structure and characteristics of fake news in order to identify fake news. As a result, numerous approaches for automatic fake news detection have been proposed in the literature. Most of these approaches transform the fake news detection into a binary classification task, where each statement, "i.e., news" is labeled as true or false using various machine learning techniques (e.g., [13][14]) or deep learning based techniques [16]. These approaches require data corpus to correctly detect fake news. Rubin [15] introduced three criteria used to determine the quality of created text corpus for identifying fake news, i.e., all facts included in the dataset must be verified; all facts occur in a specific period (e.g., during US election); the way used to observe the facts must be similar, and the facts must have a different level of impact on society. The text corpus has an advantage that the pre-processing is straightforward and simple. However, it suffers from the following limitation, i.e., the only text analysis will reveal limited clues that are not enough to effectively detect fake news. Therefore, current approaches have integrated information based on the propagation network of news that captures how they spread. Ruchansky [16] introduced a new approach, called CSI that encapsulates three modules: Capture, Score, and Integrate. The capture module captures the temporal patterns of the users with the textual information of the news. Score Module exploits users' profiles to learn their vector representation and computes a score for each user engaged in spreading news. It then combines the output of the previous two modules to classify the news as fake. Singhanian et al. [17] propose an approach that is based on deep learning techniques. In their approach, three layers are used to explore the different levels of the text of news separately (i.e., word level, sentence level, and title level). Liu and WU [18] propose an approach that aims to detect fake news at an early stage by exploiting the propagation network of news. In their approach, news are modeled as a multivariate sequence whose elements represent users involved in spreading news. Users are represented as a vector based on features extracted from their profile. Recurrent Neural Networks (RNN) and Convolutional Neural Networks (CNN) are applied to learn vector representation of news's sequence, which feeds into a multi-layer neural network to classify news as fake or not. Wu and Liu [19] propose an approach based on tracing network of news and using the LSTM-RNN model for classification. Instead of identifying helpful features that apply to detect fake news, Vo and Lee [20] identify people, called guardians, who are interested in correcting fake news and propose a recommendation system that recommends URLs of fact checking to guardians to integrate with fake news. Karimi et al. [21] propose an approach that considers fake news detection a multi-class classification task. In their approach, CNN and LSTM are used to automatically extracted feature vectors from each textual source of news and used an interpretable multi-source fusion model to integrate the learned feature vectors into one vector.

Then, the Multi-class Discriminative Function component is used to determine the class of the fakeness of news. Aghakhani et al. [22] propose an approach called FakeGAN, which uses GAN algorithms to detect false reviews. Goldani et al. [23] propose a method that uses capsule neural networks with word embedding representation to enhance fake news detection performance. Two widely used datasets, i.e., ISOT and LIAR, evaluate the proposed method's performance. Wang et al. [24] propose an end to end approach that teaches common features representation among events and uses them to detect fake news on new events.

## III. RESEARCH OBJECTIVE

There is no doubt that the current political events have led to an increase in fake news circulation. In fact, humans are inconsistent and very poor in detecting fake news. Thus, researchers have exerted their efforts to automate the process of identifying fake news. The most well-known attempts blacklist authors and sources that are unreliable. However, we need to consider more complex cases where reliable authors and sources publish fake news to have a reliable, fully automated detecting solution. Machine learning proves to be useful in detecting language patterns. Hence, this research aims to use different machine learning models to detect language patterns that distinguish between real and fake news.

## IV. BACKGROUND

### A. Naive Bayes

The Naïve Bayes classifier is a classifier based on Bayes' Theorem:

$$P(A|B) = \frac{P(B|A) \times P(A)}{P(B)} \quad (1)$$

Where A and B denote two conditions. The naive Bayes classifier considers each semantic feature as a condition and classifies the samples with the highest occurring probability. Noteworthy, it assumes that the semantic features are independent. Naïve Bayes is considered to be one of the most efficient and effective classification algorithms. It can work on small sample sizes and produces an accurate classification result [25].

### B. Random Forest

A decision tree comprises parents with different conditions branch, in which, each node represents a class for classification. The random forest classifier is an ensemble method for classification which construct a multitude of decision trees. We set parameters such as `n_estimators`, `min_samples_split`, `random_state`, `max_depth` to obtain the best performance. In which, `n_estimators` represents the number of decision trees in the random forest, `min_samples_split` represents the minimum amount of samples to split an internal node, and `max_depth` represents the maximum depth of a decision tree.

### C. Decision Trees

A decision tree is a set of decision nodes that start at the root. The benefits of utilizing a decision tree include easy interpretation, efficient handling of outliers, no need for the linear separation of classes, dependent features. Nevertheless,

the existence of so many sparse features could lead a decision tree to overfit, and thus it performs poorly.

#### D. Artificial Neural Network

A neural network is made of interconnected processing nodes known as neurons that work together to solve very particular problems. Using deep neural networks is considered one of the most successful methods of machine learning. Lately, the new advents of neural networks and pre-trained word embedding have become the main basis of new rich ideas of NLP tasks. Nevertheless, the current model treats all words as a network of input and does not take into account the function of keywords. Consequently, redesigning the neural network model by combining the advantages of the two methods and increasing the weight of keywords in the network could lead to a remarkable improvement.

#### V. DATASET

We use a public dataset in [26], which comprises 12.8K human-labeled short statements from PolitiFact through its API. POLITIFACT.COM editor was applied to each statement to evaluate its validity. Six fine-grained labels for news truthfulness are considered into multiple classes, including false, true, pants-fire, mostly-true, half-true and barely-true. The distribution of labels in this dataset is as follows: a range between 2,063 to 2,638 for multiple labels and 1,050 pants-fire labels. Moreover, the dataset comprises of different metadata. These metadata contain valuable information about the speaker, total credit history count of the speaker, state, subject, party, and job. The total credit history count, including the half-true counts, false counts, pants-fire counts, barely-true counts, mostly-true counts. The statistics of the dataset are listed in Table I. Some selection samples from the dataset are presented in Table II.

TABLE I. THE STATISTICS OF LIAR DATASET [26]

| LIAR Dataset Statistics        |        |
|--------------------------------|--------|
| Training set size              | 10,269 |
| Validation set size            | 1,284  |
| Testing set size               | 1,283  |
| Avg. statement length (tokens) | 17.9   |
| Top-3 Speaker Affiliations     |        |
| Democrats                      | 4,150  |
| Republicans                    | 5,687  |
| None (e.g., FB posts)          | 2,185  |

TABLE II. TWO RANDOM EXCERPTS FROM THE LIAR DATASET [26]

| Sample 1   | Sample 2   |
|--|--|
| <b>Statement:</b> The last quarter, it was just announced, our gross domestic product was below zero. Who ever heard of this? It's never below zero.   | <b>Statement:</b> Under the health care law, everybody will have lower rates, better quality care and better access.   |
| <b>Speaker:</b> Donald Trump   | <b>Speaker:</b> Nancy Pelosi   |
| <b>Context:</b> presidential announcement speech   | <b>Context:</b> on Meet the Press  |
| <b>Label:</b> Pants-fire   | <b>Label:</b> False  |
| <b>Justification:</b> According to Bureau of Economic Analysis and the National Bureau of Economic Research, the growth in the gross domestic product has been below zero 42 times over 68 years. That's a lot more than never. We rate his claim Pants on Fire! | <b>Justification:</b> Even the study which Pelosis staff cited as a source of that the statement suggested that some people would pay more for health insurance. Analysis at the state level found the same thing. The general understanding of the word everybody is every person. The predictions do not back that up. We rule this statement False. |

#### VI. EVALUATION

##### A. Experimental Setup

The experiments of this paper were conducted on a server having 32 GB RAM, GeForce GTX 1080 GPU of 8 GB GDDR5X memory, and 2560 NVIDIA CUDA cores. We used Keras library for implementing the proposed models.

##### B. Evaluation Metrics

To evaluate the models, training and validation accuracy are reported for the data partitions. Accuracy is calculated based on the following mathematical representation. Apart from accuracy, other performance measures, that is, True Positive Rate (TPR) also known as Recall, Precision (Pre), and F1 measures, are calculated based on equations 2, 3, 4 and 5, respectively.

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} \quad (2)$$

$$TPR = \frac{TP}{TP+FN} \quad (3)$$

$$Pre = \frac{TP}{TP+FP} \quad (4)$$

$$F1 = 2 * \frac{TPR * Pre}{TPR+Pre} \quad (5)$$

where FP, TN, TP, and FN denote false positives, true negatives, true positives and false negatives, respectively.

## VII. METHODOLOGY

The LIAR dataset is a well-known dataset in the realm of detecting fake news. It consists of 12,836 human-labeled short statements. The chosen instances in this dataset are from more natural contexts such as political debates, Facebook, tweets posts, etc. Furthermore, it contains 12787 news items. In each item, the following features are provided:

- News statement
- Barely true counts
- Subject of news
- False counts
- Half true counts
- Speaker name
- Speaker's job title
- Mostly true counts
- Pants on fire
- State information
- Party affiliation
- Venue

### A. Data Preparation

We split the categorical (text) features and numerical features into two categories:

- Numerical features are (false counts, mostly true counts, pants on fire, barely true counts, mostly true counts, and half true counts). As we know, we do not need to do pre-processing on the numerical features because these features contain true counts and false counts, so we will use these counts for each news item.
- Categorical features are (Party affiliation, Venue, Subject of news, Speaker name, Speaker's job title, State information, and News statement).

Our primary focus was on feature engineering; if we could add some other features or fine-tune the features, detecting news accuracy can be much efficient. Therefore, we explore all categorical features to extract the best feature that distinguishes true and fake news.

In the party affiliation feature, we extracted the party affiliation's unique parties and then replaced all the different parties into four categories named as republican, democrat, unknown, and others. Consequently, we intend to make feature values closer to the class label of news.

While in other features, we tokenize all the words to work on each word separately. After tokenization, we removed the stop words of English because stop words are not good words that cannot distinguish between true and fake news. Thus, we removed these stop words because they fall into both: true or fake news. Next, we applied stemming on the words/tokens because we wanted to use only the (base/root) form of words.

Although stemming does not work well on all the words, our goal was to convert all these features into categories.

There were unique words in Label/class, speaker, and state info features, so we encoded them into unique numbers (half-true as 0, false as 1, mostly-true as 2, barely true as 3, true as 4, and pants-fire as 5). We encoded the other features categories into unique numbers as well.

Finally, in the Statement of news feature, we removed the punctuation from the Statement's sentences. We then removed the repeating characters; we also clean hyperlinks and other special characters from the text of statement news. After cleaning the news statement, we applied unigram feature extraction, bigram features extraction, and trigram features extraction. We observed that the trigram feature yields good results, among others (see Fig. 1).

### B. Missing Data

We conduct an investigation to check the missing value because it can affect the overall performance of algorithms. Thus, we found that 3565 speaker job-titles were missing, 2747 state information missing, and 129 venues. Initially, we replaced the missing values with NaN, and after that, we replaced these with unknown words. It is worth mentioning that, based on our observation, using our method to handle the missing value does not make a significant difference or even any difference.

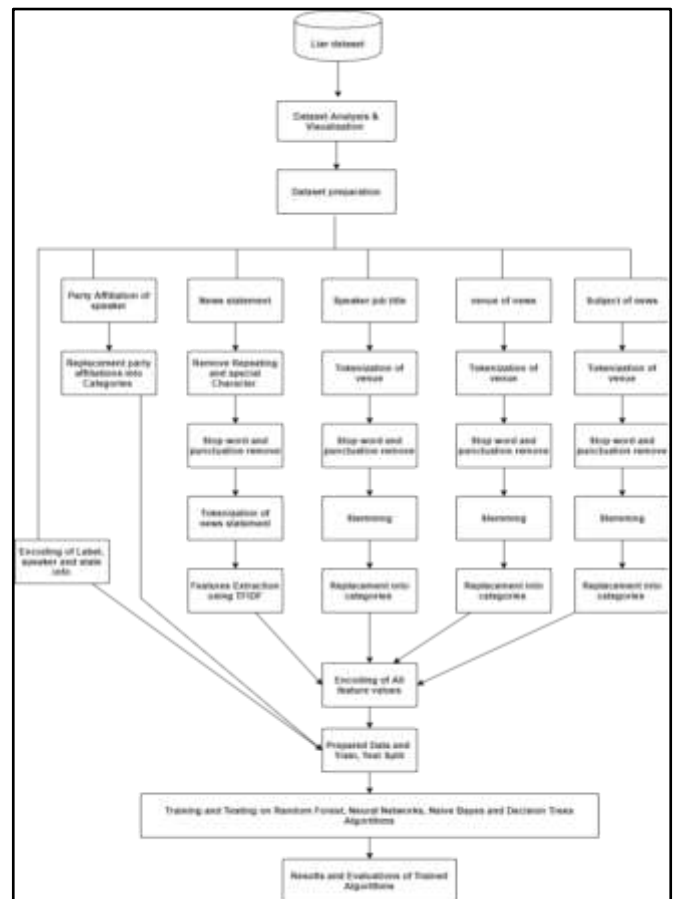


Fig. 1. The High-Level Architecture of the Entire Process.



C. Overfitting and Cross Validation

Overfitting is one of the central problems in machine learning. It arises when the model performs poorly on unseen data while giving excellent results on training data. Cross-validation is a way to overcome such an issue; it aims to test the model's ability to correctly predict new data that was not used in its training. Cross-validation shows the model generalization error and performance on unseen data. K-fold cross-validation is one of the most popular versions. In our experiment, we use k-fold cross-validation to ensure we avoid overfitting.

VIII. EXPERIMENTAL RESULT

A. Training and Testing

There were three separate files in this dataset. We combined all three files and pre-processed the data. Then, we prepared it for the training and testing sets. We divided the data 70-30% for training and testing sets. Next, we used cross-validation 5 times, with 80-20% split every time for training and testing. It is important to note that we shuffled all the rows with random state 6 to avoid any train models' biases. We employed four different machine learning algorithms, and we used Python 3.6.5 as our programming language for the implementation. The classification models that we implemented are the random forest, Naïve Bayes, neural network, and decision trees algorithms to explore further how well our data fit into the models. These algorithms are suitable for several classifications as they have their own properties. We observed different variations during the training of machine learning algorithms. We tune the parameters of each algorithm and get different results. Still, k-fold validation techniques can be applied. Accordingly, we can use the data in different folds to train the algorithms each time on a different training set. After training, we used the trained models of each algorithm and tested them on the testing set.

B. Result

To evaluate the classification process of each algorithm, standard metrics that measure the overall performance were considered. The number of predictions (whether correct or incorrect) with each class are shown in the confusion matrix (see Fig. 2). Fig. 2 shows the confusion matrix of each classifier on the test set. Naïve Bayes classifies all classes with more accuracy. Other classifiers encounter some difficulties classifying mostly true and pants-fire. For these labels, detecting the correct label is more challenging, and many pants-fire texts are predicted as false.

Nevertheless, it is challenging to distinguish between false and barely-true and between mostly-true and true. To have a thoroughly detailed analysis, we evaluated each model's performance within and across different K-fold (see Tables III to VII). This allowed us to further study the overall performance of each algorithm and see the generalization as well.

Fig. 3 presents the overall comparison of all algorithms. Among all the classifiers we have implemented, Naïve Bayes gives the highest accuracy. Other classifiers, such as random forest, appeared to be vulnerable to overconfidence due to the

usage of independent variables to predict outcomes. Noteworthy, it requires that each data point needs to be independent of all other data points. In this dataset, the news statement determines the feature word length.

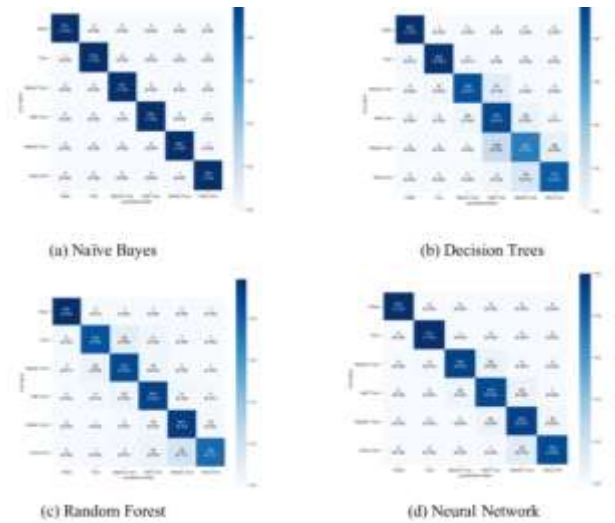


Fig. 2. Confusion Matrix of Classification using different Machine Learning Algorithms for LIAR Dataset.

TABLE III. K-FOLD-1 RESULTS 80% TRAINING AND 20% TESTING

| Evaluation measures | Random Forest | Naïve Bayes | Neural Network | Decision Trees |
|---------------------|---------------|-------------|----------------|----------------|
| Accuracy            | 0.90%         | 0.99%       | 0.97%          | 0.94%          |
| Precision           | 0.90%         | 100%        | 0.97%          | 0.93%          |
| Recall              | 0.89%         | 100%        | 0.97%          | 0.95%          |
| F1-Score            | 0.89%         | 100%        | 0.97%          | 0.94%          |

TABLE IV. K-FOLD-2 RESULTS 80% TRAINING AND 20% TESTING

| Evaluation measures | Random Forest | Naïve Bayes | Neural Network | Decision Trees |
|---------------------|---------------|-------------|----------------|----------------|
| Accuracy            | 0.91%         | 0.99%       | 0.89%          | 0.91%          |
| Precision           | 0.91%         | 100%        | 0.91%          | 0.92%          |
| Recall              | 0.90%         | 100%        | 0.90%          | 0.91%          |
| F1-Score            | 0.91%         | 100%        | 0.90%          | 0.92%          |

TABLE V. K-FOLD-3 RESULTS 80% TRAINING AND 20% TESTING

| Evaluation measures | Random Forest | Naïve Bayes | Neural Network | Decision Trees |
|---------------------|---------------|-------------|----------------|----------------|
| Accuracy            | 0.93%         | 0.99%       | 0.94%          | 0.96%          |
| Precision           | 0.92%         | 100%        | 0.95%          | 0.96%          |
| Recall              | 0.91%         | 100%        | 0.95%          | 0.96%          |
| F1-Score            | 0.91%         | 100%        | 0.95%          | 0.96%          |

TABLE VI. K-FOLD-4 RESULTS 80% TRAINING AND 20% TESTING

| Evaluation measures | Random Forest | Naïve Bayes | Neural Network | Decision Trees |
|---------------------|---------------|-------------|----------------|----------------|
| Accuracy            | 0.91%         | 0.99%       | 0.94%          | 0.95%          |
| Precision           | 0.92%         | 100%        | 0.93%          | 0.97%          |
| Recall              | 0.90%         | 100%        | 0.95%          | 0.94%          |
| F1-Score            | 0.91%         | 100%        | 0.94%          | 0.95%          |

TABLE VII. K-FOLD-5 RESULTS 80% TRAINING AND 20% TESTING

| Evaluation measures | Random Forest | Naïve Bayes | Neural Network | Decision Trees |
|---------------------|---------------|-------------|----------------|----------------|
| Accuracy            | 0.92%         | 0.99%       | 0.86%          | 0.76%          |
| Precision           | 0.93%         | 100%        | 0.90%          | 0.81%          |
| Recall              | 0.91%         | 100%        | 0.84%          | 0.77%          |
| F1-Score            | 0.92%         | 100%        | 0.85%          | 0.75%          |

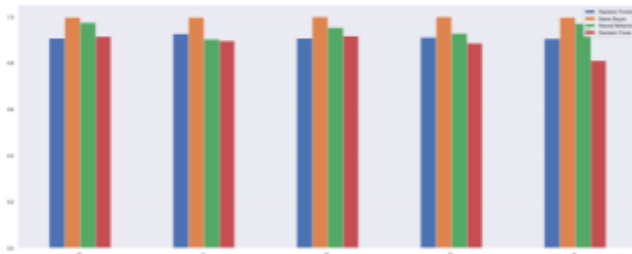


Fig. 3. Overall Accuracy Comparison of different Machine Learning Algorithms for LIAR.

Consequently, the model will tend to overweight the significance of observations when such observations are related. It is worth mentioning that, as we can see from the results, there is no overfitting problem because, in every split, the model performance is consistent. Likewise, we have evaluated the trained models, and we could see the fine details of each algorithm's confusion matrix and that there is no false rate of predictions on class level.

### IX. DISCUSSION

We examined the performance of several machine learning algorithms: the random forest, the Naïve Bayes, the neural network, and the decision trees algorithms. In general, our obtained results verify the pros and cons of the compared different machine learning algorithms when they have been used in detecting fake news. In the following few lines, we analyze the results and give an insight into Naïve Bayes performance:

- In this research work, we intend to train the dataset on different algorithms to determine which algorithm performs well. The reason for the Naïve Bayes algorithm's good performance is that it works well on the text, based on Bayes theorem. Naïve Bayes computes conditional probabilities of two events on the basis of text occurrence individually and differentiates each event/class accordingly.
- Naïve Bayes algorithm is better than other algorithms in this dataset. The accuracy is good as we evaluated the trained model with different evaluation measures, and the converging/training time of the Naïve Bayes is excellent (see Table VIII).

TABLE VIII. THE RUN-TIME FOR THE DIFFERENT MACHINE LEARNING ALGORITHMS (IN SECONDS)

|          | Random Forest | Naïve Bayes | Neural Network | Decision Trees |
|----------|---------------|-------------|----------------|----------------|
| Run Time | 54.3 s        | 1.3 s       | 420.2 s        | 540.1 s        |

- When applied to this dataset, the computational runtime and accuracy comparisons led us to conclude that the Naïve Bayes is the best method in general.

### X. CONCLUSION

Understanding the rationale of specific fake news items infers many details about the different involved factors. Recently, a rapidly increased number of models were proposed in the literature to automatically detect fake news. The two influential factors that significantly impact these models' accuracy are the datasets and a set of explicit classes. Our experiment posits that good models should require a reasonable number of fine-tuning when tested on different datasets. This paper investigates four machine learning classifiers' performance, namely, the random forest, the Naïve Bayes, the neural network, and the decision trees algorithms for identifying fake news. We used a publicly well-known dataset, i.e., LIAR. Based on our results, we observed good performance of the Naïve Bayes algorithm because of the computation of conditional probabilities of two events on the basis of text occurrence individually and the differentiation between each event/class accordingly.

### REFERENCES

- [1] Zhang, X.; and Ghorbani, A. A. (2019). An overview of online fake news: Characterization, detection, and discussion. *Information Processing & Management*, 57, 102025.
- [2] Dale, R. (2017). NLP in a post-truth world. *Natural Language Engineering*, 23, 319-324.
- [3] Newman, N.; Dutton, W.; and Blank, G. (2013). Social media in the changing ecology of news: The fourth and fifth estates in Britain. *International journal of internet science*, 7.
- [4] Bondielli, A.; and Marcelloni, F. (2019). A survey on fake news and rumour detection techniques. *Information Sciences*, 497, 38-55.
- [5] Granik, M.; and Mesyura, V. (2017). Fake news detection using naive Bayes classifier. *Proceedings of 2017 IEEE First Ukraine Conference on Electrical and Computer Engineering (UKRCON)* (pp. 900-903).
- [6] Balmas, M. (2014). When fake news becomes real: Combined exposure to multiple news sources and political attitudes of inefficacy, alienation, and cynicism. *Communication Research*, 41, 430-454.
- [7] Horne, B. D.; and Adali, S. (2017). This just in: Fake news packs a lot in title, uses simpler, repetitive content in text body, more similar to satire than real news. In *Eleventh International AAAI Conference on Web and Social Media*.
- [8] Google News Initiative (2018). <https://newsinitiative.withgoogle.com/>, Retrieved March 25, 2020.
- [9] Google Announcement. (2018). Google announces google news initiative to help quality journalism in digital. Retrieved March 21, 2020, from <http://cnaoe.com/technology/google-announces-googlenews-initiative-to-help-quality-journalism-in-digital-age>.
- [10] Tin, P. T. (2018). A study on deep learning for fake news detection.
- [11] Meel, P.; and Vishwakarma, D. K. (2019). Fake News, Rumor, Information Pollution in Social Media and Web: A Contemporary Survey of State-of-the-arts, Challenges and Opportunities. *Expert Systems with Applications*.
- [12] Vosoughi, S.; Roy, D.; Aral, S. (2018). The spread of true and false news online. *Science* 359, 1146-1151.
- [13] Shu, K.; Sliva, A.; Wang, S.; Tang, J.; and Liu, H. (2017). Fake news detection on social media: a data mining perspective. *SIGKDD Explor.* 19(1), 22-36.
- [14] Zannettou, S.; Sirivianos, M.; Blackburn, J.; and Kourtellis, N. (2018). The web of false information: rumors, fake news, hoaxes, clickbait, and various other shenanigans.
- [15] Rubin, V.; Conroy, N.; Chen, Y.; and Cornwell, S. (2016). Fake news or truth? using satirical cues to detect potentially misleading news.

- Proceedings of the Second Workshop on Computational Approaches to Deception Detection. 7–17.
- [16] Ruchansky, N.; Seo, S.; and Liu, Y. (2017). CSI: a hybrid deep model for fake news detection. Proceedings of the 2017 ACM Conference on Information and Knowledge Management, CIKM. Singapore, 797–806.
- [17] Singhanian, S.; Fernandez, N.; and Rao, S. (2017). 3han: A deep neural network for fake news detection. Proceedings of International Conference on Neural Information Processing. Springer, Cham, 572–581.
- [18] Liu, Y.; and Wu, Y. B. (2018). Early detection of fake news on social media through propagation path classification with recurrent and convolutional networks. Proceedings of the Thirty Second AAAI Conference on Artificial Intelligence. New Orleans, USA.
- [19] Wu, L.; Liu, H. (2018). Tracing fake news footprints: characterizing social media messages by how they propagate. Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining, WSDM. Marina Del Rey, USA, 637–645.
- [20] Vo, N.; and Lee, K. (2018). The rise of guardians: fact checking url recommendation to combat fake news. Proceedings of the 41st International ACM SIGIR Conference on Research & Development in Information Retrieval. Ann Arbor, USA, 275–284.
- [21] Karimi, H.; Roy, P.; Saba-Sadiya, S.; and Tang, J. (2018). Multi-source multi-class fake news detection. Proceedings of the 27th International Conference on Computational Linguistics, COLING. SantaFe, USA, 1546–1557.
- [22] Aghakhani, H.; Machiry, A.; Nilizadeh, S.; Kruegel, C; and Vigna, G. (2018). Detecting deceptive reviews using generative adversarial networks. Proceedings of 2018 IEEE Security and Privacy Workshops (SPW). 89-95.
- [23] Goldani, M. H.; Momtazi, S.; and Safabakhsh, R. (2020). Detecting Fake News with Capsule Neural Networks. arXiv preprint arXiv:2002.01030.
- [24] Wang, Y.; Ma, F.; Jin, Z.; Yuan, Y.; Xun, G.; Jha, K.; Su, L.; and Gao, J. (2018). Eann: Event Adversarial Neural Networks for Multi-Modal Fake News Detection. Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. 849–857.
- [25] Mansour, A. M. (2018). Texture Classification using Naïve Bayes Classifier. IJCSNS International Journal of Computer Science and Network Security.
- [26] Wang, W. Y. (2017). ” Liar, Liar Pants on Fire”: A New Benchmark Dataset for Fake News Detection. Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics. 422-426.

# Cloud-Based Outsourcing Framework for Efficient IT Project Management Practices

Mesfin Alemu<sup>1</sup>, Abel Adane<sup>2</sup>, Bhupesh Kumar Singh<sup>3</sup>  
Faculty of Computing and Software Engineering  
Arba Minch University

Durga Prasad Sharma<sup>4</sup>  
AMUIT  
MOSHEFDRE under UNDP

**Abstract**—The optimum utilization of human resources is one of the crucial exercises in IT organizations. To provide a well-organized and cohesive working environment, organizations need to review their work culture in reference to newly evolved tools and techniques. To reduce the development cost of the IT projects and the optimum utilization of human resources, organizations need to review and redesign the project development processes. The significant challenges faced by IT organizations are the rapid switch-over (attrition) of IT professionals, physical migration or deployment, and redeployment of the human resources. This research paper is an effort towards the multilateral exploration of the techniques to adapt and improve the ICT enabled project management practices in an outsourced environment. This research is an effort with special reference to developing countries such as Ethiopia, where an acute shortage of high skilled IT human resources and their physical migration from one project location to another project location is a costly and challenging task. Ethiopia as a developing country and its IT industry is challenged by several issues like the capacity of ICT infrastructure and the skilled human resources. In such situations, IT projects are either challenged, impaired, or completed failed due to lack of IT human resources with desired skills and ultramodern up to date IT infrastructure. In this research paper, cloud computing technology is assumed as a key to the solution. For this, a systematic and careful investigation using mixed data analysis approach was used to adopt cloud-based outsourcing in IT project management practices i.e. design, development, and testing over outsourced systems by outsourced IT human resources. The major findings of this paper are to investigate and analyze how these cloud-based resources can be explored without physical movement or migration. For the novel improvement in the existing IT project management practices, the salient stakeholders' views were collected and analyzed for designing cloud-based outsourcing IT project management framework for the Ethiopian IT industry. The framework was functionally tested over the cloud-based Bitrix24 platform.

**Keywords**—Outsourcing; project management; cloud; IT industry; framework

## I. INTRODUCTION

Cloud computing is established as one of the computing technology which can provide the IT resources as and when needed and supports the real-time availability, scalability, and reliability using pay per use model. It is a model that enables convenient, and on-demand access to the IT resources over a wide area network. These resources are auto-configurable pooled computing resources such as networks, servers, storage,

applications, and services that can be rapidly provisioned and released with minimal management efforts or interaction with the service providers. It involves shifting the costs from capital expenditures (i.e. buying and installing servers, storage, networking, and related infrastructure) to an operating expense model, where one can pay for the usage of these resources. As a general standard a project can be defined as a temporary endeavor designed to produce a unique product, service or result with a defined beginning and end undertaken to meet unique goals and objectives, typically to bring about beneficial change or added value. The temporary nature of the projects stands in contrast with repetitive, permanent, or semi-permanent functional activities to produce products or services. In the context of computing or IT, project management is the application of expertise skills, IT infrastructure tools, and advanced techniques to a set of interconnected activities to meet the requirements of the IT projects [1] [2]. In general, the main goal of any project management effort is to manage the resources assimilated in such a way that the project is completed on time, within budgeted cost, and according to the desired functionalities or scope with promised quality expectations of the sponsor.

IT Projects have a terrible track record of their success and failure in the past couple of decades [3]. The 1995 Standish Group study (CHAOS) found that only 16.2% of projects were successful in meeting only scope, time, and cost goals, and over 31% of projects were canceled before completion [4]. A Price water house coopers study revealed that half of all the projects were failed and the success was only 2.5% where they met their targets for scope, time, and cost goals.

The IT industry of Ethiopia has been facing several challenges such as acute shortage of skilled human resources with the latest technology skill sets, modernized ICT infrastructure, platforms, tools, and techniques [5] [6].

How to ensure the optimum utilization of IT resources via emerging Medias such as the cloud is only an assumption? IT industry organizations of the developing countries are still lagging behind in sharing, deploying, and redeploying the IT human resources without physical migration from one physical location to another. The principal question in the mind is; how to explore and ensure the optimum utilization of IT human resources and infrastructure through a cloud based outsourcing techniques in project design, development, test, and management? How to investigate and explore the possibility of the cost reduction, time optimization, skills or expertise outsourcing, flexible outsource partner selection with frequent

provisioning and re-provisioning towards enhancement of the quality and the success of the projects? Ethiopia is a developing country located in the horn of Africa where most of the IT industries are focusing only on localized business and related IT projects and partners. These projects have salient limitations like quality, timeliness, scalability, robustness, flexibility, and high security of the project code and solutions with current practices of IT project management. Based on the thorough observation and preliminary study, the following research questions are set for the aforementioned problems:

- What are the basic issues and challenges that often make IT projects fail in the IT industry of Ethiopia?
- Which emerging tools and technologies can be explored to ensure the success of IT project management in the Ethiopian IT industry?
- Which Cloud-enabled IT project management outsourcing practices the framework can be a key instrumental to alleviate such issues and challenges?

To answer the aforementioned questions and design a solution framework; this research paper proposed to investigate and analyze the issues and challenges that lead to fail the IT projects before completion and design a cloud-enabled IT project management outsourcing practices framework for Ethiopian IT industries.

The following specific objectives of the study were also formulated with intermediate activities to achieve the goal of the research.

- 1) To investigate and analyze the issues and challenges affecting the success of the IT Projects in the IT industry of Ethiopia.
- 2) To identify, and explore the possible applications of cloud-enabled outsourced IT project management tools and technologies for improving success and the productivity of the IT projects.
- 3) To design and develop a contextualized Framework for Outsourcing in IT project management over cloud platforms.
- 4) To evaluate the Performance and Productivity of the Framework towards optimum utilization of IT resources in outsourced environments.

#### A. Scope of the Research Study

The main boundary of this study was to delimit the discovery of cloud-enabled tools, techniques, and their applications in outsourcing practices of IT projects in Ethiopian IT industry. The final contribution proposed was to design, develop, and demonstrate a framework for the IT project management related outsourcing practices overcloud. The major management practices considered were resource allocation and re-allocation, monitoring, control, deployment, re-deployment, follow-up meeting, instant reminders, rewards, computing, communication, collaboration, and ensuring the optimum utilization of resources over cloud-enabled platforms. The study covers only Ethiopian IT industries and IT-related projects.

#### B. Significance of the Study

The proposed research is significantly important for effective resource utilization, success rate enhancement, quality improvement, and cost reduction for the IT industry of Ethiopia. It explores the possible usage of cloud-enabled tools and techniques in computing, communication, and collaboration practices and opportunities worldwide for IT professionals. The IT professionals can share their knowledge, intelligence, and skills in the worldwide collaborative environment even without physical migration. The proposed cloud-enabled framework will be an essential instrumental to facilitate the outsourcing of almost everything in an inter or intra-organizational environment.

The main contribution of this research paper is the cloud-enabled framework which tried to advance the alleviation mechanism of the issues and challenges in traditional IT project management practices. The prototype developed, demonstrated, and evaluated with selected features evidently justified the improvements in cloud-enabled IT Project Management practices. The prototype demo and the user acceptance clearly validated the new knowledge contribution of cloud-based outsourcing of IT Project Management. Practices. Thus the framework promises to improve the salient features of project activities such as computing, communication, collaboration, monitoring, and control of human resources in anytime, anywhere over any-device with cost-effectiveness, reliability, scalability, optimum utilization, and all-time availability. The framework devised a new idea to establish a new pattern for IT Project managers where they can be on-site away from the site or in their offices away from the offices.

The paper is framed based on a scientific sequence of steps and pedagogy. The introduction of this research paper covers the basic background of the research domain, problem statement, research gap, research questions, objective, and contributions. In the review of the literature part, the selected concepts related to the problem domain are covered with some background reports of the world agencies such as the World Bank. The review of literature critically reviewed the selected papers to find out the research gaps, and to justify the research initiative with worth solving claim. In the research methodology part, basic ideas about research design type, research tools, and research method selection criteria are explained with parametric suitability assessment and analysis. In the data collection section the sample size, sampling technique, and the types of primary data collection methods and tools are explained. After the collection of data using three methods i.e. survey, interview, and technical observation, the framework is designed and explained with detailed functionality of each component. To demonstrate and validate the framework, a prototype is designed using the cloud-based Bitrix24 platform. The designed prototype is demonstrated before salient stakeholders and users to collect user acceptance and explained in detail using tables and charts. Finally, the summary of the research findings and contributions are covered in the conclusion part and the recommendation are also forwarded in the last section for future research directions.

## II. REVIEW OF LITERATURE

The rigorous review of literature was done to understand the gaps in the existing state of art researches towards addressing the issues and challenges in IT project management practices in general and developing countries such as Ethiopia as a specific case. The main focus of the review of literature was to understand and analyze the existing issues and challenges in the IT Project management practices in IT industry of the world vs. Ethiopia. The paper for the rigorous review were selected from the peer review journals, conference proceedings, and research project reports.

A study of Stephen Cacciola and Robert Gibbons [7] was conducted to investigate the outsourcing IT to improve the organizational performance. This study is an exploratory research which reveals how the cloud systems facilitate the services and improve the productivity that was envisioned but often not realized by organizations implementing old computing models, and allow for further productivity improvements in organizations that have benefited from previous technologies. In this research, the main focus was on the performance improvements parameters without any consideration of the project management practices in detail, specifically an exploration of offshoring or outsourcing in developing countries.

Research of Muhammad Younas Imran Ghani et al. [8] contributed the major benefits due to the amalgamation of agile software development methodology and cloud computing. This research tried to explore the efficient facilitation of global agile software development in the cloud environment. The researchers tried to explore the infrastructure features required for agile development in a distributed environment. This research is relevant to the proposed dimension of our research paper but the tools, methodologies, approach, and analysis didn't outline anything in the outsourcing of the agile development approach during project execution and management practices in IT industries or organizations.

Another important study of Mihret Abeselom Teklemariam and Ernest Mnkandla [9] was done for Software project risk management practice in Ethiopia. The major findings of this research study were focused on identifying uncertainties by project managers on risk management processes that whether they are carried out in the project implementation or not. This articulates a gap in the ability of project managers to adequately manage project activities. This research describes the insignificant relationship between risk management practices and project success. It suggests the presence of other factors that can play significant roles in the success or failure of projects but they are not definitely outsourcing related parameters. The proposed study tries to investigate the factors affecting the success of IT projects during the development phase. The study does not cover issues such as how outsourcing via the cloud can support in minimizing the failure possibilities in general and developing countries like Ethiopia as a special case where capacities of skilled human resources and ultramodern ICT infrastructure are typically challenged.

Research by Faith Shimba [10] was confined to investigate the successful adoption of cloud computing as a key to the realization of benefits promised by cloud computing

technologies. As organizations need the high processing capabilities, large storage capacity, IT resource scalability, and high availability, at the lowest possible cost. In this context, cloud computing becomes an attractive alternative media. The study explains that how an emphasis on collaboration between clients and vendors is essential for the successful adoption of cloud computing. This research is relevant to the proposed research because one of the dimensions of the proposed study is to adopt cloud computing in IT project management for outsourcing activities. This study motivates how outsourcing practices can be migrated over the cloud in the next generation of IT project management practices.

Jianfeng Wang, and Xiaofeng Chen [11] conducted a survey on Efficient and Secure Storage for Outsourced Data. This research study is the best benchmark for the proposed research. The study explored the support by providing the concept of outsourced efficient and secure infrastructure for the storage. This research focuses only on the efficient and effective use of cloud-based storage.

Gyöngyvér Husztáné Acsai [12] reveals the new knowledge about the qualitative inquiry of the project management of the Virtual Teams. This research study is focused on the advantages and challenges in the project management of virtual teams compared to the academic cyclorama. The study revealed four new advantages. The two challenges not yet identified and studied in the project management of virtual teams. Furthermore, a research gap at the cross-section of virtual project management and cloud computing is investigated studied in this study. The results of this study indicate that cloud computing tools are indispensable for virtual collaboration and benefits have been gained due to the adoption and usage of cloud computing tools in virtual project management. This research paper recommended that the future work of virtual project management and its connection with cloud computing can offer several choices for further research areas and our proposed study is focusing on that. The researcher recommended that widening this research, a case study about the perceived effects of a new, cloud-based project management software in virtual teams can give better outputs. And this is what our proposed study contributed a new idea for next-generation services towards better support from the global community of experts and services providers in the true globalization of the IT industrial revolution. The researcher focused only on qualitative analysis of advantages and challenges but didn't develop or propose a solution as a new knowledge contribution.

A paper of Sunil Patil and Y.S. Patil [13] conducted a review on outsourcing with a special reference to telecom operations. This paper was focused on outsourcing IT management and explores relevance to telecom operations. In the case of telecom operators, it is observed that the basic set of parameters influencing the decision of outsourcing is the same as the rest of the industry. It is observed that telecom operators have extended this model by outsourcing the management of network infrastructure, management of towers, billing systems, marketing, etc. This is creating new working models and relationships. This research is an effort that provides a direction for outsourcing concepts with a new dimensional thrust towards the exploration of the possibilities to share different

infrastructure now and then. And also helps with the trend in outsourcing is multi-sourcing, collaborative innovation needs to happen where all the vendors work together, innovate together with the client team and implement innovations. Innovations can be in different domains such as technology, processes, products and services, and forward-looking areas. This study didn't address the challenges faced by IT Project Management leaders during the development phases like scope creep, technology creeps, and scalability issues both in terms of infrastructure and human resources in the developing country industries.

Another research of Muhic et al. [14] reviewed was related to the next generation outsourcing overcloud. This case study exposed new knowledge that cloud sourcing reduces cost and complexity in the advantage of increased labor productivity. This study reveals some motivational benefits such as cost, complexity, and increase in productivity of human resources.

Stephan Schneider and Ali Sunyae [15] conducted research and discussed the determinant factors that are inherent to the particular sourcing option such as the risk of losing access to data and the benefits of increased scalability. Researchers have investigated a rich array of technical characteristics as determinant factors of IT Sourcing decisions, predominantly concerning the risks or benefits of the desired sourcing option. This study is an important review but focuses on contribution to the practice, as the determinant factors of cloud-sourcing decisions. This serves as a basis for practitioner-oriented guidelines and best practices regarding how to select and offer cloud services rather than discussing the management of the outsourcing service over the cloud environment.

A research study [16] discusses the Open Clouds for Research Environments consortium by putting in place an easy adoption route. It was estimated that numerous European research and education institutes will be able to directly consume these offerings via the European Open Science Cloud service catalog, through ready-to-use agreements. This research provides a new dimension of open source cloud usage. Other researches tried to encourage trust, security and transparency using different techniques such as MLP neural network and particle swarm optimization algorithm to detect intrusion and attacks. Such research is the efforts towards strengthening the security tire of the computing systems when we talk about the project management over virtualized cloud platforms [17] [18] [19] [20].

The rigorous analysis of several research studies focused and relevant to the proposed research, it was clearly identified and observed that "designing a framework for outsourced IT project management practices over the cloud" can be a new and innovative idea for new knowledge contribution to the domain.

### III. RESEARCH DESIGN AND METHODOLOGY

#### A. Research Design

The proposed research study is the mixed version of constructive and applied research design. The research paper used a mixed research approach i.e. qualitative & quantitative both for data collection and analysis. The detailed data collection methods and tools are illustrated in Fig. 1.

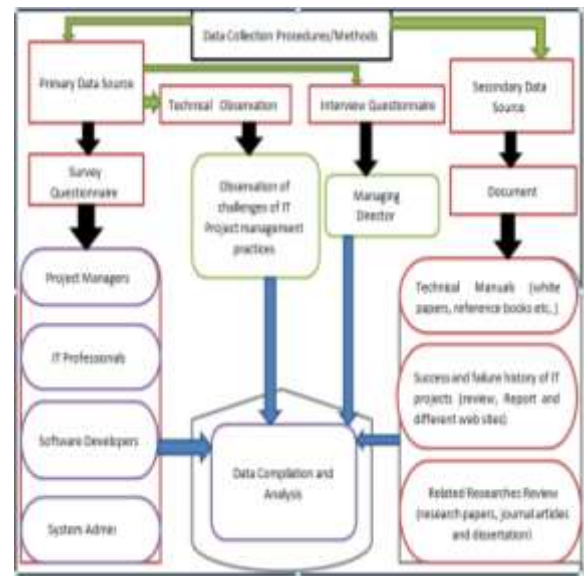


Fig. 1. Data Collection Procedures.

#### B. Data Sampling Strategy

The sample size for this research was seventy (70) and determined based on the mixed version of the online sampling tool i.e. Rao soft and the purposive sampling technique. The sample size was selected based on certain criteria set and the conditions. The open-ended Interview and Online Questionnaire-based Survey were considered for the detailed factual findings. Since the research study focuses on the scientific inputs from numerous stakeholders and observations of the researchers and therefore applied and constructive design strategy was followed. This strategy implied that the sample size seventy (70) is sufficiently representative for generalization of the results in the domain-specific user community.

#### C. Survey Questionnaire

The survey research questionnaires' were prepared and distributed to the numerous stakeholders such as project manager, IT professionals, software developer, system admin, and the end-users of the project management. This process was done for collecting the real facts about issues and challenges in existing IT Project Management practices in the IT industry/software companies. The responses of the respondents were collected by the researcher in a single folder. The collected data were processed using the Google data analysis tool for revealing the hidden insights.

#### D. Technical Interview

An Open-ended Interview questionnaire was designed for professionals/expert stakeholders. An open-ended interview questionnaire for IT Industry Professionals was distributed via a cloud-based Google Form to collect the detailed professional inputs from the professionals. Therefore, this interview questionnaire was considered to collect the general and managerial facts in detail for cross-validation of the input facts about the project management practices. In this process, the features and benefits of the cloud-enabled outsourced IT project management practices were compared with classical/traditional management practice and the research questions were framed accordingly.

### E. Technical Observation

In this section, a detailed technical observation was done by the researchers themselves. The technical observation was based on a checklist to collect and cross-validate the primary facts collected via survey and interview about the issues, challenges, features, performance, and other attributes of the existing state of art practices in the IT Project Management domain.

### F. Selection of Research Demo and Validation Tools

Cloud technologies have salient platforms/tools available in the IT market for computing, communication, and collaboration. As presented in Fig. 2, this research paper selected the tools based on the suitability assessment. These tools were used for framework designing, prototype development, and the functional demonstration of the research outcomes of the outsourced IT project management practice framework over the cloud. The outcomes of the research i.e. Framework was validated using two-fold methods i.e. 1) Functional demonstration, and 2) user/professional acceptance with the selected parameters as presented in the chart.

1) *The Bitrix24*: Bitrix24 is a free cloud service technology platform that provides over 30 handy tools, including online file storage and sharing, document management, real-time communications, and human resources management system. The best of all Bitrix24 is available as a self-hosted software platform for on-premise deployment that comes with API and open-source code. This can migrate from cloud to the user's server any time the user wants. Bitrix24 comes with free online workflow automation and business process management tools that can shoot the productivity of the users through the roof while eliminating the need to perform routine tasks manually. The system is industry independent and can be used to establish, standardize, and monitor processes and workflows in any IT department of the industry. After the overall suitability assessment, the Bitrix24 was found to be the most suitable tool and platform for the functional demo of the prototype of the outsourced IT Project Management Framework over the cloud.

2) *The Only-office*: The interface of the only office is divided into several modules: Documents, CRM, Projects, Mail, Community, Calendar, and Talk. The mail module combines a mail server for creating own-domain mailboxes and a mail aggregator for centralized management of multiple mailboxes. The calendar module allows planning and monitoring of personal and corporate events, task deadlines in Projects and CRM, sending and receiving invitations to events. A calendar can be integrated with the third-party calendars that support it. The community module offers corporate social network features: polls, corporate blogs and forums, news, orders, announcements, and messenger. According to this study, it is one of the supporting tools for project management practice over the cloud but it has limited features related to bitrix24. Because of this, the researcher preferred to use the bitrix24 for some demonstrations of the project management activities.

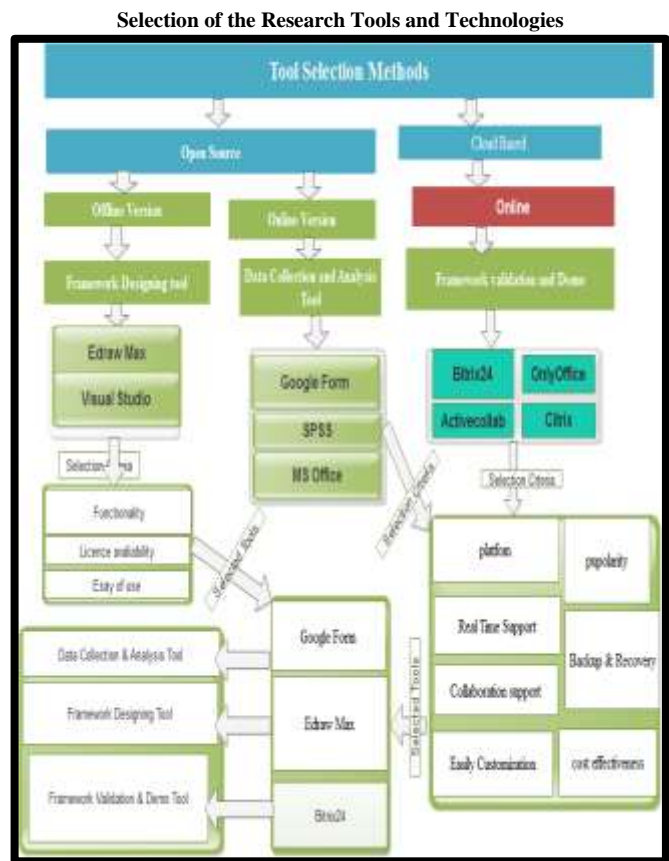


Fig. 2. Tool Selection Procedures.

Finally, parameter based suitability assessment of the different tools and technologies, the following tools and methods were selected for the different activities carried out in this research:

- Data Collection Methods: Survey Questionnaire, Interview Questionnaire, and Technical Observation.
- Data collection & Analysis: Google Form (GF) & SPSS but the research preferred to use the GF.
- Framework designing: Edraw Max Platform.
- Framework Demo & Validation: Bitrix24 Platform.

## IV. DATA ANALYSIS AND THE DISCUSSION OF RESULTS

### A. Primary Data Collection and Analysis

In this phase, responses were collected using a structured questionnaire for survey and interview from the domain-specific professionals and the stakeholders (i.e. Project manager, IT professionals, Network Admin, Software developers, and related professional knowledge holders). The collected facts were then analyzed and summarized for investigating the issues and challenges in the current status and the practices in the existing state of the art of IT project management. It was critically analyzed in comparison with outsourced IT project management over the cloud.



Most of the developing countries like Ethiopia and its organizations having a lack of IT professionals were critically reviewed and analyzed for the existing state of art practices. The organizations status were also analyzed for managing the services effectively and efficiently to encourage the intervention or the adoption of cloud computing in the IT industry, and software development. Also, human resource, time, and cost management towards the betterment of the business process, project control and effective communication were considered in the fact-finding and analysis process.

One of the anomalous issues observed during the fact-finding phase was the usage of emerging technologies for the computing and management of the resources. Also, it was observed that the cost management, and schedule management with follow up of the activities within the prescribed timeframe were found to be delayed and the projects were failed. The study used a purposive sampling technique for the survey, interview, and technical observation. The sample size selected was 70, for the survey and 6 for the professional's interview. The researcher also conducted a self-technical observation using a checklist. From the target sample of seventy, only forty respondents participated in the fact-finding phase and forwarded their responses. The fact findings data analysis were as follows.

1) What is the status of cloud technology adoption in Ethiopia?

The prime aim of this research was to investigate the current status of IT project management practices in the IT industry of Ethiopia. It was aimed to investigate and analyze the cost-effectiveness of IT projects, issues, and challenges in terms of time, quality, resource capacity, remote computing, communication, ease of anytime collaboration, discussion, agility, setting priorities, smartness in IT project management practices.

Also, it was envisioned to know how to adapt the newly evolved technologies like a cloud in project management practices to improve the above-mentioned features and create a new environment of outsourcing practices so as to alleviate the acute shortage of high skilled human resources in developing countries such as Ethiopia.

The responses of the respondents revealed that the adoption or intervention of such kinds of cloud-enabled IT Project Management practices is very low in the Ethiopian IT industry. As presented in Fig. 8, the 62.5% of responses indicate that the project management practices in the Ethiopian IT industry is still very low. As presented in Fig. 3, the 37.5% responses indicated that the cloud adoption or intervention in the IT industry project management practices is very low. The same facts were validated during the interview of the technical/managerial experts in the selected IT organizations and technical observation of the researcher. Maximum IT project managers accepted that they still not explored the possible usage of cloud-based project management practices. During the technical observation, it was verified that the managers are equipped with only managerial skills i.e. they are lagging behind in technology adoption tactics for becoming a techno-savvy professional. It was revealed that they still use traditional IT project management practices. These two

analytical facts indicate that there is either an acute shortage or lack of access to the cloud-based IT project management practices in Ethiopian IT organizations or IT project managers are not aware of such technologies. They need awareness of the smart adoption of the newly evolved technologies like a cloud for facilitating the exchange of information/data, advanced management, and real-time communication in IT project management activities.

Also, it was observed that the adoption of cloud-enabled IT project management practices can minimize the cost of the project, minimize the time of development, and improve the quality of the project outcomes. The analysis of the responses discovered that cloud-enabled technology if properly adopted and practiced; can improve the probability of the project success and mitigate the higher possibility of the project failures before completion in the IT industry of Ethiopia.

2) What types of remedial action/strategy have you adopted to overcome such challenges?

Currently IT-related technologies have created salient types of emerging management practices such as how to control and communicate in remote environments and information distribution within a short period of time in IT companies/industries.

As presented in Fig. 4, maximum i.e. 55.6% of respondents indicated that they are looking for an alternative mechanism or technology to resolve such aforementioned issues and challenges to optimizing the possibility of success and to minimize the failure rates. The 22.2% respondents responses indicated that projects are challenged or over-budgeted because of traditional project management practices, and only 11.1% revealed a critical question on success and responded that IT projects were canceled before completion and the reasons were unknown. The rest of the 11.1% responses indicated that they refuse to take projects because of the shortage of skilled human resource on the latest technologies and the unavailability of the infrastructure capacities in Ethiopia. As a matter of fact, IT project management requires focused and deep-rooted technical skills and knowledge in the areas of specializations.

This implies that; the Ethiopian IT industry needs a wide range adoption of cloud supported IT project management practices that are not yet adopted in Ethiopia. It can definitely help in alleviating such issues and challenges which are the main and root causes of the high rates of project failure or project rejection. This will make the IT project management practices more effective and efficient.

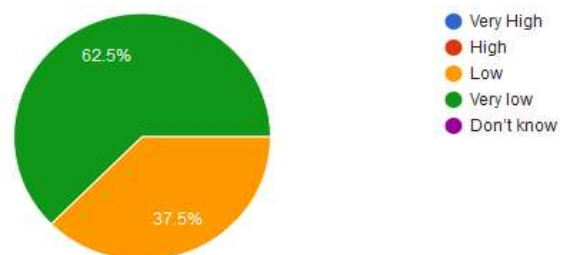


Fig. 3. Status of the cloud-based Technology Adoption in the Ethiopian Organization.

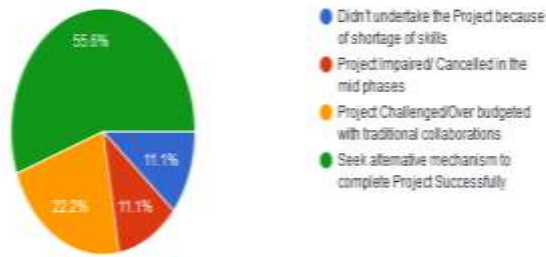


Fig. 4. What Types of Remedial Action/Strategy have you Adopted to Overcome Such Challenges?

3) Which technology from the following will be secured, robust, and most suitable for IT Project management?

Features such as suitability, robustness, and other managerial issues are very important to be assessed for the adoption of cloud in IT project management practices. As presented in Fig. 5 the maximum i.e. 66.7% responses of the respondents revealed that cloud-based outsourcing is better and suitably robust for improved IT project management activities. This indicates that the cloud can be one of the most suitable platforms for the assurance of a high success rate and to minimize failure possibilities. The 22.2% responses of the respondents revealed that traditional project management practices are better and easy as Ethiopia has poor internet connectivity. Only 11.1% responses of the respondents advised that cloud-based offshoring instead of outsourcing can be better of project management. The data analysis clearly justifies that there is an acute shortage of such practices and Ethiopia is lagging behind. Also, there is an acute shortage of research studies that can help in enquiring the suitability assessment of cloud-based outsourcing techniques in IT project management.

In the personal interview phase, when the same questions were asked to the IT managerial and technical staff for collecting the subjective inputs; their responses were somehow similar but they add few points like issues of the internet the speed which is the key to ensure the success of a cloud-based IT project management practices. The researcher also observed that technical and managerial staff are not well aware with such kinds of outsource practices over cloud but they are very much inspired and motivated to adopt.

4) Which is the most Challenging Factors that leads to fail the IT projects before completion in IT industries of Ethiopia?

To investigate the existing status of the IT resource management and paybacks of the cloud-enabled outsourcing in IT projects, the selected participants, professionals, and experts were asked to participate in the survey and detailed interview process along with technical observation of the researcher. As presented in Fig. 6 the maximum i.e. 50% respondent's responses revealed some hidden facts about the resources. It was revealed that the success of the IT projects is typically affected by the availability of the essential resources. As presented in Fig. 6, 25% of respondents were concerned about the budget as a major issue while the other 25% on the schedule of the project. Thus, the maximum i.e. 50% of respondents provided a clear picture and recommended that the cloud-based outsourcing of the resources (Human and IT), and

their management practices can significantly improve the IT project's success. The researcher's personal observation not only verified but strongly recommended that not only resources but the cloud-based resource sharing can significantly reduce the overall budget and time of the IT projects. This can lead to higher success in the project.

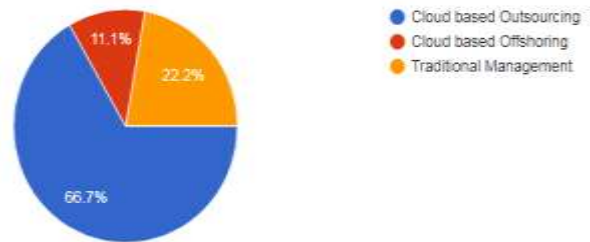


Fig. 5. Suitability, Robustness and Managerial Issues with the cloud Computing.

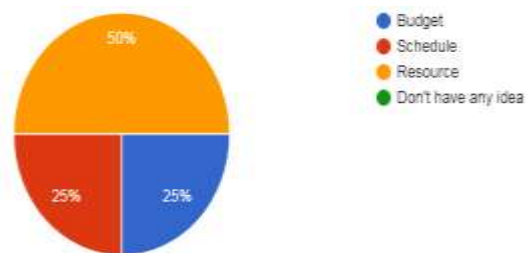


Fig. 6. The Resource Management in IT Project Management Practice Paybacks to Outsourcing.

## V. PROPOSED FRAMEWORK FOR OUTSOURCED IT PROJECT MANAGEMENT PRACTICES OVER CLOUD

In this phase, the gathered data and their analysis was used as an input to design the most viable framework for outsourced IT project management practices overcloud.

Moreover, this research answered the research questions: What are the basic issues and challenges that often make IT projects fail in the IT industry of Ethiopia? Which emerging tools and technologies can be explored to ensure the success of IT project management in the Ethiopian IT industry? Which Cloud-enabled IT project management outsourcing practices the framework can be a key instrumental to alleviate such issues and challenges?

To manage and control the resources during the project control and management process, different types of computing, communication, and collaboration technologies like email, telephone, postal services, traditional boards, and filing cabinets with paper-based heavy weight manual files are used. In today context, it was clearly revealed that these techniques are obsolete and there is a strong need to exploit and adopt advanced technologies such as the cloud in the developing countries like Ethiopia to improve the IT project management practices.

To migrate towards the advanced practices in IT project management, the researchers proposed a new contextualized framework- Durga Prasad Sharma-Mesfin Alemu-Abel Adane (DPS-MA-AA) for the effective utilization of the outsourced IT resources and support services over cloud. This framework

as presented in Fig. 7 was proposed to improve or replace the current state of art practices in IT project management. To check the validity of the localized contextual framework for the Ethiopian IT industry, we validated the framework by two-fold methods i.e. 1) Functional demonstration with limited features, and 2) stakeholder validation after the demo to check that how it can be a great instrumental towards the alleviation of the identified issues and challenges that adversely affect the success of the IT projects.

An outsourcing environment can support IT projects with a wide variety of services like 1) ease of access to domain-specific experts worldwide without their physical migration, 2) the entirety of the IT function, 3) support to easily defined functions of the project (designing, coding, testing, disaster recovery etc.), 4) network services, and 5) software component design, development, and testing. Different organizations realize to adopt the outsourcing for a number of reasons, and the most of them are based on the effective management and business profitability to the high rate of success of the projects.

## VI. DESIGN OF THE CLOUD BASED DPS-MA-AA FRAMEWORK FOR OUTSOURCED IT PROJECT MANAGEMENT PRACTICES

To explore and exploit the knowledge of the emerging technologies; the features of proposed next generation systems need to be studied and analyzed for their pros and cons.

This study designed a Cloud-enabled IT Project Management Outsourcing Practice framework named as DPS-MA-AA as presented in Fig. 7. This framework is a dynamic enabler to the IT project management practices and activities. This framework enables users to perform all the project related activities such as assigning tasks, monitoring, management, control, coordination, communication, sharing, deployment, redeployment, migration, computation, calendaring, and scheduling overcloud. These all the services can be supervised, managed, and controlled by the IT project manager. The cloud service providers (CSPs) can provide infrastructure and software platforms to transform this framework into reality.

To access the service over the cloud, clients need to sign-up first with their genuine and verifiable digital credentials overcloud. Afterward, the client can be invited by the organization (outside/inside) and follow the link of the cloud governance instructions. They can find the tasks or activities assigned to perform within a given time frame in relation to the agreement of the organization or assignee with the client.

In the case of IT organization; if users want to access the service or need to use the infrastructure of the cloud service provider (CSP); they should specify that what the users want to access it anytime, anywhere over any device.

The main aim of the proposed localized and contextualized Outsourcing IT Project Management Framework was to help and support the project management practices overcloud. These new ways of performing project-related tasks and rendering services are measured as the next generation task dynamics for the improvement of IT Project management services.

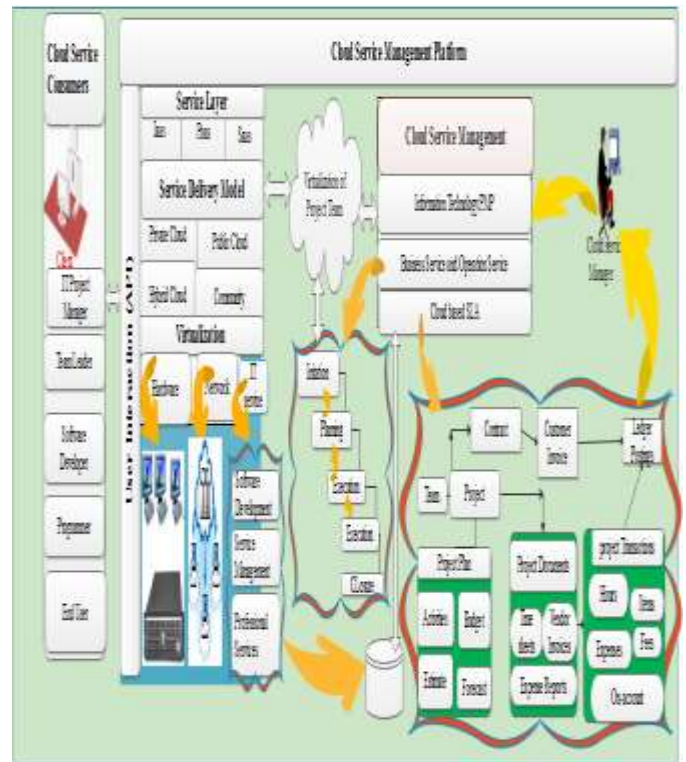


Fig. 7. The Proposed cloud-enabled Outsourcing IT Project Management Framework- DPS-MA-AA.

## VII. COMPONENTS OF THE FRAMEWORK

### A. Cloud Service Consumer Layer for Cloud Service Community

This is the Client Community Layer. In this layer, important players can be anyone from the Service Community in general but as a case study of Ethiopian industry, the service community consumers' are- IT Project Manager, Team Leader, Software Developer, Programmer, and the End Users. The service community consumers can access the cloud services provided by the cloud services providers (CSPs) in the selected business cases. Also, the IT organizations can get the services to advance or upgrade their IT capacity, internal work cultures by hiring external professionals from outside to perform the tasks that are out of the local capacity in terms of high-end IT skills. This implies that human resources can be outsourced (hired) using a pay-per-use model overcloud. In this manner, IT companies can fill the gap of "haves and have nots" in terms of IT human resource capacity. Here the product quality can also be improved or maintained.

### B. Cloud Service Management Platform Layer

This is the main Processing Layer of the Framework. In the IT project management, the resources, and the skills the gap with scope creep creates a typical challenge as investigated in a rigorous review of related works. These challenges or problems are the roots causes of the project failure, and/or impairments, and the same is verified by review of different research literature. To resolve such challenges or problems; this framework provides a new baseline solution and the direction about how to achieve the project goals even if such human resource challenges co-exists and the organization is lagging

behind in the capacity. The framework clearly shows the different pathways that how an IT project manager can resolve such issues and challenges in terms of scarcity of resources, the capacity of skills and the quality assurance of the project by applying the cloud-enabled technology-based tools and techniques to facilitate/ hire the human resource from outside without physical hiring or migration using the pay-per-use model. Thus the monitoring and control of the project cost, and time, through tracing and management systems over the cloud can support project success.

1) *Cloud service layer*: This layer is designed for the services in which the client requests are forwarded for the computing / communication/collaboration services required by cloud service providers (CSPs). These cloud-enabled services are readymade services. These services can be used/requested by any type of user to support and serve their professional activities, functions, and operations. The service models in this Layer may be anything like SaaS, IaaS, or PaaS and users can use it in terms of required software, infrastructure, and platform. These services are provided overcloud via the network with negotiable prices or a cost-sharing basis or free of the charge. To make such services easy and user friendly, the online training to the target users and employees are also facilitated via live media platforms and chat boxes. This Service Layer can help in connecting the users to SaaS (software as a service), PaaS (platform as a Service), and IaaS (Infrastructure as a services).

2) *Service delivery model*: This is the deployment layer in the framework. This layer offers its benefits through four types of service delivery/deployment models namely Private, Public, Hybrid or Community Models. A private cloud is built and managed within a single organization for the mission-critical secret services. Private clouds enable an organization to use cloud computing technology as a means of centralizing access to IT resources by different parts, locations, or departments of the organization. When a private cloud exists as a controlled environment, the problems described in the Risks and Challenges section do not tend to apply. A public cloud is a set of computing resources provided by third-party organizations. A hybrid cloud is a mix of computing resources provided by both private and public clouds. A community cloud shares computing resources across several organizations, and can be managed by either organizational IT resources or third-party CSPs. Similar to a public cloud except that its access is limited to a specific community of cloud consumers. The community cloud may be jointly owned by the community members or by a third-party cloud provider that provisions a public cloud with limited access.

3) *Virtualization layer (hardware, software & IT platform services)*: This is an interface layer of the framework in which all the resources are virtualized and made available to the users without physical movement. The clients or the users can access the resources like hardware, software, and supporting IT services to minimize the resource problem in IT project

management exercises, and scale-up or scale-down based on the capacity of the CSPs.

4) *Virtualization layer of project teams (Interface Layer)*: This layer is the interface/link layer between Service Delivery Model and the Cloud Service Management Layer. This layer facilitates the virtual teams to work remotely by deployment or redeployment and assigning or reassigning the responsibilities by project managers or project leaders. The team can co-operate with each other to perform collaborative tasks from anywhere, anytime over any device in a virtualized environment. This can alleviate and overcome the challenges faced in manual practices of IT project management.

5) *Cloud service management layer*: This layer of the framework consists of three major components.

a) *IT Project Management Practices*: These practices include Management, Control, Collaboration, Communication, Computations, Hiring, Discussion, Meeting, Coding, Testing, and Submitting. All of these activities are proposed to be done over the cloud platforms in an easy and convenient manner with low cost and high performance in anytime, anywhere over any device.

b) *Business Service Operations*: The cloud service operations for business where every system needs three activities i.e. 1) Management, 2) Control, and 3) Monitoring/Metering of the consumed services. Business service and operation component of the framework serves the process of project management in overall activities and practices in IT projects.

c) *Service Level Agreement (SLA)*: Service Level Agreement component of the framework is used for monitoring or metering of the consumed cloud services and helps in preparing the bills.

### C. Cloud Service Manager

This Component of the Framework is responsible for the activities like 1) IT Project Management Practices, 2) Business Service Operations, and 3) Service Level Agreement (SLA). Also, this component includes the services related to the Contract, Customer Invoice, and Ledger Postings. It decides about the billing/business operations of consumed cloud services with proper management, control, and monitoring/metering. In this layer, the project contains different financial and managing components like project plan, project document, project transactions, project team, project contract, project customer, vendor invoices, and ledger postings. The project plan can have activities, estimates, budgets, and forecasts. Project documents also have timesheets, vendor invoices, and expense reports of the projects. Project transactions refer to the hours of the activities, items, fees, and on-account. The contract is the agreement between the services user and the service provider. Customer and vendor invoices is an invoice may be created before or after the product or service is received. It's common for an invoice to be included with products being delivered, so the recipient can check off the items to make sure they are all their in-service stack. Ledger postings are the summary of all of the contract agreement of the organization and make the correction in each accounting system. Also, it is the process of business

transactions by recording information about the account. These all the activities are done in this layer. In this framework, the User Information APIs are used for integration and interaction among different components used by the IT project teams and members.

#### D. Validation through Functional Demonstration of the Framework over Bitrix24 Platform

According to the framework, the IT project management practices can be implemented, managed, and controlled overcloud. This phase of the research used a cloud-based Bitrix24 platform for designing a functional Prototype with limited features and functionalities. The Bitrix24 was selected to showcase the way how to manage the IT project management practices of the organization in an outsourced environment. It was based on the operations i.e. how to manage, and control the flow of information in the entire organization inside and outside both and share the project data/information within the timeframe. To explore and exercise the IT project management practices over the cloud, the end-user or client can access the tools and techniques through a well-structured registration process with a premium account. For the full professional transformation of the IT Project Management practices over the cloud; the designers can use the premium version with full permissions under SLAs. Thus this functional demo with prototype tried to evaluate the Framework based on the existing challenges in the current state of art IT project management practices in traditional/classical manners.

#### Access of Cloud-Based Bitrix24 Network

##### Rules:

- Register Before Sign into the Platform Network of Bitrix24.
- Login using the login provided by the Bitrix24.
- Create the private/public/hybrid network to make communication and collaboration overcloud.

1) *User/Manager:* As presented in Fig. 8, we can create a Bitrix24 network for private/public/hybrid cloud for communication. Now, select the option “My Bitrix24” above the left corner of the login page.

Then, Click on the link provided by the Bitrix24 platform network.

2) *Task management:* Fig. 9 presents how to control and monitor the project management practices over the cloud in real-time manner.

3) *How Create the Workgroup and Project Management Groups?*

In this activity as presented in Fig. 10 there are different options provided to the user to create the workgroup based on his/her choice to use the network. This is decided based on nature of the project for facilitating the entire communication and protection of sensitive information from outside access or loss of information based on the given workgroups.



Fig. 8. Login Provided form by the Bitrix24.

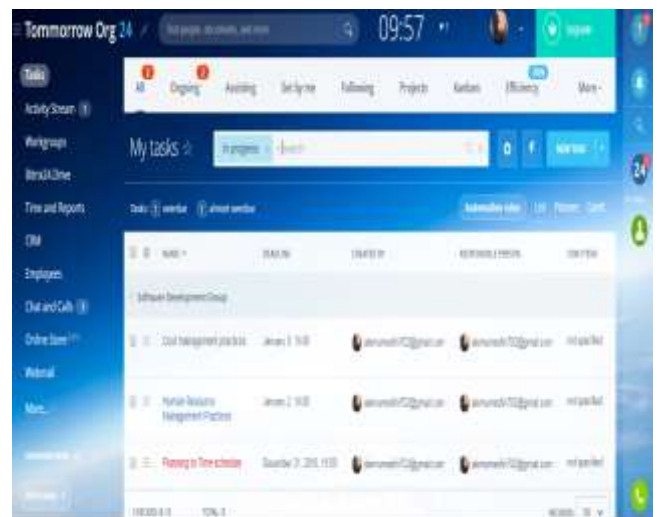


Fig. 9. Task Management Activity.

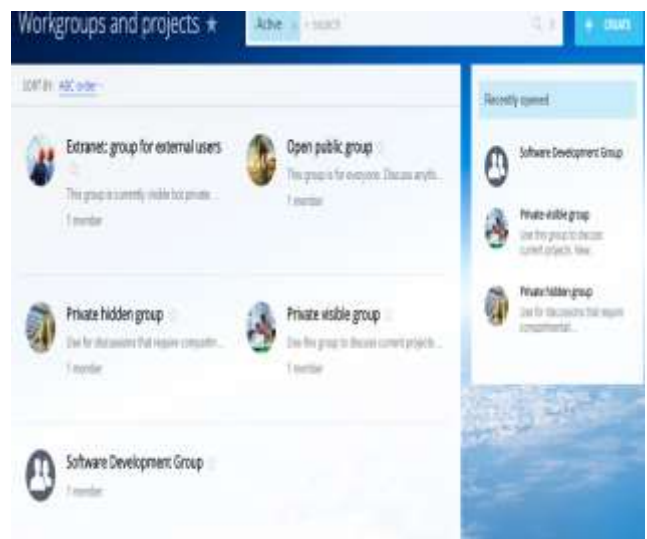


Fig. 10. Workgroups of the Project.

### VIII. THE VALIDATION OF THE DPS-MA-AA FRAMEWORK

To evaluate or validate the framework’s functionalities, this research paper selected only the effectiveness and efficiency of the framework in terms of computing, communication and collaboration in IT project management practices. The cloud-enabled a functional prototype of the IT project management framework clearly demonstrated the effectiveness and the efficiency improvement in the control and management of the activities, and resources used for computing, communication and collaboration purposes. This situation clearly justifies the economic improvement of the organizational activities. The factors evaluated were cost effectiveness, timeliness, resource utilization, accessibility, and availability of IT projects, information sharing along with effective collaboration.

1) *Cost effectiveness*: Based on the measurements of the resources and their cost like the hardware and software were compared with the traditional vs. cloud-based. In the case of traditional, the organization usually purchases all the required resources for their activities to perform the tasks and use them to do the required activities. But these resources consume the high capital of the organization and still limited in size i.e. scalability is limited, all the resources are centralized at a single location, fast obsoleting devices, high cost and underutilized. But in contrary to this traditional framework, the cloud-based resources are highly scalable, on-demand available, distributed over salient locations, no worry about obsoleting devices, rent based low cost and proper utilization with operational expenditure.

2) *Time complexity*: The time is the most important asset in the IT Projects and their Managements. The cloud-based IT project management practices are time efficient because they are managed in a virtualized environment i.e. without any physical movement of any hardware, software, or human resources. The virtual teams can be deployed, redeployed, transferred virtually, managed, controlled from anywhere at any time using any device.

3) *Evaluation of the resource utilization based on the measurement criteria*: The resource utilization overcloud was confirmed as better performing than traditional IT Project Management practices as is presented in Table I.

4) *Accessibility and availability of IT projects*: Since high uptime assurance under the SLA provisions is the basic promise of the CSPs. As presented in Table II, the accessibility and availability of the resources over the cloud also confirmed better than traditional IT Project Management practices.

5) *Information sharing and effective collaboration*: The design artifacts and features as presented in Table III of the cloud-based resources themselves justify the better opportunity in Information Sharing and effective collaboration than traditional IT Project Management practices.

6) *Security and privacy assurance under SLA*: The design artifacts and features of the cloud as presented in Table IV such as Security as a Service Models themselves justify the better opportunity of Information Security and Privacy than traditional IT Project Management practices.

TABLE I. EVALUATION OF THE RESOURCE UTILIZATION BASED ON THE MEASUREMENT CRITERIA

| SNO | Measuring Criteria    | Resource utilization |        |      |             |        |      |
|-----|-----------------------|----------------------|--------|------|-------------|--------|------|
|     |                       | Traditional          |        |      | Cloud based |        |      |
|     |                       | Low                  | medium | high | Low         | medium | high |
| 1   | Customer satisfaction | ✓                    | ✓      | ✓    | ✓           | ✓      | ✓    |
| 2   | Quality improvement   | ✓                    | ✓      | ✓    | ✓           | ✓      | ✓    |
| 3   | Product Improvement   | ✓                    | ✓      | ✓    | ✓           | ✓      | ✓    |

TABLE II. EVALUATION OF ACCESSIBILITY AND AVAILABILITY OF IT PROJECT BASED ON THE MEASURING CRITERIA

| S. No | Measuring criteria | Anywhere Accessibility            |               |                                   |               |
|-------|--------------------|-----------------------------------|---------------|-----------------------------------|---------------|
|       |                    | Traditional Project accessibility |               | Cloud based Project Accessibility |               |
|       |                    | Availability                      | Accessibility | Availability                      | Accessibility |
| 1     | Inside             | Yes                               | Yes           | Yes                               | Yes           |
| 2     | Outside            | No                                | No            | Yes                               | Yes           |
| 3     | Replica            | No                                | No            | Yes                               | Yes           |

TABLE III. EVALUATION OF INFORMATION SHARING AND EFFECTIVE COLLABORATION BASED ON MEASURING CRITERIA

| S. No | Measuring Criteria             | Information sharing and effective collaboration |                        |
|-------|--------------------------------|---|------------------------|
|       |                                | Traditional                                     | Cloud based            |
|       |                                | Effectiveness                                   | Effectiveness          |
| 1     | Two-way collaboration          | Less effective- Low                             | Highly effective –High |
| 2     | Platform support collaboration | Less effective-Low                              | Highly effective –High |
| 3     | Effective sharing of data      | Less effective-Low                              | Highly effective –High |

TABLE IV. SECURITY AND PRIVACY ASSURANCE UNDER SLA

| S. No | Measuring Criteria/Requirements                     | Security and Privacy |         |                    |   |
|-------|---|----------------------|---------|--------------------|---|
|       |   | Traditional          |         | Cloud based        |   |
|       |   | Security             | Privacy | Security           | Privacy   |
| 1     | Platform form                                       | Relatively good      | No      | Promised under SLA | High based on the user                          |
| 2     | Infrastructure                                      | Relatively good      | No      | Promised Under SLA | High based on the management perspective        |
| 3     | Inside and outside user access security and privacy | Relatively good      | No      | Promised Under SLA | High based on the nature of communication model |
| 4     | Application security                                | Relatively good      | No      | Promised Under SLA | High based on the nature of apps.               |

### IX. CONCLUSION

The research study concludes that the current status of IT Project management practices in developing countries such as Ethiopia is quite obsolete. Research clearly investigated, observed, and analyzed the critical gaps between the traditional IT project management practices and the modernized outsourcing of IT Project Management practices. The survey, interview, and technical observation clearly revealed that the current state of art IT Project Management practices are still not aligned with the latest IT tools, techniques, and practices. It was also discovered that IT companies are facing an acute shortage of IT human resources with high-end desired skill sets and up-to-date hardware and software resources with legal licenses. It was observed that there is an urgent need for critical review and redesign of the IT Project Management Practices in developing countries like Ethiopia. The judicious intervention or adoption of the modernized tools and techniques like the cloud is the need of the IT industry. The newly proposed framework DPS-MA-AA forwarded the alleviation mechanism of the issues and challenges in traditional IT project management practices. The prototype with selected features evidently demonstrated the improvements in IT Project Management practices. The prototype of the framework and the user acceptance clearly cross-verified and validated that the cloud-based IT Project Management Practices are the better and instrumental option for improving the efficiency and effectiveness of the IT Project Management activities like computing, communication, collaboration, monitoring, control, access, usage of remotely available hardware, software and human resources, etc. in anytime, anywhere over any-device with cost-effectiveness, high reliability, high scalability, and all-time availability. The framework coined a new idea to establish a new paradigm for IT Project managers where they will be on-site away from the site or maybe in an office away from the office. A new community of mobile IT project managers can be created where they can be accessible at anytime, anywhere over any device.

### X. RECOMMENDATIONS

There are many opportunities to perform additional depth and breadth research building on these findings. This study provides a starting point for the development of a more comprehensive conceptual framework of the capabilities that lead to outsourcing effectiveness and efficiency.

This study also provides a baseline for initiating further quantitative, survey-based research to get support for the findings from cloud-based project management software implemented with virtual teams.

Also, the complex cloud computing tools with premium privileges in the IT project management can be a better option to test the performance of the framework in the future. Finally, the research paper recommends adopting the framework in real-world environments to evaluate the performance.

### REFERENCES

- [1] M. Harwardt, "Criteria of Successful IT Projects from Management's Perspective," *Open Journal of Information Systems (OJIS)*, vol. 3, no. 1, p. 26, 2016.
- [2] D. P. Sharma, "Convergence of Intranetware in Project Management for Effective Enterprise Management," *Journal of Global Information Technology (JGIT)-USA*, vol. 4, no. 2, pp. 65-85, 2008.
- [3] Attarzadeh, "Project Management Practices: The Criteria for Success or Failure," *Communications of the IBIMA*, vol. 1, p. 8, 2008.
- [4] T. Clancy, "The Chaos Report," Standish Group International, 2018.
- [5] L. S. Njuguna Ndung'u, "The Fourth Industrial Revolution," *Capturing The Fourth*, Washington, D.C., 2000.
- [6] P. D. Authorized, "THE Change Nature Of Work," *World Development Report*, Washington DC, 2019.
- [7] S. C. a. R. Gibbons, "How and when can outsourcing IT improve organizational," 2012.
- [8] D. N. A. J, M. M. K. Muhammad Younas Imran Ghani, "A Framework for Agile Development in Cloud Computing Environment," 2016.
- [9] M. T. a. E. Mnkandla, "Software project risk management practice in Ethiopia, (2017)," 2017.
- [10] F. Shimba, "Cloud Computing: Strategies for Cloud Computing Adoption, (2010)," 2010.
- [11] J. Wang1, "Efficient and Secure Storage for Outsourced Data," *SPRINGER*, vol. 1, no. 3, p. 178-188, 2016.
- [12] G. H. ACSAI, "Project Management Of Virtual Teams: A Qualitative Inquiry (May 2016)," 2016.
- [13] Y. P. Sunil Patil, "A review on outsourcing with a special reference to telecom operations".
- [14] M. J. B. Muhic, "Cloud Sourcing- Next Generation Outsourcing," Elsevier.
- [15] A. S. Stephan Schneider, "Determinant factors of cloud-sourcing".
- [16] J. M. D. Blasco, "The H2020 OCRE Project Opens the Gates of the Commercial Cloud and EO Services Usage to the Research Community," *Emerging Science Journal*, vol. 4, no. 2, pp. 89-103, 2020.
- [17] A. S. Saljoughi, "Attacks and Intrusion Detection in Cloud Computing Using Neural Networks and Particle Swarm Optimization Algorithms," *Emerging Science Journal*, vol. 1, no. 4, pp. 179-191, 2017.
- [18] DP Sharma, N. Shekhawat, "Cloud Computing Security through Cryptography for Banking Sector," in *INDIACom-2011*, INDIA, 2011.
- [19] DP Sharma, H. S. Shekhawat, "Hybrid cloud computing in e-governance: Related security risks and solutions," *Research Journal of Information Technology*, vol. 4, no. 1, pp. 1-6, 2012.
- [20] DP Sharma, Bright Keshwani, Dharmveer Yadav, "Study of Intranet over Cloud," *International Journal of Innovation in Engineering and Technology*, vol. 7, no. 2, vol. 7, no. 2, 2017.

# A Clustering Hybrid Algorithm for Smart Datasets using Machine Learning

Dar Masroof Amin<sup>1</sup>

Research Scholar  
MMICT and BM Maharishi Markandeshwar  
(Deemed to be University)  
Mullana, Haryana 133203, India

Dr. Munishwar Rai<sup>2</sup>

Professor  
MMICT and BM Maharishi Markandeshwar  
(Deemed to be University)  
Mullana, Haryana 133203, India

**Abstract**—In the field of data science, Machine Learning is treated as sub-field which primarily deals with designing of algorithms which have ability to learn from previous information and make future predictions accordingly. In traditional computational world the Machine Learning was generally performed on highly performance servers and machines. The implementation of these concepts on Big Data analytics algorithms has high potential and is still in its early stages. So far as machine learning is concerned, performance measure is an important parameter to evaluate the overall functionality of the algorithms. The data set is a different entity and the measuring of performance on a data which is unseen is also called as test set, and training set is a Data set which is training itself. The Data Mining is extensively using learning algorithms for data analysis and to formulate future predications based on archived data. The research presented provides a step forward to make smart data sets out of training data set by evaluating machine learning algorithm. The research presented a novel hybrid algorithm that attempts to incorporate the feature of similarities in Random Forest machine learning algorithm for improving the classification accuracy and efficiency of working.

**Keywords**---Random Forests (RF); Jaccard Similarity (JS); triangle; smart data; root mean square error; mean absolute error; machine learning

## I. INTRODUCTION

Big Data terminology is generally applied to the data that grows exponentially and which cannot be accessed by using conventional database systems. The size of data sets involved in big data cannot be handled by traditional software technology and database. The common tools, storage systems cannot store, process and manage the size of datasets [1]. The big data analytics is changing the overall life on the globe in various aspects, viz. health care, marketing, etc. [2]. The new technology and techniques are getting imbibed into out day to utility, Internet of Things (IoT) devices. The technologies being used generates valuable data which can in turn be used for making important decisions [3]. The data that is generated out of these devices can be either result of conscious intervention or unintentional. This involvement of the human in creating of data is at same time creating opportunities for analyzing the data for various purposes. The number of devices that were connected to internet and were generating the data was double in 2008 than general human beings. The expectation is that it will reach up to 50 billion by the end of 2020 and hence the creation of data can be seen as

exponentially growing [4]. In a similar manner by 2025 the economy is like to grow \$11.1 trillion a year [5]. This is the reason that multinational corporations' are moving towards big data technologies to improve their skills for making additional profit out of their investments [6]. The increase in general technology is also playing a vital role in meeting the demands of growing data [7]. The solution provided by the IoT is by combining the information technology with hardware and software. The Big Data analytics generally provides the soft solutions to handle the exponential growth of data. The physical functionality of digital devices is then accessed locally and globally [8]. The creation of human sense while providing soft solutions to the problems arising out of growth of data, the cost of security system and other functions can be minimized at optimal functionality. The application of distributed system in data analytics will minimize the load over the system. The data generated from IoT should be explored for using it in a formal process. As the traditional analysis of data was providing reports or models on the basis of data available in the system similarly the Big Data coming out of IoT is providing analysis of data generated on real time basis. The system should handle this real time input properly and should provide an optimal solution to the user of the data so far as decisions are concerned. This analytics enables smart decision making and means of quantification and goal tracking. And the traditional modulation provides analysis of static data analysis of unstructured data [9]. The case of Big Data is complex where large data is involved and which needs finding correlations between various types of input in real time. In outmoded analysis of data, the archival data is used for extracting and establishing relationship among various variables. Machine learning instead begins with the result of variables involved and thereon uses the interaction of the predictor variables. The Google's machine learning application is reducing the use of energy and making cool environment for their data center. By adopting machine learning technology, Google will save millions of dollars in creating a favorable environment for their servers. The machine learning algorithms have capability to learn various physical environmental changes and looking hidden patterns in data and later on making smart decisions to tackle any odd situation. The use of the technology have improved the services of smart homes, healthcare, agriculture etc. In the upcoming years billions of machines and devices will be connected to the internet and therefore data generated will be huge. The gigantic growth in the data have to be tackled by



smart machine learning algorithms in order to reap the overall benefit of this growing data. The prime objectives of this proposed research is to provide a hybrid algorithm for applying on smart data. And feasibility of running machine learning algorithms on Big Data Frameworks and optimization of algorithms for big data [9].

The prime focus of the proposed research is to provide insight how data sets are handled by learning algorithms. The quantification of results and learning pattern of various algorithms provides a source to make to future decisions. In addition, the research provides understanding of machine learning algorithms for comprehension of smart data sets generated by IoT. At the very first instance the smart data is generated by IoT devices through inbuilt applications. This smart data generate with specific domain can be fetch to a machine learning model for resolving the issues arising out of the growing use of devices and data. In case of real time applications, there should be consideration of response time and reliability. This can be achieved only if a learning algorithm would be used properly with prior compatibility with the speed of data creation. The accuracy level should also be a prime criteria for handling the data. Because the accuracy can only provide us better results after stage of data execution. The machine learning algorithms reveal may insights regarding the data characteristics. To explore the more insight into the smart data, the data patterns must be looked into detail. This pattern finding and extraction of data will help in enhanced accuracy score, event will be responded on real time basis and it will consequently affect the decision making. The Random Forest has less training time and multiple trees minimizes risk of overfitting. Moreover, this machine learning algorithm performs better on big data, for data sets with large size, highly accurate predictions are produced. In addition Random Forests can maintain accuracy when a large portion of data is missing.

Random Forest or Random Decision Forest is method that operates by constructing multiple Decision trees during training phase. The Decision of the majority of the trees is chosen by the random forest as final solution.

The other details of this proposed research follows as. Section 2 discusses the related work in the area of handling data sets. The section 3 provides the basic concept of smart data technologies. The proposed hybrid algorithm is discussed in Section 4. The results and observation from experimental work can be viewed in Section 5. The research is concluded with future scope in Section 6.

## II. LITERATURE REVIEW

Oscar D. Lara et al. surveyed on wearable sensors for the human activity recognition in the state-of-the-art domain. The parameters used where learning scheme and response time for introducing the organization of human activity recognition systems in two-level-taxonomy format. The scheme qualitatively compared 28 systems with regard to parameters viz general design issues, response time, flexibility, obtrusiveness, recognition accuracy and other issues. The important parts like machine learning and extraction is included as components of human activity recognition

systems. The authors have provided further directions to explore in more pervasive and realistic scenarios [10].

The researchers have used smart phones to present an efficient way to classify the activity of humans on daily basis. The basic design has been simulated by considering the fixed-point arithmetic of Support Vector Machines methodology. The research has using principles of Structural Risk Minimization in which complex approaches are neglected instead simpler techniques are used which provides equivalent attributes to learn. The use of the proposed research is to update the ambient applications using current technology such as in smart and remote patient monitoring environments. The properties like minimal use of resources and real time processing which saves the energy overhead for maintaining recognition is providing advantage over traditional approaches [11].

A novel approach has been proposed by the authors based on sensor worn on body [12]. A two-phase algorithm with abnormality detection has been proposed for looking into abnormal activities when there is scarcity of training data. SVM one-class in phase first is built of normal activities for filtering the normal instances of the activities. The adapted abnormal activity models are used to track suspicious traces using KNLR. The researchers claimed that proposed approach provides a better results and tradeoff between false alarm rate and detection rate. The effectiveness of the approach is demonstrated for real data obtained from human body using sensors. The authors have also described drawbacks of their work that if abnormal activities will become normal, then there is risk in general abnormal models. The type of situation is arise when a user repeats the things alternatively after fixed instances of time [12].

A detailed aspects of human activity recognition for offline and real time processing has been presented by Bishoysefen et al. [1].The researcher explained best features of classification algorithm for achieving realtime optimization for recognition accuracy and computational complexity. The data from more than 10 sources has been collected based on day to day activities and exercises. The result of analysis showed that machine learning algorithms have better performance with regard to the efficiency and accuracy.

The big data is amazing source of supply of useful knowledge and information for different end-users and varied systems. In order to handle such kind of inflow of information, the automation is a reliable source and that can be achieved using processing through machine learning. The information and communication technology is serving in various analysis sectors by providing specified tools and platforms for making professionals enable get valuable predictions. These ICT based techniques are developed by prominent firms viz. IBM, Microsoft, Google, etc. The research published provides concepts of Machine learning algorithms in Big Data Analytics [13].

The authors researched the work done on Internet of Things using the big data mining concept. The researchers provided three tiers for stage wise analysis of data in future [14].

The researchers [15] designed a random forest algorithm with improved classification for multiple classes related to a disease. The algorithm provides better classification of individual variable. The improved method has increased accuracy and general process of work. The percentage increase for the classification accuracy was achieved upto 97.80% for multi-class dataset.

The data missing in data sets is creating large number of classification errors. There is a technique called imputation technique that helps to complete data which is having missing datasets. The researcher [16] developed an approach for incorporation of feature selection of genetic-based method and imputation for enhancement of classification of missing or incomplete data.

The internet among various objects is called (IoT) using sophisticated and complex communication technology without human intervention. The researchers put forth a big data based analytical healthcare system using Random Forest algorithm. The methodology proposed shows better accuracy for classification than traditional logistic regression and Gaussain method [17].

The research presented in [18] focuses on smart network fault analysis prediction. The proposed research used a modified RF algorithm for providing enhancement in accurate analysis prediction. The methodology proposed improves accuracy of overall system.

The study provides details about how machine learning techniques can become base for smart data analysis for IoT. The deployed methodology provides high velocity data processing. In addition fast training and classification of datasets have been achieved [19].

The presence of noise during classification causes incorrect labeling in data. The situation becomes disastrous as it changes the basic variables and instances. The researchers [20] put forth an ensemble method which is iterative in nature for restricting noisy instances. The proposed methodology effectively contributes in transforming simple big data to smart one.

The authors [21] proposed a novel approach to deal with big data in which Random Forest algorithm and Support Vector Machine is used. The feasibility and robustness check is performed using parameters like confusion matrix, recall, precision, specificity and sensitivity. The result shows 95% accuracy with big data.

The researchers [22] proposed a big data framework with scalability which collects data from smart devices and stores that in NoSQL. The machine learning algorithms has been imbibed with framework for future predictions. Various machine learning algorithms have been used to load forecast domain specific environment.

The research in [23] focuses on making smart data out of large chunks of raw data. The researchers have used two big data libraries BigDaPspark and BigDaPFlink for extracting

smart data from big data. These libraries are built on Apache Spark and Apache Flink big data frameworks for cleaning, discretization and other things.

The authors in [24] have integrated smart persistence algorithm, past production data, irradiance along with Random Forest machine algorithm for enhancing capability of producing forecasts accurately. The methodology has shown incredible results.

The researchers have proposed a framework for extracting hidden patters with data stored in database. The research was carried out for looking into future possibility of treating vulnerable diseases like breast cancer, diabetes, heart diseases etc. Machine learning algorithms such as Random Forest and others have been used for prediction analysis. The results showed that Random Forests provide better and accurate results than others [25].

### III. SMART DATA TECHNOLOGIES

The growing demands of control and co-ordination among various sectors and programs has created need to make the novel systems to be smart enough to take care of even a minute detail. The general architecture of smart exigencies is shown in Fig .1. The system has to co-ordinate with various domain specific data stores and provide an optimal solution in the form of analytics reports [26].

Fig. 1 provides a detailed pictorial representation of the smart data technology program. The use of Big Data analytics and techniques to extract patterns and solutions for a domain specific system incorporated with issues related to security, social and economic aspect, legal aspect and many more. In order to deal with problems that depends on large data sets and technology uses follow some constraints that are used in all digital solutions provided. The program associated with smart data generally provides three issues described that should be incorporated in addition to achievement of smart analytics of big data through machine learning.

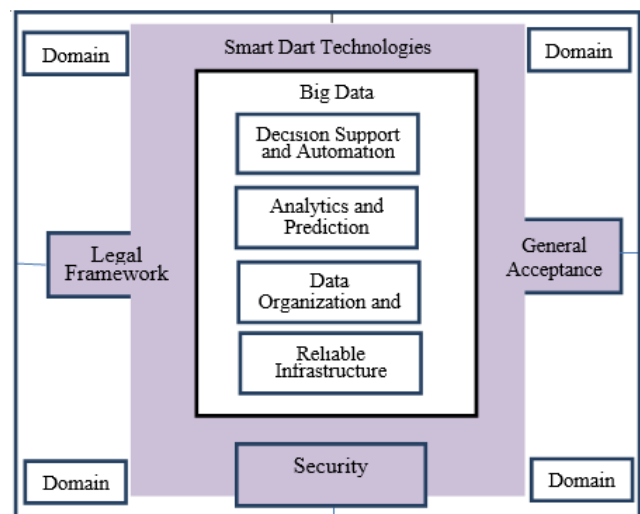


Fig. 1. Smart Data Technology Program [26].

#### IV. HYBRID ALGORITHM

The research proposed here is based on hybrid algorithm that is used to implement the machine learning at node level. Each node is an analytical unit it is related domain [27]. The node is a mode which provides a storage unit and the related functionality to analyses the data present in node. In current context Random Forest (RF) is an important concept to convert a simple data set to a smart one. This machine learning algorithm is flexible algorithm. The prime functionality of this algorithm is to provide accurate and efficient results, with using any tune of type hyper-parameter [25]. The RF is researcher's prime importance because of its simplicity and mainly to provide definite classification and regression. The Random Forests algorithm comes under the category of supervised learning algorithm. The algorithm creates a random forest decision trees which are trained with "bagging" methodology. The Random Forests provides multiple trees to classify on basis of attributes a new object [26]. Bootstrapping the data plus using the aggregate to make a decision is called "Bagging". The general flow of algorithm is shown in Fig. 2 reference node architecture. The data is collected from reliable source. The source can either be online source for current data or it can be archival data from any storage medium. The data is fetched on basis of some query to be supplied for analytics purpose. The algorithm proposed is integration of similarity formula with Random Forest Machine Learning algorithm.

The section describes the general working of the proposed algorithm that is based on machine learning concept random forests. The trees constructed using this ML algorithm works independently. A rectangular matrix is used to represent nodes of a tree which in turn represent the data stored in each domain node, and at each step of the construction the cells associated with leafs of the tree form a partition of matrix. The root of the tree corresponds to all of matrix.

At each step of the construction a leaf of the tree is selected for expansion. In each tree we partition the data set randomly into two parts, each of which plays a different role in the tree construction. We refer to points assigned to the different parts as structure and estimation points respectively. The shape of a given tree is influenced by allowing structure points. The internal node of the tree is determined by the split points and dimensions. The predictions made by the leafs of the tree are not by any way intervened by the structure points. There is a dual role played by the estimation points. There is no effect by estimation points on shape of partition of trees but the estimation points fit the values of estimators in each leaf. The assignment of points to estimations or structures with original equivalent probability has data that is randomly distributed among each trees. The partitions ensures the real time consistency in the trees formed out of data. There is no need to add parts for fitting each and every subset of data with the tree because making more samples lowers the performance of the system. The construction of tree is kind of parameterized which provides minimum estimations that should appear in a leaf. The size of the training set is taken into consideration while setting the parameters. The training set size also provides information about the minimum size of leaf [28].

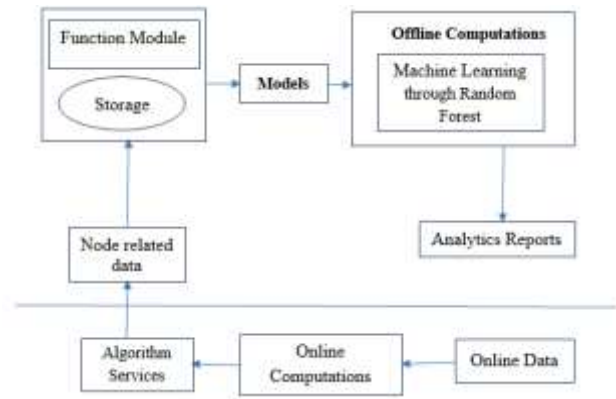


Fig. 2. Reference Node Architecture.

There is random selection of minimum of  $(1+\text{Poisson}(\lambda), D)$  on basis of distinct candidate dimensions when there is selection of leaf for expansion. The candidate split points are traversed for choosing split point in a leaf in case of each dimensions of candidate. In case of standard algorithm of random forest, the splits points are in a specified range and the candidate dimensions is projected into these points. At time of searching process, the range is declared by selecting specified points and process of search is done only over the defined set of points. However, the hybrid algorithm changed the traditional process by restricting the range to the specified set of structure points. And restring rage enables the trees to be in balanced state [28].

$$Error(L) = \frac{1}{N^s(L)} \sum_{I_j=S}^{Y_j=A} (Y_j - \bar{Y}^L)^2 \quad (1)$$

$$I(S) = Error(L) - Error(L') - Error(L'')$$

The L represents the leaf to split and L', L'' represents the children which splits L at S. The empirical mean is denoted by notation  $\bar{Y}^L$  for structure points falling in the cell L and the number of structure point counts in L is denoted by.

$N_s(L)$ . whether point  $(X_j, Y_j)$  is an estimation point or structure is denoted by indicator variables  $I_j \in \{e, s\}$ . When  $I(S)$  maximizes without creating any further children with less than  $k_n$  estimation points, the split is chosen as candidate and for non-existence of any such candidate stops the expansion [28].

$$f_n^i(x) = \frac{1}{N^e(A_n(x))} \sum_{I_j=e}^{Y_j \in A_n(x)} Y_j \quad (2)$$

The predictions of each tree are averaged by:

$$f_n^{(P)}(x) = \frac{1}{P} \sum_{j=1}^P f_n^j(x) \quad (3)$$

The various similarities has been implemented and an integrated method has been put forth for minimization of errors. The similarities has been defined through random mathematical functions. The Jaccard index has been used when the data is of discrete nature to check the errors [29]. The formulae specification are:

$$\text{Jaccard Similarity } (j_{i,j_p}) = \frac{|I_i \cap J_p|}{|I_i \cup J_p|} \quad (4)$$

Where

$$J_i = \{u \in U | r_{u,j} > 0\} \text{ and } J_p = \{u \in U | r_{u,p} > 0\}.$$

Triangle Similarity

$$(j_i, j_p) = 1 - \frac{\sqrt{\sum_{u \in C_{j_i, j_p}} (r_{u,i} - r_{u,p})^2}}{\sqrt{\sum_{u \in C_{j_i, j_p}} r_{u,i}^2} + \sqrt{\sum_{u \in C_{j_i, j_p}} r_{u,p}^2}} \quad (5)$$

[0,1] is the range value, and 0 indicates  $C_{j_i, j_p} = \emptyset$  [4]

### A. Algorithm Description

Algorithm( Nodes, N, S)

Nodes represent the set of nodes in a domain

N is number of Nodes with varied sets of clusters &

a) For each nodes N in Nodes

If the Key matches the node N

For each sub node S

Create a fuzzy random forest using formula(1-3) for different variations

For each decision tree in forest

If the input variable is incessant

For each split point

Create a partition using formula-5;

If there is definite variable

Calculate the similarity using formula-4

b) then choose the value with the optimal index;

c) accordingly provide child nodes based on the output produced using the indexes;

d) the calculate the gravity of association of each value with the next level nodes;

e) Repeat all the steps a-d until for every node in a tree

**end**

The general declarations of similarities in the coded form is written as follows:

```
public class GenerateCosine {
    public double cosval()
    {
        double cv=Math.random();
        return cv;
    }
    public double jaccardcal()
    {
        double jv=Math.random();
        return jv;
    }
}
```

The definitions are generally coded in a format shown below.

```
GenerateCosine gn = new GenerateCosine();
for(int i=0; i< similaritycount;i++)
{
    if(allcosine[i]==0)
    {
        allconsine[i]=(float)gn.cosval();
        System.out.println(similarityvalue[i][1]+""
+ "" +similarityvalue[i][2]+""
+ "" +allcosine[i]);
    }
}
```

## V. RESULTS DISCUSSION

The proposed algorithm gathers information from various data sources. The data files are distributed into the various domains through predefined algorithm [27]. The further process is done in internal nodes where the actual acquisition, management and mining is performed for analytics process. The new smart data set based algorithm proposed performs the optimal analytics of data by following the concepts of machine learning. The improved machine learning algorithm through incorporation of multiple similarity index is providing the environment for error free analytics in real time manner. The nodes present inside the domain are maintaining the security while loading data for mining purpose [27]. The proposed algorithm enhances the services and gives generalized machine learning Big Data techniques and various different protocols.

The algorithm provides the services in three phases. In phase first, the data mining process is done and the smart data clusters are fetched into the program for making the trees using enhanced Random Forest algorithm. Also in the same phase, the hybrid algorithm has been tested on twelve distinct datasets. The datasets has been obtained from Kaggle online data repository and is tabulated in Table I. The granularity based approach for data stream mining is good method so far as computational intelligence is concerned [30]. The datasets contains varied number of input attributes and instances. The varied folded validation has been performed with more than two trials with different basic attributes for the machine learning algorithm. In order to check the exact result the number of folds and validations has been kept same in each dataset in order to check the overall functionality. This phase generally provides the classification accuracy of proposed algorithm. Fig. 3 shows the classification accuracy comparison of the traditional and proposed algorithm and which clearly signifies the outstanding performance of the system in consideration. The red dots are showing the enhanced results so far as classification accuracy is concerned.

TABLE I. CLASSIFICATION EFFICIENCY OF HYBRID ALGORITHM

| Disease Type     | Instances | Attributes | Classification Efficiency | Classification Efficiency of Hybrid Algorithm |
|------------------|-----------|------------|---------------------------|---|
| COVID-19         | 457       | 10         | 94.44                     | 98.32   |
| EBOLA            | 145       | 8          | 87.06                     | 95.3  |
| MERS             | 121       | 5          | 78.32                     | 86.15   |
| H1N1             | 175       | 10         | 86.06                     | 96.06   |
| SARS             | 208       | 12         | 72.45                     | 89.35   |
| HIV/AIDS         | 193       | 8          | 92.36                     | 99.12   |
| H3N2 HonKong Flu | 113       | 10         | 84.16                     | 92.34   |
| H2N2 Asian Flu   | 241       | 6          | 89.02                     | 96.56   |
| Spanish Flu      | 243       | 5          | 78.12                     | 87.63   |
| ZIKA             | 54        | 6          | 79.11                     | 89.6  |

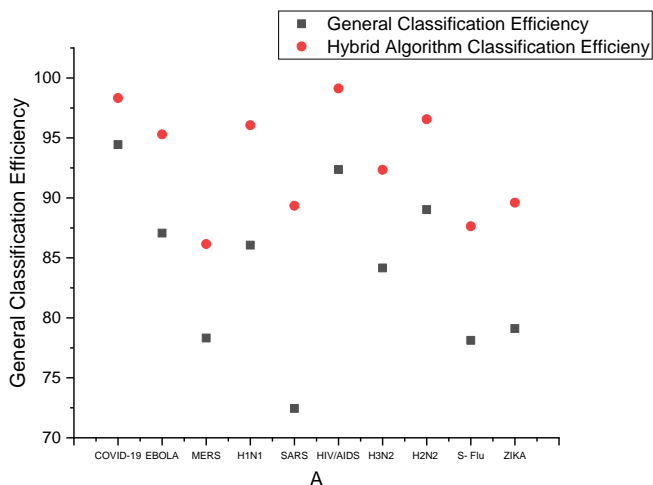


Fig. 3. Classification Accuracy of Hybrid Algorithm.

The phase-II provides verification details of error measures enhancement of proposed algorithm and details of various kinds of error measures. The error specification of results general machine learning algorithm and the proposed algorithm has been compared by incorporation of new rule called integrated method of Traingle, Jaccard similarity. The classical machine algorithm was executed with smart datasets proposed in [28]. Then, the analysis of data present in the leaf nodes was done for root mean square error and absolute error. Fig. 4 shows the performance of all the datasets for different data set when the attributes are selected less in number. The procedure proposed in [31] has been adopted to measure the difference in errors. The least error has been shown by the grey color integrated method. Table II shows the extent of data and its relative measures of error reduction for various similarities.

Table III provides the data related to error reduction when the number of attributes selected have been increased by 10% than earlier in k range. In Fig. 4, there is comparison of Mean Absolute Error using various similarity concept of measures.

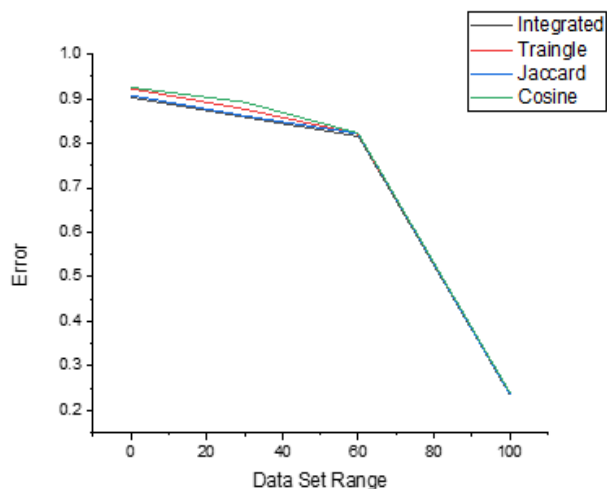


Fig. 4. Mean Absolute Error by Proposed Algorithm using Integrated Formula.

TABLE II. MEAN ABSOLUTE ERROR ON VARIED INPUTS

| Measure /Dataset Size | 10 K  | 50 K  | 100K  | 1M    |
|-----------------------|-------|-------|-------|-------|
| Cosine                | 0.732 | 0.696 | 0.625 | 0.187 |
| Jaccard               | 0.711 | 0.674 | 0.617 | 0.18  |
| Triangle              | 0.724 | 0.688 | 0.621 | 0.183 |
| Integrated            | 0.707 | 0.671 | 0.614 | 0.179 |

TABLE III. ROOT MEAN SQUARE FOR HYBRID INTEGRATED ALGORITHM FOR K DATA (HIGHLIGHTED)

| Measure/Dataset   | Flu Data Set 100K |
|-------------------|-------------------|
| Cosine            | 0.748             |
| Jaccard           | 0.729             |
| Triangle          | 0.721             |
| <b>Integrated</b> | <b>0.713</b>      |

The general systems cannot complete the algorithm with a measurable period of time whereas, integrated similarity achieves the optimal/best values for mean absolute error. The datasets used showed lowering of the percentage errors by a significant value about 6% in two data sets, 7%, 14% and 24% in other data sets respectively than the results from general methods. The results are shown in Fig. 5 for various similarity measures and the proposed method always performs best among others. The percentage change in error depends on accuracy of input data set.

Table IV provides the data related to error reduction when the number of attributes selected have been increased by 10% than earlier in m range. Fig. 6 shows details about lowering Root Mean Square Error. The reduction of errors as shown by using novel approach is 7.2 %, 17%, 21% and 12% for smart datasets. The traditional similarity measures cannot get completed within measuring time unit whereas, same is accepted by the new integrated measure. The results obtain are optimal/best for the root mean square error. The results clearly shows that optimal machine learning is adapting to the changes made in the input data sets.

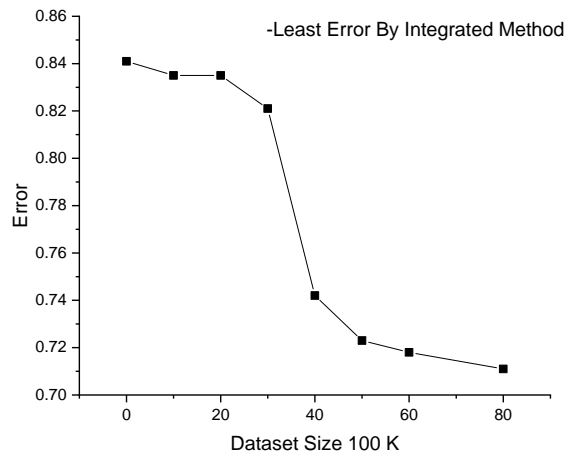


Fig. 5. Root Mean Square Error b for Dataset of Size 100k.

TABLE IV. ROOT MEAN SQUARE FOR HYBRID INTEGRATED ALGORITHM FOR M DATA (HIGHLIGHTED)

| Measure/Dataset   | Flu Dataset for 1M |
|-------------------|--------------------|
| Cosine            | 0.681              |
| Jaccard           | 0.662              |
| Triangle          | 0.679              |
| <b>Integrated</b> | <b>0.659</b>       |

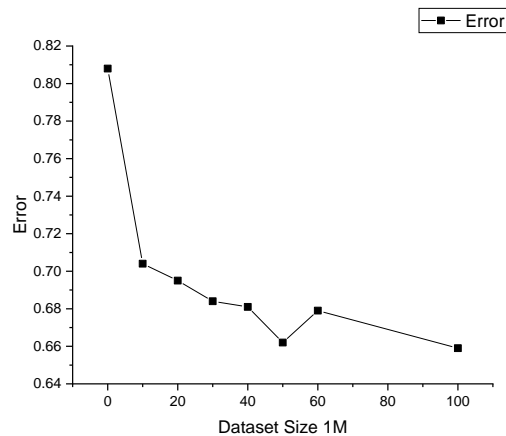


Fig. 6. Root Mean Square Error for Dataset of Size 1M.

## VI. CONCLUSION AND FUTURE SCOPE

The inflow of smart applications is having lot of impact on general overall development. The presence of such applications will improve the standard of living. The incorporation of Big Data analytics into difference applications will provide decision making and effective intelligent policies. The research proposed provides a hybrid algorithm for using smart data sets from IoT systems for future trend analysis, general systems co-ordination and accuracy. The proposed algorithm provides a layered model for various operations viz. processing, aggregation, filtering and transmission of big data. The prime role of this proposed algorithm is to make the decision trees classified by machine learning Random Forest concept. The proposed scheme also provides optimal throughput with greater variance over time. To reduce the variation of the throughput, a cross layer model will be considered in future work. The model will look further into heterogeneous challenges.

### REFERENCES

[1] Sefen, S. Baumbach, "Human Activity Recognition Using Sensor Data of Smart phones and Smart watches,"/ ICAART 2016.12.  
[2] H. Green, "The Internet of Things in the Cognitive Era: Realizing the Future and Full Potential of Connected Devices," Technical Report; IBM Watson IoT: New York, NY, USA, 2015.1.  
[3] J. Gubbi, R. Buyya, S. Marusic, Palaniswami, "M. Internet of Things (IoT): A vision, architectural elements, and future directions. Future,"Gener. Comput. Syst. 2013, 29, 1645–1660.  
[4] D. Evans, "The Internet of Things: How the Next Evolution of the Internet is Changing Everything", CISCO White Paper 2011. 1, 1–11.  
[5] J. Manyika, M. Chui, P. Bisson, J. Woetzel, R. Dobbs, J.A.D. Bughin, , "Unlocking the Potential of the Internet of Things; Technical Report; McKinsey Global Institute: New York, NY, USA, 2015.4.  
[6] I. Lee, K. Lee, "The Internet of Things (IoT): Applications, investments, and challenges for enterprises," Bus. Horiz. 2015, 58, 431–440.

[7] Y. Yoo, O.Henfridsson, K. Lyytinen, "The New Organizing Logic of Digital Innovation: An Agenda for Information Systems Research," Inf. Syst. Res. 2010, 21, 724–735.  
[8] F. Wortmann, K. Fluchter, "Internet of Things: Technology and Value Added," Bus. Inf. Syst. Eng. 2015.57, 221–224.  
[9] M. T. Yazici, S. Basurra, M. MGaber, "Edge Machine Learning: Enabling Smart Internet of Things Applications," Big Data Cognitive Computing. 2018, 2, 26; doi:10.3390/bdccc2030026.  
[10] O. D. Lara and M. A. Labrador, "A Survey on Human Activity Recognition using Wearable Sensors,"/ IEEE 2013.  
[11] D. Anguita, "Energy Efficient Smartphone-Based Activity Recognition using Fixed-Point Arithmetic,"/ JUCS 2013.10.  
[12] J Yin, "Sensor-Based Abnormal Human-Activity Detection,"/ IEEE 2008.11.  
[13] K. Sree Divya, P. Bhargavi, S. Jyothi, "Machine Learning Algorithms in Big data Analytics," International Journal of Computer Sciences and Engineering, Vol.6, Issue.1, pp.63-70, 2018.  
[14] S. Shadroo, A.M. Rahmani, "Systematic survey of big data and data mining in the internet of things," (2018) Comput Netw 139:19–47.  
[15] A. Paul, S. Rho, "A probabilistic model for M2M in IoT networking and communication," (2016) Telecommun Syst 62(1):59–66.  
[16] C.T Tran, M. Zhang , P. Andreae, B. Xue , L.T. Bui, "An effective and efficient approach to classification with incomplete data," (2018) Knowl Based Syst 154:1–16.  
[17] K. Lakshmanaprabu, K. Shankar, M. Ilayaraja, N. A. Wahid, V. Vijayakumar and N. Chilamkurti, " Random forest for big data classification in the internet of things using optimal features," (2019) International Journal of Machine Learning and Cybernetics. 10. 10.1007/s13042-018-00916-z.  
[18] R. Lin, Z. Pei, Z. Ye, B. Wu, G. Yang, "A voted based random forests algorithm for smart grid distribution network faults prediction", (2019) Enterprise Information Systems. 14. 1-19. 10.1080/17517575. 2019.1600724.  
[19] M. H. Alsharif, A. H. Kelechi, K. Yahya, S.A. Chaudhry, ". Machine Learning Algorithms for Smart Data Analysis in Internet of Things Environment," (2020) Taxonomies and Research Trends. Symmetry. 12. 88. 10.3390/sym12010088.  
[20] D. G. Gil, F. L.Sánchez, J. Luengo S. García and F. Herrera, "From Big to Smart Data: Iterative ensemble filter for noise filtering in Big Data classification," (2019) International Journal of Intelligent Systems. 34. 10.1002/int.22193.  
[21] B. Devi, S. Kumar, Anuradha and V.G. Shankar, " AnaData: A Novel Approach for Data Analytics Using Random Forest Tree and SVM," (2019) In: Iyer B., Nalbalwar S., Pathak N. (eds) Computing, Communication and Signal Processing. Advances in Intelligent Systems and Computing, vol 810. Springer, Singapore. https://doi.org/10. 1007/978-981-13-1513-8\_53.  
[22] S. Oprea and A. Bâra, "Machine Learning Algorithms for Short-Term Load Forecast in Residential Buildings Using Smart Meters, Sensors and Big Data Solutions," in IEEE Access, vol. 7, pp. 177874-177889, 2019, doi: 10.1109/ACCESS.2019.2958383.  
[23] D. G. Gil, A.A Barros, J.Luengo, S. Garcia and F. Herrera, "Big Data Preprocessing as the Bridge between Big Data and Smart Data: BigDaPSpark and BigDaPFLink Libraries. 324-331. 10.5220 /0007738503 240331.  
[24] J. H. Tato, "Using Smart Persistence and Random Forests to Predict Photovoltaic Energy Production," (2018) Energies. 12. 100. 10.3390/en12010100.  
[25] P. Kaur, R. Kumar and M. Kumar, " A healthcare monitoring system using random forest and internet of things (IoT). Multimed Tools Appl 78, 19905–19916 (2019). https://doi.org/10 .1007 /s 11042-019- 7327-8.  
[26] BMWi, "Smart-data-technologien (German)," Smart Data Accompanying Research, Tech. Rep., 2015. [Online]. Available: https://www.digitale-technologien.de/DT/Redaktion/DE /Downloads /Publikation/smartdata\_brochure\_english.pdf?\_\_blob=publicationFile&v =18.

- [27] D. Masroof, Munishwar Rai, "A Novel Framework for Enhancing QoS of Big Data", (IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 11, No. 4, 2020.
- [28] M. Denil, D. Matheson and N. D. Freitas, "Narrowing the Gap: Random Forests In Theory and In Practice", Proceedings of the 31st International Conference on Machine Learning, Beijing, China, 2014. JMLR: W&CP volume 32.
- [29] P. Jaccard, "Nouvelles recherches sur la distribution florale," Bull Soc Vaud Sci Nat. 1908; 44:223±270.
- [30] M. M. Gaber, "Data stream mining using granularity-based approach. In Foundations of Computational, Intelligence," Volume 6; Springer: Berlin/Heidelberg, Germany, 2009; pp. 47–66. 8.
- [31] S-B. Sun, Z-H. Zhang, X-L Dong, H-R Zhang, T-J. Li, L. Zhang, F. Min, "Integrating Triangle and Jaccard similarities for recommendation," (2017) PLoS ONE 12(8): e0183570. <https://doi.org/10.1371/journal.pone.0183570>.

# Grid Search Tuning of Hyperparameters in Random Forest Classifier for Customer Feedback Sentiment Prediction

Siji George C G<sup>1</sup>

Research Scholar  
Department of Computer Science  
CMS College of Science and Commerce  
Coimbatore, Tamilnadu, India

B.Sumathi<sup>2</sup>

Associate Professor  
Department of Computer Science  
CMS College of Science and Commerce  
Coimbatore, Tamilnadu, India

**Abstract**—Text classification is a common task in machine learning. One of the supervised classification algorithm called Random Forest has been generally used for this task. There is a group of parameters in Random Forest classifier which need to be tuned. If proper tuning is performed on these hyperparameters, the classifier will give a better result. This paper proposes a hybrid approach of Random Forest classifier and Grid Search method for customer feedback data analysis. The tuning approach of Grid Search is applied for tuning the hyperparameters of Random Forest classifier. The Random Forest classifier is used for customer feedback data analysis and then the result is compared with the results which get after applying Grid Search method. The proposed approach provided a promising result in customer feedback data analysis. The experiments in this work show that the accuracy of the proposed model to predict the sentiment on customer feedback data is greater than the performance accuracy obtained by the model without applying parameter tuning.

**Keywords**—Classification; grid search; hyperparameters; parameter tuning; random forest classifier; sentiment analysis

## I. INTRODUCTION

The Classification is a text mining tasks in which class of a particular input is identified by using a given set of labelled data. Both supervised and unsupervised methods are used for classification. In the first method, learning is done through predefined labelled data. In this, a set of labelled input documents are given to the model by the end-user. The two main categories of supervised learning are parametric and non-parametric classification. The probability distribution of each class is the base of parametric classification. If the density function is known, it will be better to use non-parametric classification. Recently, people are using this classification process especially supervised classification to develop multiple interesting platforms for business. Sentiment analysis is the most attractive platforms which make use of the advantages of supervised classification methods.

Sentiment can be described as a person's feeling about a particular thing. It includes the task of binary classification in which documents are classified into two different classes such as positive sentiment or negative sentiment. Due to the fast popularity of social networks [1], people are using it for

sharing their views, opinion and ideas. Social networks provide a platform for the people to create a virtual civilization [2]. Sentiment analysis is a mining process based on user-generated comments to identify positive or negative feelings. Opinions are always important to a business. Most of the business decision is performed based on customers' reviews. The analysis of customer or product review involves the extraction of sentiment from product document [3]. Business organizations are very conscious to know whether customers like their product or service, what customers feel about the product, which type of product or service customers like or dislike, etc.

Sentiment analysis is usually applied text input which help to identify the sentiment in a particular document and thus it is considered as the main part of text mining. Other than text classification, it requires more knowledge of the language. Generally, machine learning algorithms are considering the occurrence of the words in a document, so it tough to recognize the supreme attitude in that specific document. The sentiment analysis should be the process of identifying the polarity present in the given text or document i.e., positive or negative.

There are number of supervised machine learning algorithms are used for sentiment analysis. The performance of these classification algorithm is depending on its specific domain [4]. Random Forest classifier is largely used for this purpose. It is considered as an ensemble method [5] which generates many classifiers and finally aggregates their result for prediction. This will create a number of decision trees in the training phase [6]. The risk of noise and outliers will be high when having a single tree in classifier and it will definitely reduce the output of the processing. Due to the randomness property of Random Forest classifier, it is highly robust to outliers and noises. This classifier can handle missing values also.

One better approach to increase the outcome of any classifier is to tune the hyperparameters of that classifier [7]. The parameters that are set by the data analysts before the training process is called hyperparameters and it is independent of the training process. For example, in a random forest, a hyperparameter would be how many trees have to be



included in the forest or how many nodes each tree can have. Optimizing these hyperparameters for the classifier is the key to the perfect prediction of unlabeled data. These can only be achieved through trial and error methods. Different values of hyperparameters are used, then compare their result and finally find the best combination of them. The tuning process of hyperparameters is mainly depended on experimental results and not the theoretical result.

In this work, the Grid Search approach is applied for tuning Random Forest classifier and tried to identify the best hyperparameters. The implementation of Grid Search is simple [8]. A set of hyperparameters and their values are feed to it first and then run an exhaustive search overall all possible combination of given values then training the model for each set of values. Then Grid Search algorithm will compare the score of each model it trains and keeps the best one. A common extension of Grid Search is to use cross-validation i.e., training the model on several different folds with different hyperparameter combinations to find more accurate results.

The rest of the paper is organized as follows. In Section II, previous work in these research topics are discussed. Section III explains the proposed system model and architecture. The experimental results are discussed in Section IV and it is followed by a conclusion in Section V.

## II. RELATED WORK

Rafael G. Mantovani et al. [9] made an investigation on random search and grid search methods. They aimed to tune the hyperparameters of the classifier called Support Vector Machine (SVM). Their experiment was performed by using a huge dataset, finally, they compared the performance of Random Search with four methods such as Particle Swarm Optimization, Genetic Algorithm, Grid Search method and Estimation of Distributed Algorithm. The result of this work reveals that the predictive power of SVM classifier with Random Search is same as the other four techniques used and the advantage of this combination was the lowest computational cost of the model.

Xingzhi Zhang et al. [10] effort was to propose an optimized novel of Random Forest Classifier for credit score analysis. For optimizing Random Forest Classifier, the authors developed a system called NCSM which uses grid search and feature selection. The developed model has the capability to overcome the problem of irrelevant and redundant features and got good performance accuracy. The model used the information entropy to select the optimal features. From the UCI database, two sets of data are selected as input to examine the performance of developed model. Their experiments show that proposed system has dominating the performance of some other methods.

A hybrid approach based on Random Forest and Support Vector Machine is proposed by Yassine Al Ambrani et al. in 2018 [11] for identifying Amazon product reviews. Cross-validation method with fold value 10 has been used for this work. Both Support Vector Machine and Random Forest Classifier are used by authors to do classification of product reviews. The classification result of both classifiers is with the

hybrid method. The result shows that the hybrid method of random Forest and SVM outperforms the individual methods.

An ensemble-based customer review sentiment analysis is done in 2019 [12] by Ahlam Alrehili and Kholood Albalawi. The proposed method used a voting system which combines five classifiers Random Forest, Naive Bayes, SVM, bagging and boosting. Six different scenarios are performed by authors to measure the result of the proposed model against five used classifiers. They are using unigram (with/without) stop words removal, bigram (with/without) stop words removal and using trigram stop word removal. Among this, the highest accuracy of 89.87% is given by the Random Forest classifier.

Sentiment analysis on the blogs are carried by Prem Melville et al. in 2009 [13]. They combined classification of text with lexical knowledge. A unified framework is proposed by the authors and the framework used lexical information to filter information for a specific domain. The combination of training examples using Linear Pooling with background knowledge is performed well and had an accuracy of 91.21%.

The neural network has more hyperparameters which have to be set by hand. Nauria Rodriguez-Barroso et al. worked on these neural network parameters in 2019 [14]. They used SHADE evolutionary algorithm to perform optimization of different deep learning hyperparameters to perform twitter sentiment analysis. The Spanish tweets are selected as dataset for their work. The findings reveal that hyperparameters selected by SHADE algorithm help to improve the proposed model's performance.

Airline data sentiment analysis is performed by Bahrawi in 2019 [15]. Six airline tweet data from Kaggle is used for this study and Random Forest classifier is used for sentiment prediction. Classifier predicted 63% of tweets as negative, 21% as neutral and 16% as positive. The accuracy achieved by the Random Forest algorithm was only 75%. The author suggested to build model by using some other machine learning algorithms to get a better result.

A new credit scoring model called NCSM is proposed by Xingzhi et al. [16] in 2018. Grid search method and feature selection are applied for this model in order to optimize the Random Forest classifier's performance. This proposed model achieved high prediction accuracy as compared with some other commonly used methods.

## III. PROPOSED MODEL

The architectural diagram of the proposed model is depicted in Fig. 1. The collected customer feedback data go through several processing stages and feature extraction is performed. After extracting the necessary features, it is given as input to the Random Forest classifier. Finally, parameter tuning by Grid Search method is applied to increase the classifier's performance. This section gives the detailed description of the proposed model.

### A. Pre-Processing of Data

As an initial step, the original input data is examined in pre-processing stage and make the raw data convenient for using in classification process. It is the first and crucial step in creating a model. While creating a machine learning model, it

is not always possible to get clean and formatted data. For this, the data pre-processing task is used. The real-world data may be in an unusable format and contains missing values, noises, etc. This type of data is impossible to use directly for machine learning model. Data pre-processing is an important task to clean the original data for the machine learning model and thereby increase model accuracy and efficiency. The following steps are used for pre-processing:

- Tokenization- Divided the customer feedback input into a number of individual words called tokens.
- Removal of special characters, numbers, stop words and punctuations since it does not any sentiment.
- Stemming- It involves normalizing the input data. For example, reducing words like loves, loving and lovable into its root word i.e., love is often used in the same context.

### B. Feature Extraction

In this step, new features are extracted from existing dataset and thereby reduce the count of features used for processing task. The new reduced feature set will be capable to represent majority information in the initial feature set. This text feature extraction will directly influence the accuracy of the classification. The two techniques used in this work for feature extraction are.

- CountVectorizer
- Term Frequency-Inverse Document Frequency (TF-IDF)

Count Vectorizer: The text data need special preparation before using it for predictive modelling. The number of occurrences of every word in a given document can be identified by using Count Vectorizer. It will provide a vector with frequencies of each token in the given document. Term Frequency-Inverse Document Frequency: TF represent the result after dividing the occurrence of a word in a particular document by the total count of words present in that document. IDF is used to find out the weight of rare words across the entire document in the corpus. When TF is multiplied by IDF it will result in TF-IDF.

### C. Algorithms Applied

1) *Sentiment classification-random forest classifier:* Random Forest classifier is a flexible supervised algorithm which can be used for text classification. The working of this algorithm is based on tree collection in which every tree depends on different random variables [17]. It uses Divide-and-conquer approach. Forest represents a collection of many trees. From random subsets of input data, this algorithm will generate several small decision trees. Consider a random vector of dimension n, where  $A = (A_1, A_2, \dots, A_n)^T$  is a set of real-valued input variables and a random variable B which represent the real value response, then we assume an unknown joint distribution  $P_{AB}(A, B)$ . The goal of this algorithm is to find a prediction function  $f(A)$  for predicting B. A loss function  $L(B, f(A))$  is used to find the prediction function and it should minimize the expected value of the loss.

$$E_{AB}(L(B, f(A))) \tag{1}$$

$L(B, f(A))$  is used to represent how prediction function  $f(A)$  is close to B. Zero-one loss is the choice of L for classification. Minimizing  $E_{A,B}(L(B, f(A)))$  for zero-one loss gives

$$L(B, f(A)) = I(B \neq f(A)) = \begin{cases} 0 & \text{if } B = f(A) \\ 1 & \text{otherwise} \end{cases} \tag{2}$$

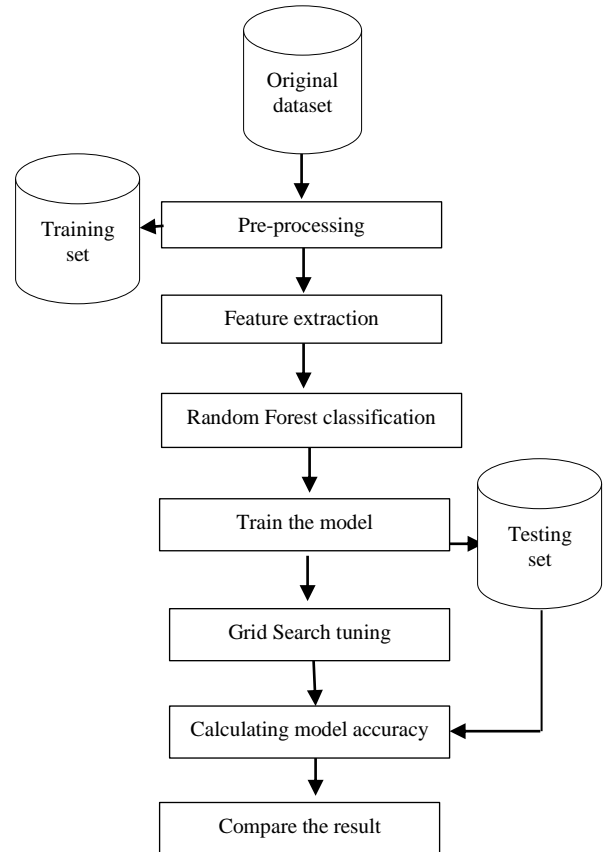


Fig. 1. System Architecture.

The procedure for the Random Forest classifier is given as follows:

**RandomForestProc()**

Input: Training Data D, Number of Trees T, Total Features TF, Subset of Features tf

Output: labelled classes for the input dataset

- Repeat the followings for each trees in the forest T:
- Consider and choose a bootstrap sample S with size D from training data
- Recursively repeat the followings for generating the tree t
- Randomly select tf from total features TF
- Choose the better feature among TF
- Node to be split
- After generating trees, test instance should be given to every tree then based on majority votes, class label will be assigned.

2) *Hyperparameter tuning- grid search method:* Machine learning model has many parameters [18] to tune and by tweaking these parameters, the performance of the model can improve. Hyperparameter tuning is the best method to execute a different number of parameter combinations to assess a classifier’s performance. Assessing a classifier by using training data will cause a fundamental machine learning problem called overfitting. The overfitting is the situation in which a model performs poorly on test data and highly on training data. Therefore, cross-validation is used with the grid search method for hyperparameter optimization.

The grid search method is an approach used to identify the optimum parameters of a classifier so that a model can accurately predict some unlabeled data. The Grid Search method is used to tune some hyperparameters which cannot directly learn from the training process. The classification model has many hyperparameters and finding the best combination of these parameters is a challenging process. One of the best methods used for this purpose is the Grid Search method. Suppose, a machine learning model X has hyperparameters h1, h2 and h3. The Grid Search method defines a range of values for each hyperparameter h1, h2 and h3. It will construct many versions of X with all possible combinations of h1, h2 and h3. This range of hyperparameter values is known as a grid.

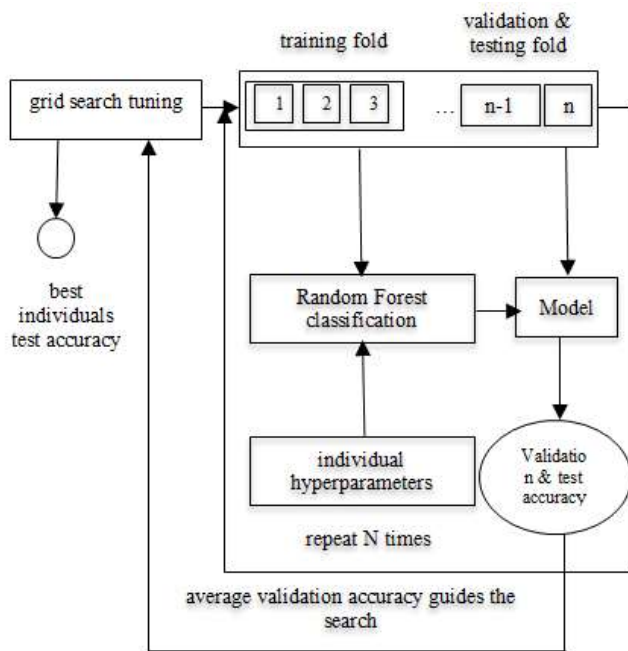


Fig. 2. Hyperparameter Tuning Architecture.

The hyperparameter tuning architecture is depicted in Fig. 2. The input data is divided into a training set, testing set and validation set. The tuning process is executed by separating the data set into n different portions. Then, the Random Forest Classifier trained in n-2 portions for each candidate solution selected by the tuning technique. The validation set is used to validate the developed model and the last portion is used to test the model. The test accuracy and validation accuracy are evaluated by using the model. Then

the model is instigated by the training set and the hyperparameter value determined by the tuning technique. These steps are repeated for N times. To guide the search process, the average validation accuracy is used as the fitness value. Finally, it will return the individual with the highest accuracy and the performance of the method is the average test accuracy of that individual. The procedure for the proposed hybrid model of Random Forest classifier and Grid Search method for sentiment prediction is as follows.

**RandomForest+GtridSearchProc( )**

- Consider binary classification dataset of product reviews N samples and split it in a train and test set
- Define a pipeline with Random Forest Classifier
- Setup a grid for total number of features to be used and total number of trees to be constructed in the forest
- Define a function to run Grid Search method which takes the input such as defined pipeline, parameter grid and train and test set
- Define the objective function which takes set of hyperparameters and output the accuracy score

$$accuracy = f(hyperparameters)$$

- Select a random combination
- Define the number of search iteration
- Iterate through all possible combination of values specified in the grid one at a time
- Pass these values to the objective function
- Repeatedly execute the objective functions for each and every combination of hyper-parameter values
- Evaluate the best hyper-parameter which maximize the accuracy

$$hyperparams *= argmax f(hyperparams)$$

**IV. EXPERIMENTS AND RESULTS**

Labelled customer feedback data on electronic items collected from UCI database and it is used as input for this work. It includes 1500 reviews (750 positive and 750 negative reviews). This work aims to classify these customer feedbacks into two different categories such as positive feedback and negative feedback. 7-fold cross-validation is used for calculating the model’s accuracy. First, customer feedback data analysis is performed by using Random Forest classifier with default hyperparameters and achieved 84.53% of accuracy. Table I gives the result of customer feedback analysis using Random Forest classifier.

From the Table I, it is clear that 1268 customer reviews are classified correctly among 1500 and 232 are wrongly classified by the model.

TABLE I. RANDOM FOREST CROSS VALIDATION RESULT

|          | Positive | Negative | Total |
|----------|----------|----------|-------|
| Positive | 682      | 68       | 750   |
| Negative | 164      | 586      | 750   |
| Total    | 846      | 654      | 1500  |

To increase the accuracy of the classifier, parameter tuning using Grid Search method is used in this work. The Random Forest classifier has several parameters, which can be adjusted to get optimal performance. Two of those parameters are the number of trees constructed for classifying new data and the maximum number of variables used in individual trees. The class GridSearchCV available in Scikit Learn is used for this study. The GridSearchCV evaluates, all possible combinations of parameter values and finally, the best parameter combination is retained. This work mainly concentrates on two parameters of Random Forest classifier.

The GridSearchCV uses max\_features for denoting the maximum number of variables used in independent trees and n\_estimators for denoting the total number of trees to be constructed in the forest. The Table II provides the result of parameter tuning of Random Classifier on customer feedback data. The score in Table II represents the accuracy of the classifier using the 7-fold cross-validation method. sqrt and log2 are the two options tried for max\_features.

According to Table II, the highest accuracy of the Random Forest Classifier is 90.02% at the parameters 'max\_features'='sqrt' and 'n\_estimators'=400.

The Table III shows the result of the proposed method (with best parameters) which uses the Grid Search approach for hyperparameter tuning. By using the proposed method, among 1500 reviews, 1353 reviews of customers are classified correctly and 147 are not. The accuracy comparison of two used methods is given in Table IV.

Fig. 3 depicts the total number of instances classified by Random Forest Classifier and proposed system.

Fig. 4 depicts the most frequent words identified by the proposed system. It shows the top 15 words appeared in the customer feedback data and from the figure it is clear that the most frequent word is product.

A detailed comparison of Random Forest classifier and proposed system which uses the Grid Search method for parameter tuning is depicted in Table V.

TABLE II. BEST PARAMETERS IDENTIFIED BY GRID SEARCH

| No | Tuning Parameters  | Score(7-foldcross validation) | Best parameter                            |
|----|--|-------------------------------|---|
| 1  | n_estimators=[10,100,1000,1500]<br>max_feature=[sqrt , log2] | 87.04                         | {n_estimators=1000<br>max_features=sqrt } |
| 2  | n_estimators=[600,1000,1300]<br>max_features= [sqrt]         | 89.56                         | {n_estimators=600<br>max_features= sqrt}  |
| 3  | n_estimators=[400,600,800,900]<br>max_features= [sqrt]       | 90.02                         | {n_estimators= 400<br>max_features=sqrt,} |
| 4  | n_estimators=[300,400,500,550]<br>max_features= [sqrt]       | 90.02                         | {n_estimators= 400<br>max_features=sqrt,} |
| 5  | n_estimators=[325,350,400,450]<br>max_features= [sqrt]       | 90.02                         | {n_estimators= 400<br>max_features=sqrt,} |
| 6  | n_estimators=[340,380,400,425]<br>max_features= [sqrt]       | 90.02                         | {n_estimators= 400<br>max_features=sqrt}  |

TABLE III. CROSS VALIDATION RESULT OF PROPOSED METHOD

|          | Positive | Negative | Total |
|----------|----------|----------|-------|
| Positive | 721      | 29       | 750   |
| Negative | 118      | 632      | 750   |
| Total    | 839      | 661      | 1500  |

TABLE IV. COMPARISON BASED ON CLASSIFIED INSTANCES

|                 | True classification | Wrong classification | Accuracy (%) | Time taken (Seconds) |
|-----------------|---------------------|----------------------|--------------|----------------------|
| Random Forest   | 1268                | 68                   | 84.53        | 6.70                 |
| Proposed Method | 1353                | 147                  | 90.02        | 8.02                 |

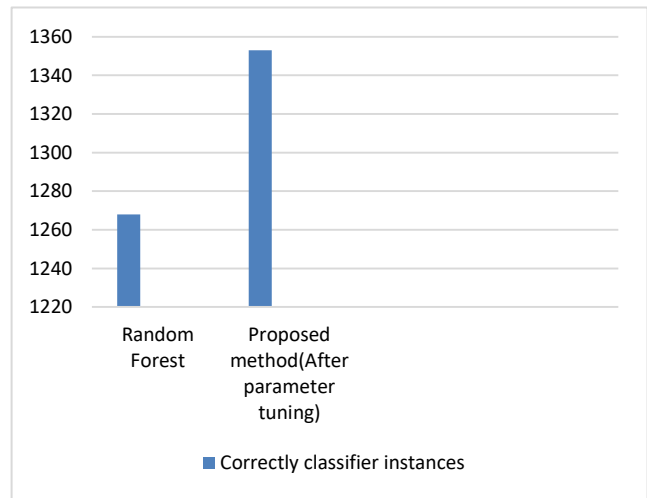


Fig. 3. Number of Instances Classified Correctly.

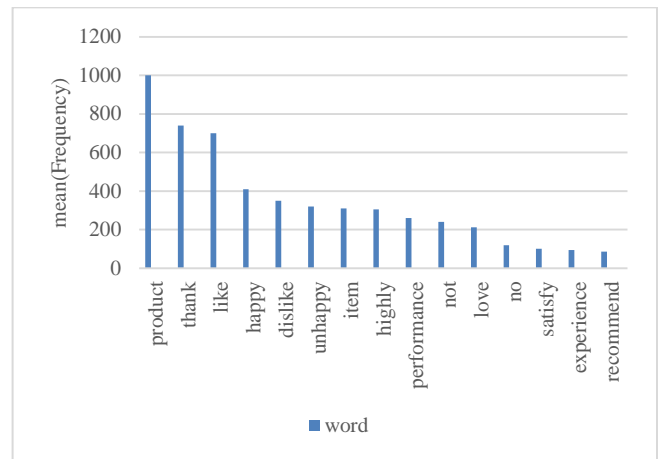


Fig. 4. Most Frequent Words.

TABLE V. COMPARISON

|             | Random Forest | Proposed Method(Random Forest + Grid Search) |
|-------------|---------------|--|
| Sensitivity | 90.93         | 96.13  |
| Specificity | 78.13         | 84.26  |
| Recall      | 90.93         | 96.13  |
| Precision   | 80.61         | 85.93  |
| F-Measure   | 85.45         | 90.74  |

## V. CONCLUSION

Sentiment analysis is essential for a business organization to perform decision making. It can be used for different tasks such as calculating or expressing sentiment on any product or service. In this work, the best parameters are tuned by Grid Search method for Random Forest classifier. Experimental results on customer feedback data show that Random Forest provides the best result with an accuracy of 84.53%. But, by tuning number of maximum trees in the forest and depth of trees, the accuracy of the developed model increases to 90.02%. The result shows that parameter tuning has successfully helped to generate the best model to classify new data. At the same time, the Random Forest classifier take more execution time when the number of trees in the forest is increased. In the future work, the proposed model can use for multi-class sentiment prediction since it concentrated binary classification only.

## REFERENCES

- [1] M. Ahmad, S. Aftab, S.S Muhammad, and S. Ahmad, 2017. Machine learning techniques for sentiment analysis- A review, *International Journal of Multidisciplinary Science and Engineering*, vol. 8, no. 3, pp. 27-32.
- [2] S. H. Yadav, and P. M. Manwatkar, 2015. An approach for offensive text detection and prevention in social network, 2015 IEEE International Conference Innovations in. Information, Embedded and Communication Systems (ICIECS), pp. 3-6.
- [3] Bagus Setya Rintyarna, Riyanarto Sarno, and Chastine Fatichah, 2020. Enhancing the Performance of Sentiment Analysis Task on Product Reviews by Handling Both Local and Global Context, *International Journal of Information and Decision Sciences*.
- [4] R. Xia, C. Zonga, and S. Li, 2011. Ensemble of feature sets and classification algorithms for sentiment classification, *Information Sciences*, Elsevier, vol. 181, pp.1138-1152.
- [5] Yashaswini Hegde, and S.K. Padma, 2017. Sentiment Analysis Using Random Forest Ensemble for Mobile Product Reviews in Kannada, 7<sup>th</sup> International Advance Computing Conference(IACC), IEEE.
- [6] Shahnoor C. Eshan, and Mohammad S Hasan, 2017. An Application of Machine learning to Detect Abusive Bengali Text, 20<sup>th</sup> International Conference of Computer and Information Technology(ICCIT).
- [7] Muhammad Murtadha Ramadhan, Imas Sukaesih Sitanggaang, Fahredi Rizky Nasution and Abdullah Ghifari, 2017. Parameter Tuning in random Forest Based on Grid Search Method for Gender Classification Based on Voice Frequency, *International Conference on Computer, Electronics and Communication Engineering*.
- [8] J. Bergstra, and Y. Bengio, 2012. Random search for hyper-parameter optimization, *Journal of Machine Learning Research*, vol. 13, pp. 281-305.
- [9] Rafael G. Mantovani, Andre L. D. Rossi, Joaquin Vanschoren, Bernd Bischl and Andre C. P. L. F., 2015. Effectiveness of Random Search in SVM hyper-parameter Tuning. *IEEE Proceedings of the 2015 International Joint Conference on Neural Networks*, July 2015.
- [10] Xingzhi Zhang, Yan Yang, and Zhurong Zhou, 2018. A Novel Credit Scoring Model based on Optimized Random Forest. 8<sup>th</sup> Annual Computing and Communication Workshop and Conference (CCWC).
- [11] Yassine Al Ambrani, Mohamed Lazaar, and Kamal Eddine El Kadiri, 2018. Random Forest and Support Vector Machine Based Hybrid Approach to Sentiment Analysis. *The First International Conference on Intelligent Computing in Data sciences*, vol. 127, pp. 511-520.
- [12] Ahlam Alrehili and Kholood Albalawi, 2019. Sentiment Analysis of Customer Reviews using Ensemble Method, *International Conference on Computer and Information Sciences (ICCIS)*, IEEE.
- [13] Prem Melville, Wojciech Gryc, and Richard D. Lawrence, 2019. Sentiment Analysis of Blogs by Combining Lexical Knowledge with Text Classification, *Knowledge Discovery and Datamining (KDD'09)*, Paris France.
- [14] Nuria Rodriguez-Barroso, Antonio R. Moya, Jose A. Fernandez, Elena Romero, Eugenio Martinez-Camara, and Francisco Herrera, 2019. Deep Learning Hyper-Parameter Tuning for Sentiment Analysis in Twitter Based on Evolutionary Algorithms, *Proceedings of federated Conference on Computer Science and Information Systems*, pp. 255-264.
- [15] Bahrawi, 2019. Sentiment Analysis using Random Forest Algorithm Online Social Media Based, *Journal of Information Technology and Its Utilization*, vol. 2, issue 2.
- [16] Xingzhi Zhang, Yan Yang, and Zhurong Zhou, 2018. A Novel Credit Scoring Model based on Optimized Random Forest, 8<sup>th</sup> Annual Computing and Communication Workshop and Conference (CCWC).
- [17] Adele Cutler, D. Richard Cutler, and John R. Stevens, 2012. Random Forests, *Ensemble Machine Learning*, pp. 157-175.
- [18] Hitesh H Parmar, Sanjay Bhandari, and Glory Shah, 2014. Sentiment Mining of Movie Reviews using Random Forest with Tuned Hyper-parameters, *International Conference on Information Science*.

# Aspect-Based Sentiment Analysis and Emotion Detection for Code-Mixed Review

Andi Suciati<sup>1</sup>, Indra Budi<sup>2</sup>  
Faculty of Computer Science  
Universitas Indonesia  
Depok, Indonesia

**Abstract**—Review can affect customer decision making because by reading it, people manage to know whether the review is positive, or negative. However, positive, negative, and neutral, without considering the emotion will be not enough because emotion can strengthen the sentiment result. This study explains about the comparison of machine learning and deep learning in sentiment as well as emotion classification with multi-label classification. In machine learning comparison, the problem transformation that we used are Binary Relevance (BR), Classifier Chain (CC), and Label Powerset (LP), with Decision Tree (DT), Random Forest (RF), Support Vector Machine (SVM), and Extra Tree Classifier (ET) as algorithms of machine learning. The features we compared are n-gram language model (unigram, bigram, unigram-bigram). For deep learning, algorithms that we applied are Gated Recurrent Unit (GRU) and Bidirectional Long Short-Term Memory (BiLSTM), using self-developed word embedding. The comparison results show RF dominates with 88.4% and 89.54% F1 scores with CC method for food aspect, and LP for price, respectively. For service and ambience aspects, ET leads with 92.65% and 87.1% with LP and CC methods, respectively. On the other hand, in deep learning comparison, GRU and BiLSTM obtained similar F1- score for food aspect, 88.16%. On price aspect, GRU leads with 83.01%. However, for service and ambience, BiLSTM achieved higher F1-score, 89.03% and 84.78%.

**Keywords**—Sentiment analysis; emotion; multi-label classification; machine learning; deep learning

## I. INTRODUCTION

Review is an evaluation to entities such as product, restaurant, place, etc. that can be used by customers or owner as product input. This review usually contains several aspects such as in laptop [1], the aspects that can be evaluated are hardware, price, etc. This evaluation can affect the decision making from customer. For instance, when people want to go to trip, they will read the review of several places and compare them. One of domain examples that usually get many reviews is restaurant. There are several platforms in internet for restaurant review, such as Zomato<sup>1</sup> and Yelp<sup>2</sup>. In the platform, mostly people only see the ratings of the restaurant, however reading the review is very important because the customers will obtain specific information rather than only seeing the ratings. In addition, sometimes people also give ratings that are very different from the actual review. So, it can be

concluded ratings not always give the information about the quality of restaurant. Beside for decision making of customer, review also important for the product owner. Pontiki et al. [2] stated that feedback from customer will help companies measure their customer satisfaction, and for the development of their product and services they provide. For identifying the sentiment of aspect, sentiment analysis can be conducted. However, classifying the sentiment is not enough without considering the emotions from customers. Knowing the emotion can strengthen the sentiment results from a review. Furthermore, mostly a review contains two or more languages, or called code-mixed languages. This kind of review is difficult to understand by computer because computer cannot identify the languages easily like human. This also a big challenge for sentiment analysis and emotion detection. There are several classification methods that can be used, such as machine learning and deep learning. Mohammad et al. [3] used Support Vector Machine when classifying sentiment data from Twitter<sup>3</sup>. In the other hand, Stojanovski et al. [4] applied deep learning algorithm for sentiment analysis and emotion detection for Twitter data.

This research focuses to conduct sentiment analysis an emotion detection in every aspect that appeared in a restaurant review. The data were collected from Indonesian restaurant review platform, named PergiKuliner<sup>4</sup>, and this study using ‘food’, ‘price’, ‘service’, and ‘ambience’ as aspects. The sentiment polarities that were used for emotions are ‘positive’, ‘negative’, and ‘neutral’, while ‘happy’, ‘sad’, ‘surprised’, and ‘neutral’. The addition of ‘neutral’ because there is a possibility that a review contains sentiment polarity, but the emotion is difficult to detect. The method of classification that we applied is multi-label classification while the algorithms that we used are from machine learning and deep learning.

The rest of paper was organized into: in Section 2, we explained about several researches that related to our study. In Section 3, we illustrate the research steps of our experiments. For Section 4, we showed the classification results as well as analyzing them. Then in last part, we concluded the results and future work for this study.

<sup>1</sup> <https://www.zomato.com>

<sup>2</sup> <https://www.yelp.com/>

<sup>3</sup> <https://twitter.com/home>

<sup>4</sup> <https://pergikuliner.com/>

## II. RELATED WORK

There are many studies about sentiment analysis and emotion detection. Mohammad [5] did a literature studies regarding several researches about valence, emotion, and other aspects that can affect the feeling from a person. From that study, the writer describes the challenges for sentiment and emotion detection, such as language complexity, non-standardized language, lack of labeled data, subjectivity, culture differences, etc. Stojanovski et al. [4] did a sentiment analysis research using SemEval 2015<sup>5</sup> and emotion detection using Twitter data. The sentiment polarities that we used are 'positive', 'negative', and 'neutral', while for emotions, we utilized 'love', 'joy', 'surprise', 'anger', 'sadness', 'fear', and 'thankfulness'. After that, the writer applied Deep Convolutional Neural Network for sentiment and emotion detection. However, the sentiment analysis and emotion detection were conducted in separated dataset. Another study about emotion was conducted by Hassan et al. [6]. This study was emotion classification using Skip-thought Vector. Khawaja et al. [7] also did an experiment about emotion which is developing an automatic lexicon for emotion.

In Indonesia, there are also few researches about sentiment and emotion. Wikarsa dan Tahir [8] studied about emotion detection using data from Twitter, but the data were in English. Savigny and Purwarianti [9] also conducted emotion classification using YouTube<sup>6</sup> comments. For sentiment analysis, [10][11] studied it for restaurant review in Indonesia.

Several studies also have conducted for sentiment analysis and emotion detection using code-mixing data. Shalini et al. [12] studied sentiment analysis for Facebook<sup>7</sup> comments with Kannada-English languages. The experiment was done by applying Facebook's fast text, Doc2Vec with SVM, Bidirectional LSTM, and CNN. Lee and Wang [13] experimented using Chinese-English data and proposed multi-learning framework for emotion detection.

## III. RESEARCH STEPS

This section explains the methodology that applied in this research as shown by Fig. 1.

### A. Data Collection

The data were collected from PergiKuliner platform by scraping them. The collected data are the reviews for several restaurants in Jakarta, Bogor, Depok, Tangerang, and Bekasi, and the total are 20000 reviews. After filtering the data, such as deleting the duplicate and removing the spam reviews, the final data that annotated are 18908 reviews. The data were including reviews that use Indonesian, English, and code-mixed (Indonesian-English). Below are the examples of data:

1) *Indonesian*: Akhirnya cobain taichan sm martabak tipkernya Dann taichannya enak!! Hehehe Asik jg tmptnya rame. (Finally, can taste its thaichan and martabak tipker and

the taichan was delicious!! Hehehe it was fun, the place also crowded.).

2) *Mixed*: Finally got to try this current happening Korean food! Gyeran Jim (22k) Ini kaya steamed egg, yang rada di bake. Telornya ga tawar, tasty dan pinggirannya agak kering gitu. Menurut gue worth sih 22k buat ini, hehe. Probably gonna try again :) (Finally got to try this current happening Korean food! This Gyeran Jim (22k) was like steamed egg. The egg wasn't blend, tasty and the crust is bit dry. In my opinion 22K was worth for this, hehe. Probably gonna try again :))

3) *English*: Been here for several times I've been loving this place so much. The ambience is truly Japanese izakaya dining. If you eat with many people (sharing) the price would be reasonable, however if you only eat for two the price might get a little high for izakaya. Though the foods are mostly great. Cool place to hangout!

### B. Building Annotation Guidelines

After collecting data, next step is building the annotation guidelines. There are two annotation guidelines that were made. First is annotation guideline for sentiment annotation, and another one is for emotion annotation. The aspects that used 'food', 'price', 'service', and 'ambience'. The sentiment polarities that used, following [14], which are 'positive', 'negative', and 'neutral', while for emotions, we followed [15], that divided emotions into 'happy', 'sad', 'surprised', 'angry', 'disgusted', and 'fear'. We also added 'neutral' for emotion list because the possibility if the emotion is difficult to detect. Below are the definitions of the label that used.

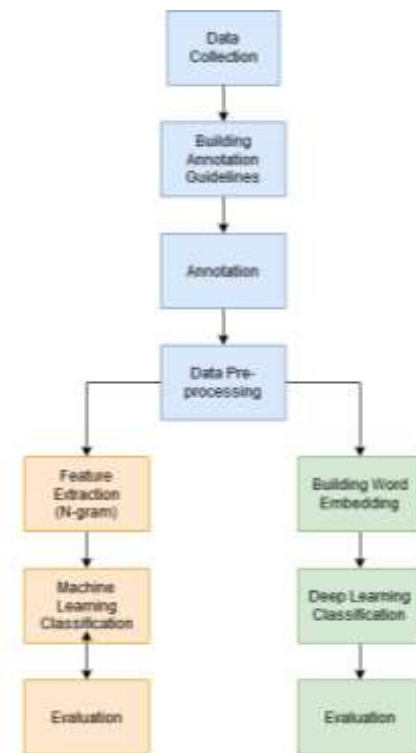


Fig. 1. Research Steps.

<sup>5</sup> <http://alt.qcri.org/semeval2015/>

<sup>6</sup> <https://www.youtube.com>

<sup>7</sup> <https://www.facebook.com/>

### 1) Sentiment labels:

a) *Positive*: Positive value can be seen by the appearance of positive terms, such as: “delicious”, “recommended”, “cheap”, “clean”, “friendly”, etc.

b) *Negative*: Negative label is given if the negative terms occur, for instance: “bad”, “horrible”, “not recommended”, “pricey”, “expensive”, “dirty”, etc.

c) *Neutral*: A review is classified as neutral if the terms that appear do not show positive or negative values. Besides, it can be noticed by the appearance of neutral terms, such as: “standard”, “so so”, “not bad but not good”, etc. In addition, the neutral label also given to the aspect that does not appear, because we assumed if an aspect does not mentioned, that means the polarity will be neither positive nor negative.

### 2) Emotion labels:

a) *Happy*: Happy emotion can be noticed by the appearance of phrases or words like: ‘I like it’, ‘really good’, ‘happy’, ‘satisfied’, ‘cool’, ‘worth’, ‘fun’, or emoticon ‘:’’, ‘:D’, etc.

b) *Sad*: Sad emotion shows the sadness or dissatisfied, and can be known by the appearance of terms ‘sad’, ‘dissatisfied’, ‘below expectation’, or with emoticon “:(”’, “:’”’.

c) *Surprised*: Surprised can be noticed by the terms like ‘I’m surprised’, ‘beyond expectation’, ‘shock’, etc.

d) *Angry*: Few terms that can be considered to label data as angry are ‘damn’, ‘angry’, ‘annoyed’, ‘annoying’, etc.

e) *Disgusted*: Disgusted emotion can be classified by the appearance of terms ‘dirty’, ‘disgusted’, etc.

f) *Fear*: Review is classified as fear if the terms like ‘afraid’, ‘worried’, etc, appears.

g) *Neutral*: Neutral label is given if the emotion in a review difficult to be interpreted. In addition, neutral emotion also will be given even though the aspects are not mentioned, like neutral definition in sentiment.

### C. Annotation

The next step is annotating the data. The annotation step consists two stages, which are sentiment annotation and emotion annotation. The method for deciding the annotator is crowdsourcing method, following a study from Sabou et al. [16]. The annotators are not linguistic experts. Besides, every review is annotated by 3 people in every stage. The method for retrieving the final label is major voting. After sentiment annotation, there are 562 data that cannot be used because the major voting results indicated that every annotator has labelled them with different labels. So, the data for the next annotation stage are 18346 reviews. However, because the limited time and number of annotators, the data that annotated for emotion label are only 15046 reviews. After applied major voting, the results of data that used are 14188. But the number of data with ‘angry’, ‘fear’, and ‘disgusted’ labels are very small, so we decided to remove those data, and the final number of data that we used for classification are 14103 reviews. Then, the labels that used are ‘positive’, ‘negative’, and ‘neutral’ for sentiment, while ‘happy’, ‘sad’, ‘surprised’, and ‘neutral’ for emotion.

### D. Data Preprocessing

After the annotation process, the next stage is data preprocessing. This stage adapted the research from [17] and consists few steps, which are:

1) *Emoticon Processing*: In this step, emoticon characters, such as :( was changed into ‘sad’, and :) into ‘happy’. This was conducted to avoid losing the information about the emoticon. Furthermore, when removing non alphabetical characters step is applied, the emoticon is not removed.

2) *Case Folding*: All of strings were changed into lowercase format to match the structures. For example, ‘Food’ was converted into ‘foods’.

3) *Abbreviation and Spelling Correction part 1*: In this part, the word spelling was corrected into formal form. For illustration, ‘I’ve visted the place, that wasn’t too crowd’ was corrected into ‘i have visted the place, that was not too crowd’. We used the abbreviation dictionary that is self-developed by [17], and contains abbreviations from Indonesian and English.

4) *Removing Non-alphabetical Characters*: After normalizing the words, then the non alphabetical characters, such as ‘.’, ‘!’, ‘@’, etc, are removed in this step.

5) *Abbreviation and Spelling Correction part 2*: In this step, the words are checked again whether all of them have been corrected. This step was applied to avoid the words that has the possibilities haven't been corrected in the third step. For instance, the phrase ‘tmptnya ga bgs!!’ was changed into ‘tempatnya tidak bgs!!’ after third step, but the word ‘bgs’ does not change into ‘bagus’ (good) because there are exclamation marks ‘!!’ that attached after words ‘bgs’. So, after the exclamation marks were removed in the fourth step, the phrase ‘tempatnya tidak bgs’, was corrected again into ‘tempatnya tidak bagus’ (the place was not good).

6) *Removing Stopwords*: In this stage, the stopwords that occur, like ‘i’, ‘you’, ‘always’, were removed. This step used dictionary built by [17] by combining NLTK<sup>8</sup> for English and Sastrawi<sup>9</sup> for Indonesian.

7) *Removing Repetitive Characters*: Sometimes, people like to express their feeling by using many unnecessary duplicated characters. These characters should be removed, and to illustrate this step, ‘happpppyyy’ is changed into ‘happy’.

8) *Stemming*: In this last preprocessing step, we removed the affixes and suffixes from the words to make them back into their base form. The functions that implemented are Snowball Stemmer by NLTK for English, and Sastrawi Stemmer for Indonesian because the data are in Indonesian and English, so, we applied two stemmers.

<sup>8</sup> <https://www.nltk.org/>

<sup>9</sup> <https://github.com/har07/PySastrawi>



E. Feature Extraction

This part explains about the feature extractions for machine learning, and the development of word embedding for deep learning.

1) *N-gram*: The features that used for classification using machine learning is n-gram language model word level. The number of gram that extracted as features are unigram, bigram, and the combination of unigram-bigram. We also applied chi-square method for feature selection.

2) *Word embedding*: For deep learning, we built our own word embedding using all scraped data from PergiKuliner. The method that implemented to build word embedding is skip-gram with dimension = 300.

IV. RESULTS AND ANALYSIS

This part explains about the experiments, results, and analysis of this research.

A. Experiments

In this study, we utilized the dataset that we made and created two scenarios for multi-label classification. Then, we compared several algorithms from machine learning and deep learning. After that, we evaluated the performances of those algorithms by comparing their F1 scores.

1) *Data*: This experiment using all data that are retrieved from annotation step. The total of data are 14103 reviews with three sentiment labels and four emotion labels. The distribution of labels for sentiment and emotion can be seen at Fig. 2 and Fig. 3, respectively. By seeing both figures, we noticed that the data have imbalanced labels for both sentiment and emotions. To illustrate, ‘food’ aspect is dominated by ‘positive’ sentiment and ‘happy’ emotion. On the other hand, all aspects beside ‘food’ is dominated by ‘neutral’ for both sentiment and emotion.

2) *Scenarios*:

a) *First scenario*: In first scenario, we employed problem transformation methods for multi-label classification in machine learning. Transformation methods that we implemented are Binary Relevance (BR), Label Powerset (LP), and Classifier Chain (CC). For machine learning algorithms, we applied are Decision Tree (DT), Random Forest (RF), Support Vector Machine (SVM), and Extra Tree Classifier (ET). The features that we used are unigram, bigram, and combination of unigram-bigram.

b) *Second scenario*: In this scenario, the deep learning algorithms that utilized are Bidirectional Long Short-Term Memory (BiLSTM), and Gated Recurrent Unit (GRU). We do not use problem transformations method like machine learning, but we assigned sigmoid as the activation function and binary cross entropy as loss function for retrieving the labels of data. The word embedding that has developed before is employed in this scenario.

3) *Evaluation*: Evaluation for both machine learning and deep learning is using kfold cross validation technique, with the number of k = 10. The scores that evaluated is f1-scores.

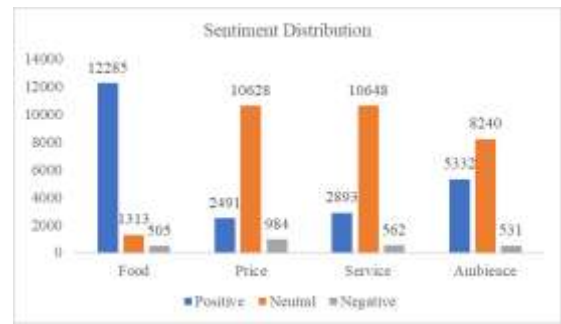


Fig. 2. Distribution of Sentiment Labels.

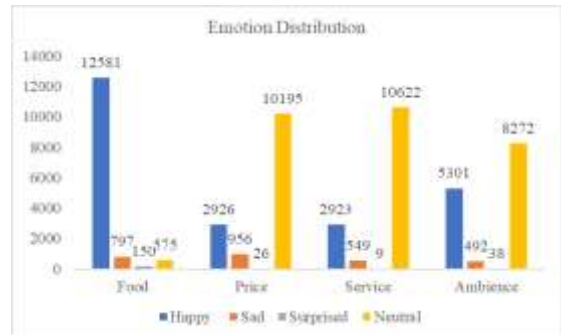


Fig. 3. Distribution of Emotion Labels.

B. Results

This section shows the performance of machine learning in first scenario, and deep learning in second scenario in every aspect of review. After that we assessed every performance in both scenarios by comparing their f1-scores.

1) *First scenario*:

a) *Label powerset*: This part presents the performance of machine learning algorithms when classified using Label Powerset (LP) as transformation method.

From Table I, it shows that ET achieved highest score, 88.17% for unigram feature. While for bigram, the highest score was acquired by RF with 87.3% for f1-score. This score was higher 0.61% compared to SVM score as second place. In the other hand, RF and ET claimed same f1 scores for unigram-bigram, which is 88.16%. By seeing the scores, it can be concluded that the best feature in this classification results is unigram.

Table II shows the performance of RF that dominated every feature in price aspect. However, for unigram-bigram feature, ET obtained same f1-score with RF, which is 89.54%. For the best feature in classification for price aspect, unigram-bigram achieved highest score compared to other two features.

TABLE I. CLASSIFICATION RESULTS FOR FOOD ASPECT

|     | Unigram       | Bigram        | Unigram-bigram |
|-----|---------------|---------------|----------------|
| SVM | 85.04%        | 86.69%        | 86.64%         |
| DT  | 82.21%        | 81.98%        | 83.40%         |
| RF  | 88.16%        | <b>87.30%</b> | <b>88.16%</b>  |
| ET  | <b>88.17%</b> | 86.10%        | <b>88.16%</b>  |

TABLE II. CLASSIFICATION RESULTS FOR PRICE ASPECT

|     | Unigram       | Bigram        | Unigram-bigram |
|-----|---------------|---------------|----------------|
| SVM | 83.16%        | 84.71%        | 85.24%         |
| DT  | 84.96%        | 83.80%        | 86.21%         |
| RF  | <b>87.17%</b> | <b>87.53%</b> | <b>89.54%</b>  |
| ET  | 86.84%        | 86.64%        | <b>89.54%</b>  |

For service aspect, Table III shows ET monopolized the scores for both unigram and unigram-bigram features. While for bigram, the highest score was led by RF with 90.88%, 0.21% higher than ET. However, the best feature for this classification in service aspect is unigram-bigram with score is 92.65% obtained by ET.

Similar to previous table, Table IV shows ET achieved highest scores for both unigram, and unigram-bigram when classifying ‘ambience’ aspect. Also, RF obtained highest score for bigram feature with 81.82%. Then, same with service aspect, in this classification results, the best feature is unigram-bigram with score is 86.98% that achieved by ET.

From Table V, we can see the highest scores in every aspect and in every feature that implemented. By seeing the table, it presents that with Label Powerset (LP), ‘food’ aspect was the only one that has highest score when it was classified using unigram with score 88.17%, while other aspects got their best performances when they were classified with unigram-bigram. Besides, ET obtained highest scores in every aspect except ‘price’ which its highest score achieved by RF. For bigram feature, all aspects were dominated by RF, but the scores are below unigram-bigram features. In addition, the aspect that has the highest score compared to other aspects is service that attained by ET with score is 92.65%.

TABLE III. CLASSIFICATION RESULTS FOR SERVICE ASPECT

|     | Unigram       | Bigram        | Unigram-bigram |
|-----|---------------|---------------|----------------|
| SVM | 88.55%        | 87.84%        | 89.80%         |
| DT  | 88.74%        | 87.89%        | 89.54%         |
| RF  | 90.64%        | <b>90.88%</b> | 91.84%         |
| ET  | <b>90.77%</b> | 90.67%        | <b>92.65%</b>  |

TABLE IV. CLASSIFICATION RESULTS FOR AMBIENCE ASPECT

|     | Unigram       | Bigram        | Unigram-bigram |
|-----|---------------|---------------|----------------|
| SVM | 81.61%        | 79.74%        | 83.13%         |
| DT  | 80.41%        | 75.53%        | 82.53%         |
| RF  | 85.82%        | <b>81.82%</b> | 86.48%         |
| ET  | <b>85.96%</b> | 81.48%        | <b>86.98%</b>  |

TABLE V. BEST PERFORMANCE OF EVERY ASPECT

|          | Unigram            | Bigram      | Unigram-bigram     |
|----------|--------------------|-------------|--------------------|
| Food     | <b>88.17% (ET)</b> | 87.30% (RF) | 88.16% (RF)        |
| Price    | 87.17% (RF)        | 87.52% (RF) | <b>89.54% (RF)</b> |
| Service  | 90.77% (ET)        | 90.88% (RF) | <b>92.65% (ET)</b> |
| Ambience | 85.96% (ET)        | 81.82% (RF) | <b>86.98% (ET)</b> |

b) Binary relevance: This part presents the performances of machine learning algorithms when classified using Binary Relevance (BR) as transformation method.

Table VI shows the performance of RF that attained highest scores in every feature in ‘food’ aspect. ET follows it by obtaining scores that not really far from RF scores. The table also shows that classification result using unigram-bigram feature is higher than other features, even though the score is only 0.01% higher than score that retrieved by using unigram feature only.

By seeing the Table VII, for the first time DT attained highest score comparing to other algorithms, with unigram feature. DT achieved 83.60%, followed by ET that got score which was 1.27% lower than DT. For bigram feature, RF achieved highest score when classifying ‘price’ aspect. However, unigram-bigram, once again, become the feature that helped ET to attain highest score for ‘price’ aspect with score 87.56%.

Similar to ‘price’ aspect results, Table VIII shows DT achieved highest score again for classifying ‘price’ aspect using unigram feature, but for this time, DT was followed by RF that was 0.39% lower than DT. RF also leads the score by classifying using bigram, and its score is 90.07%. Best feature for this aspect also obtained by unigram-bigram, with ET as classification algorithm. The score ET obtained was 91.28%, 1.21% higher compared to bigram and RF pair.

From Table IX, it can be seen that ET leads in both unigram and unigram-bigram features while classifying the ‘ambience’ aspect. While RF achieved best score when classifying using bigram feature with score id 80.12%. In addition, similar to three previous aspects, best classification score was obtained when using unigram-bigram feature by ET.

From the comparison of all machine learning algorithms that shown in Table X, we can see all best performances were attained by using unigram-bigram as feature. By applying BR method, and unigram-bigram as feature, ET successfully obtained highest scores in three aspects, which are ‘price’, ‘service’, and ‘ambience’. In other hand, RF dominates all ‘food’ aspect scores by using all features, including unigram-bigram.

TABLE VI. CLASSIFICATION RESULTS FOR FOOD ASPECT

|     | Unigram       | Bigram        | Unigram-bigram |
|-----|---------------|---------------|----------------|
| SVM | 83.24%        | 85.53%        | 84.89%         |
| DT  | 81.21%        | 80.67%        | 82.24%         |
| RF  | <b>88.17%</b> | <b>86.78%</b> | <b>88.18%</b>  |
| ET  | 88.10%        | 85.55%        | 88.04%         |

TABLE VII. CLASSIFICATION RESULTS FOR PRICE ASPECT

|     | Unigram       | Bigram        | Unigram-bigram |
|-----|---------------|---------------|----------------|
| SVM | 80.97%        | 84.16%        | 84.03%         |
| DT  | <b>83.60%</b> | 81.86%        | 85.12%         |
| RF  | 81.88%        | <b>86.25%</b> | 87.05%         |
| ET  | 82.33%        | 85.42%        | <b>87.56%</b>  |

TABLE VIII. CLASSIFICATION RESULTS FOR SERVICE ASPECT

|     | Unigram       | Bigram        | Unigram-bigram |
|-----|---------------|---------------|----------------|
| SVM | 86.34%        | 87.27%        | 87.83%         |
| DT  | <b>88.08%</b> | 86.43%        | 88.85%         |
| RF  | 87.69%        | <b>90.07%</b> | 90.45%         |
| ET  | 87.61%        | 89.82%        | <b>91.28%</b>  |

TABLE IX. CLASSIFICATION RESULTS FOR AMBIENCE ASPECT

|     | Unigram       | Bigram        | Unigram-bigram |
|-----|---------------|---------------|----------------|
| SVM | 79.66%        | 79.29%        | 82.51%         |
| DT  | 79.43%        | 73.87%        | 80.83%         |
| RF  | 83.75%        | <b>80.12%</b> | 84.72%         |
| ET  | <b>83.85%</b> | 79.48%        | <b>85.51%</b>  |

TABLE X. BEST PERFORMANCE OF EVERY ASPECT

|          | Unigram     | Bigram      | Unigram-bigram     |
|----------|-------------|-------------|--------------------|
| Food     | 88.17% (RF) | 86.78% (RF) | <b>88.18% (RF)</b> |
| Price    | 83.60% (DT) | 86.25% (RF) | <b>87.56% (ET)</b> |
| Service  | 88.08% (DT) | 90.07% (RF) | <b>91.28% (ET)</b> |
| Ambience | 83.85% (ET) | 80.12% (RF) | <b>85.51% (ET)</b> |

Furthermore, like LP, ‘service’ becomes the aspect that got highest score in Table X, which is 91.28%, compared to other aspects. Then it is followed by ‘food’, then ‘price’, and ‘ambience’ aspect, respectively.

c) *Classifier chain*: This part shows the performances of machine learning algorithms when classified using Classifier Chain (CC) as transformation method.

In Table XI, the classification results of ‘food’ aspect were dominated by RF in every feature that was used. However, in unigram-bigram feature, ET successfully gained same score with RF, which is 88.40%. Moreover, similar to LP and BR methods, by using CC, unigram-bigram still becomes the best feature of multi-label classification for ‘food’ aspect, following by unigram.

For classification of ‘price’ aspect, Table XII shows that RF attained best score in unigram, and also bigram feature. While for unigram-bigram feature, ET obtained the highest score with 89.24%, 1.62% and 3.93% higher compared to results from RF with bigram and unigram, respectively. This also means that once again, unigram-bigram is the best feature for classifying the ‘price’ aspect, similar to previous aspect.

Table XIII presents the performances of algorithms for classifying ‘service’ aspect. We can see that ET leads the score for classification using unigram and unigram-bigram, while RF achieved highest score for bigram. However, unigram-bigram still becomes the best feature for this aspect while it was classified using ET, and the f1-score is 92.09%.

Identical to previous aspect, as shown by Table XIV, ET obtained highest score for ‘ambience’ aspect in both unigram and unigram-bigram features. Best score in bigram also obtained by RF with 81.84%. Despite of it, it is still 5.26% lower than score attained by ET with unigram-bigram feature.

Again, unigram-bigram becomes the best feature for ‘ambience’ aspect.

In Table XV, it can be noticed that unigram-bigram becomes the best feature when Classifier Chain (CC) transformation method was applied. Unigram-bigram dominates all aspects, like Binary Relevance (BR). Besides, ET also attained the highest scores almost in all aspects, except ‘food’ aspect that was dominated by RF, also same with BR.

In addition, like both LP and BR results, the best score between all aspects was obtained by ‘service’ aspect when it was classified by ET using unigram-bigram. The score that ET achieved for ‘service’ aspect is 92.09%, 2.85% higher than ‘price’ aspect which was the second highest after ‘service’ aspect.

TABLE XI. CLASSIFICATION RESULTS FOR FOOD ASPECT

|     | Unigram       | Bigram        | Unigram-bigram |
|-----|---------------|---------------|----------------|
| SVM | 83.44%        | 85.74%        | 85.01%         |
| DT  | 82.34%        | 81.54%        | 83.15%         |
| RF  | <b>88.20%</b> | <b>87.21%</b> | <b>88.40%</b>  |
| ET  | 88.17%        | 86.10%        | <b>88.40%</b>  |

TABLE XII. CLASSIFICATION RESULTS FOR PRICE ASPECT

|     | Unigram       | Bigram        | Unigram-bigram |
|-----|---------------|---------------|----------------|
| SVM | 81.15%        | 84.16%        | 84.51%         |
| DT  | 84.10%        | 83.16%        | 85.55%         |
| RF  | <b>85.31%</b> | <b>87.62%</b> | 88.74%         |
| ET  | 84.93%        | 86.84%        | <b>89.24%</b>  |

TABLE XIII. CLASSIFICATION RESULTS FOR SERVICE ASPECT

|     | Unigram       | Bigram        | Unigram-bigram |
|-----|---------------|---------------|----------------|
| SVM | 86.33%        | 87.21%        | 88.44%         |
| DT  | 88.20%        | 87.05%        | 89.22%         |
| RF  | 88.13%        | <b>90.57%</b> | 91.02%         |
| ET  | <b>88.54%</b> | 90.39%        | <b>92.09%</b>  |

TABLE XIV. CLASSIFICATION RESULTS FOR AMBIENCE ASPECT

|     | Unigram       | Bigram        | Unigram-bigram |
|-----|---------------|---------------|----------------|
| SVM | 79.89%        | 79.84%        | 82.80%         |
| DT  | 79.81%        | 74.77%        | 80.78%         |
| RF  | 85.61%        | <b>81.84%</b> | 86.32%         |
| ET  | <b>85.74%</b> | 81.37%        | <b>87.10%</b>  |

TABLE XV. BEST PERFORMANCE OF EVERY ASPECT

|          | Unigram     | Bigram      | Unigram-bigram     |
|----------|-------------|-------------|--------------------|
| Food     | 88.20% (RF) | 87.21% (RF) | <b>88.40% (RF)</b> |
| Price    | 85.31% (RF) | 87.62% (RF) | <b>89.24% (ET)</b> |
| Service  | 88.54% (ET) | 90.57% (RF) | <b>92.09% (ET)</b> |
| Ambience | 85.74% (ET) | 81.84% (RF) | <b>87.10% (ET)</b> |

Table XVI shows the comparison of best performances from Binary Relevance (BR), Label Powerset (LP), and Classifier Chain (CC) with unigram-bigram as feature. We can see from the table that ‘food’ aspect got the highest score when it was classified by RF with CC as problem transformation method. Followed by BR, and LP, respectively. For ‘price’ and ‘service’ aspects, LP is better than other transformation methods when classifying both aspects, followed by CC, then BR. For ‘price’ aspect, the algorithm that obtained the highest score, which is 89.54%, was RF. While for ‘service’ aspect, the best score was achieved by ET with 92.65%. However, in case of ‘ambience’ aspect, ET attained the highest score with CC as transformation method for multi-label classification. The score that was achieved by ET in ‘ambience’ aspect is 87.1%, 0.12% higher than the score it obtained by using LP as problem transformation method.

Furthermore, it also can be noticed that BR cannot surpass both LP and CC, except in ‘food’ aspect where BR score is 0.02% higher than LP. This maybe happened because as transformation method, BR treats the labels independently before they are classified by machine learning. This means, BR does not consider the relationship between the labels. For instance, the sentiment label ‘positive’ is considered does not have relation with the emotion label ‘happy’, because both labels were classified separately. In the other hand, LP transforms the label combinations into new classes before machine learning classified them as multiclass problem. While CC transforms the labels by using the first label that obtained from first classification as a feature for classifying the next label in next classification. Thus, by seeing the way the three transformation methods work, we can conclude that LP and CC consider the relation between labels, while BR does not consider it.

Moreover, both ET and RF always obtain best score than DT and SVM in all aspects inn all transformation methods that were used in this research. It should be remembered that both ET and RF are tree-based ensemble algorithms, which means the way they work is almost similar, except the way they split the nodes and use the samples. However, by seeing Table XVI, we can see that ET dominates ‘price’, ‘service’, and ‘ambience’ aspects for all transformation methods, except for LP in ‘price’ aspect which its best score was obtained by RF. For ‘food’ aspect, all highest scores for all transformation method were attained by RF.

2) *Second scenario*: This part shows the performances of deep learning algorithms, which are BiLSTM and GRU.

TABLE XVI. COMPARISON OF PERFORMACES FROM LABEL POWERSSET (LP), BINARY RELEVANCE (BR), CLASSIFIER CHAIN (CC) WITH UNIGRAM-BIGRAM

|          | LP                 | BR          | CC                |
|----------|--------------------|-------------|-------------------|
| Food     | 88.16% (RF)        | 88.18% (RF) | <b>88.4% (RF)</b> |
| Price    | <b>89.54% (RF)</b> | 87.56% (ET) | 89.24% (ET)       |
| Service  | <b>92.65% (ET)</b> | 91.28% (ET) | 92.09% (ET)       |
| Ambience | 86.98% (ET)        | 85.51% (ET) | <b>87.1% (ET)</b> |

From the classification results of both deep learning algorithms, Table XVII shows that GRU and BiLSTM attained same scores for ‘food’ aspect. However, BiLSTM leads the scores for ‘service’ and ‘ambience’ aspects. For GRU, it obtained higher score compared to BiLSTM in ‘price’ aspect, which its score peaks on 83.01%, 0.92% higher than BiLSTM. Nonetheless, the scores from GRU in ‘service’ and ‘ambience’ are not very far from BiLSTM scores. The scores achieved by GRU are 0.33% and 0.86% lower than BiLSTM scores in ‘service’ and ‘ambience’ aspects, respectively. From this experiment, it can be concluded that GRU can compete with performances from BiLSTM, even though BiLSTM already uses future context that can help it to solve more complex classification problems.

TABLE XVII. COMPARISON OF PERFORMANCES FROM GRU AND BiLSTM

|          | GRU           | BiLSTM        |
|----------|---------------|---------------|
| Food     | <b>88.16%</b> | <b>88.16%</b> |
| Price    | <b>83.01%</b> | 82.09%        |
| Service  | 88.70%        | <b>89.03%</b> |
| Ambience | 83.92%        | <b>84.78%</b> |

In addition, like machine learning, ‘service’ aspect becomes the aspect that gotten highest score when it was classified by BiLSTM and GRU. Then, the aspect that becomes the second highest is ‘food’, followed by ‘ambience’ and ‘price’, respectively. Furthermore, it can be concluded that self-developed word embedding can work well with deep learning. Hence, the scores that obtained by deep learning algorithms are quite similar to machine learning.

### C. Analysis

Overall, the results of both scenarios show that ‘service’ aspect becomes the aspect that can be classified better than other aspects. After ‘service’ aspect, it was followed by ‘price’, ‘food’, and ‘ambience’, respectively, for machine learning. For deep learning, the second highest score was obtained when algorithms classified ‘food’, followed by ‘ambience’, then ‘price’ aspect, respectively. This may be affected by the way people express their comments towards the aspects. Usually, whenever people comment about ‘service’ aspect, people tend to use words like ‘service’ or ‘waitress’ directly in the comments, same goes with ‘price’ aspects. This kind of writing is different when people talk about ‘food’ and ‘ambience’ aspect, which can be written more creative by customers. To illustrate, people often write all the food names they ordered, and explain them in detail one by one. This can lead to misclassification by the classification program if there is a conflict occurs in an ‘aspect’. For example, the comment ‘the noodles were very good but too oily, I don’t like it’, or ‘the fried rice was delicious but the orange juice too blend’. Those kinds of reviews can create a conflict and affects the classification results. Same goes with ‘ambience’ aspect, people can explain it variatively. For instance, ‘it has beautiful decoration, but the room was full of smoke’.

For second scenario results, ‘price’ aspect become the aspect with lowest score after classification. While ‘food’ and

'ambience' aspects become two and third place after 'service' aspect that has higher score. This may be caused by the label distribution in dataset, which 'positive' sentiment and 'happy' emotions are dominant in 'food' aspect, followed by 'ambience', 'service', and 'ambience' aspect. Thus, the deep learning models learned 'positive' and 'happy' labels well, compared to other labels.

Furthermore, the features that used also affect the classification results. In first scenario, unigram-bigram feature gave more information compared to apply only unigram, or only bigram independently. When classifying, unigram can work well because in unigram, words are treated individually, and those words often appear in the dataset. To illustrate, the sentence 'I like the food but it was too pricey'. In unigram, it will be 'I', 'like', 'the', 'food', 'but', 'it', 'was', 'too', 'pricey', and for bigram, it will be 'I like', 'like the', 'the food', 'food but', 'but it', 'it was', 'was too', 'too pricey'. When classifying using bigram, the models work well but not always good compared to unigram because the combination of words in bigram are not often appear in reviews compared to unigram. Thus, if unigram and bigram are combined, the models obtain more information about word when they appear individually and when they appear as pairs. Then, for second scenario, classification with self-developed word embedding can give good results with the information especially information about semantic relations between words. Hence, it should be considered to add other features, such as POS tagging, for machine learning and deep learning to enhance their performances.

Label distribution also contributes to affect the classification results. This research has imbalanced dataset, so, it will be good to use data augmentation or apply oversampling/undersampling methods to balance the data.

## V. CONCLUSIONS AND FUTURE WORK

For this research, we made experiments and evaluated the performances of machine learning algorithms, which are Decision Tree (DT), Random Forest (RF), Support Vector Machine (SVM), and Extra Tree Classifier (ET), as well as deep learning, (Bidirectional LSTM (BiLSTM), and Gated Recurrent Unit (GRU)). We made two scenarios, which in first scenario, we applied transformation methods such as Binary Relevance (BR), Label Powerset (LP), and Classifier Chain (CC) for multi-label classification in machine learning. Then the features that used are unigram, bigram, and combination of unigram-bigram. For second scenario, we utilized sigmoid as the activation function and binary cross entropy as loss function for retrieving the labels of data in deep learning. Then, self-developed word embedding is employed in this scenario for deep learning classification. The results show RF dominates with 88.4% and 89.54% F1 scores with CC method for food aspect, and LP for price, respectively. For service and ambience aspects, ET leads with 92.65% and 87.1% with LP and CC methods, respectively. On the other hand, in deep learning comparison, GRU and BiLSTM obtained similar F1-score for food aspect, 88.16%. On price aspect, GRU leads with 83.01%. However, for service and ambience, BiLSTM achieved higher F1-score, 89.03% and 84.78%.

Since the distribution of label in our data is imbalanced, for the future, it should be considered to use balancing methods such as oversampling or undersampling. We also can apply data augmentation to retrieve new data for labels that have small numbers. Besides, we need to add more features to enhance the performance of both machine learning and deep learning.

## ACKNOWLEDGMENT

The authors would like to thank the PUTI research team for the supports and helps that provided during this research.

## REFERENCES

- [1] M. Pontiki and J. Pavlopoulos, "SemEval-2014 Task 4: Aspect Based Sentiment Analysis," no. SemEval, pp. 27–35, 2014.
- [2] M. Pontiki et al., "SemEval-2016 Task 5: Aspect Based Sentiment Analysis," pp. 19–30, 2016.
- [3] S. M. Mohammad, S. Kiritchenko, and X. Zhu, "NRC-Canada: Building the state-of-the-art in sentiment analysis of tweets," \*SEM 2013 - 2nd Jt. Conf. Lex. Comput. Semant., vol. 2, no. SemEval, pp. 321–327, 2013.
- [4] D. Stojanovski, G. Strezoski, G. Madjarov, I. Dimitrovski, and I. Chorbev, "Deep neural network architecture for sentiment analysis and emotion identification of Twitter messages," *Multimed. Tools Appl.*, vol. 77, no. 24, pp. 32213–32242, 2018.
- [5] S. M. Mohammad, "Sentiment Analysis: Detecting Valence, Emotions, and Other Affectual States from Text," *Emot. Meas.*, pp. 201–237, 2016.
- [6] M. Hassan, M. S. Bin Alam, and T. Ahsan, "Emotion Detection from Text Using Skip-thought Vectors," in 2018 International Conference on Innovations in Science, Engineering and Technology, ICISSET 2018, 2018, pp. 501–506.
- [7] H. S. Khawaja, M. O. Beg, and S. Qamar, "Domain specific emotion lexicon expansion," in 2018 14th International Conference on Emerging Technologies, ICET 2018, 2019, pp. 1–5.
- [8] L. Wikarsa and S. N. Thahir, "A text mining application of emotion classifications of Twitter's users using Naïve Bayes method," in Proceeding of 2015 1st International Conference on Wireless and Telematics, ICWT 2015, 2016, pp. 1–6.
- [9] J. Savigny and A. Purwarianti, "Emotion classification on youtube comments using word embedding," *Proc. - 2017 Int. Conf. Adv. Informatics Concepts, Theory Appl. ICAICTA 2017*, pp. 1–5, 2017.
- [10] A. Cahyadi and M. L. Khodra, "Aspect-Based Sentiment Analysis Using Convolutional Neural Network and Bidirectional Long Short-Term Memory," 2018 5th Int. Conf. Adv. Informatics Concept Theory Appl., pp. 124–129, 2018.
- [11] D. Ekawati, "Aspect-based Sentiment Analysis for Indonesian Restaurant Reviews," 2017.
- [12] K. Shalini, B. Ganesh, A. K. M, and K. P. Soman, "Sentiment Analysis for Code-Mixed Indian Social Media Text With Distributed Representation," 2018 Int. Conf. Adv. Comput. Commun. Informatics, pp. 1126–1131, 2018.
- [13] S. Y. M. Lee and Z. Wang, "Multi-view learning for emotion detection in code-switching texts," in Proceedings of 2015 International Conference on Asian Language Processing, IALP 2015, 2016, pp. 90–93.
- [14] M. Pontiki, D. Galanis, and H. Papageorgiou, "SemEval-2015 Task 12: Aspect Based Sentiment Analysis," 2015.
- [15] P. Ekman, W. V. Friesen, and P. Ellsworth, "What Emotion Categories or Dimensions Can Observers Judge from Facial Behavior? In Emotion in the Human Face," no. Cambridge University Press, pp. 98–110, 1982.
- [16] M. Sabou, K. Bontcheva, L. Derczynski, and A. Scharl, "Corpus annotation through crowdsourcing: Towards best practice guidelines," *Proc. 9th Int. Conf. Lang. Resour. Eval. Lr. 2014*, no. 2010, pp. 859–866, 2014.
- [17] A. Suciati and I. Budi, "Aspect-based Opinion Mining for Code-Mixed Restaurant Reviews in Indonesia," in 2019 International Conference on Asian Language Processing (IALP), 2019, pp. 59–64.

# Pseudo Amino Acid Feature-Based Protein Function Prediction using Support Vector Machine and K-Nearest Neighbors

Anjna Jayant Deen<sup>1</sup>, Manasi Gyanchandani<sup>2</sup>

Department of Computer Science and Engineering  
Maulana Azad National Institute of Technology, Bhopal, India

**Abstract**—Bioinformatics facing the vital challenge in protein function prediction due to protein data are available in primary structure, an amino acid sequence. Every protein cell sequence length and size are in different sequence order. Protein is available in 20 amino acid sequence alphabetic order; however, the corresponding information of the membrane protein sequence is insufficient to capture the function and structures of a protein from primary sequence datasets. A challenging task to correctly identify protein structure and function from amino acid sequence. The basic principle of PseAAC (Pseudo Amino Acid Composition) is to generate a discrete number of every protein samples. In each protein, sequence length varies due to protein functions. Some protein sequence length is less than 50, and some are large. Due to this, different sizes of the amino acid sample are chances to lose sequence order information. PseAAC feature generates a fixed size descriptor value in vector space to overcome sequence information loss and is used to further systematic evolution. Therefore machine learning computational tool synthesizes accurate identification of structure and function class of membrane protein. In this study, SVM (Support Vector Machine) and KNN (K-nearest neighbors) based prediction classifier used to identifying membrane protein and their types.

**Keywords**—Membrane protein types; classifiers; SVM (RBF); KNN; Random Forest; PseAAC

## I. INTRODUCTION

Bioinformatics is a different field of combination to solve biological problems with computational techniques dealing remarkably in extensive scale information of system biology. The amino acid residue is a part of macromolecules. The membrane protein is the type of protein residing around a cell membrane, or their subcellular locations are also defined as various types of protein. The most genome encodes a membrane protein, during the encoding process to finding genes cell membrane function perform a wide range of synthesis. Membrane cell identification hard due to lack of stability found a more flexible and hydrophobic surface part [8],[9]. However, protein structure finding is still challenging. Learning outer and non-outer membrane cells necessary to develop computational tools for new drug design and genome sequencings [10][26]. The Protein cell determines its energy and several functions; every cell component depends on protein molecules functions, the cell responsible for signalling cell system, and mobilize an intracellular response. The cell membrane of functional and structural properties targets to find disease, drug design, and novel research [15]. Membrane

cell misfold work as monitors, changing their shape and action in response to metabolic signals or information from outside the cell causes various disease like Alzheimer's disease, cardiovascular diseases, neurological disorders, and cancer [1],[20],[21]. Membrane proteins sequence combination, 45-55%, is used in protein legend docking and drug design [10]. Membrane functions are the essential element to discover new drugs and genomes [9-10]. Now capturing the features of membrane functions is responsible for the distribution of cell systems and their role. Conventional techniques used in biological experimental to predict the membrane types are costly and tedious [19]. However, a fast, automated, and effective method must be needed to identify unknown protein types. Analysis of the membrane proteins is hard, and most of them will not dissolve in ordinary solvents. Hence, very few structures of membrane protein have been found so far. Many authors were reports [11],[18],[21],[22] have shown NMR to be a powerful instrument for the detection of membrane protein structures; it is expensive and time processing. Therefore demand to construct computational methods that can predict membrane protein characteristics based on their primary sequence would be very helpful. The membrane protein is classified into mainly transmembrane or anchored protein attached to inside and outside the cell. These membrane cells are further classified into eight subtypes: Type-1, Type-2, Type-3, Type-4, Multi-pass are transmembrane protein and were Peripheral, Lipid chain, and GIP are anchored protein [2], [5]. Currently, different kinds of feature extractions and classification methods have been built to be used to predict membrane types. ACC (Amino acid composition) is used in predicting membrane protein types [3],[5],[13], first used by the article [3], but sequence order of information can't store during implementing amino acid composition. Therefore Chou's suggests PseAAC (Pseudo Amino Acid Composition) evaluates the composition value of amino acid in fixed combinations and saves them. Further, many authors [1],[4],[5],[6],[18] have been processed and suggested different techniques to depict protein samples to overcome. PseAAC, feature extraction method, followed by many latest article [13],[18],[5]. Various computational methods based on learning classifiers and ensemble methods have been used for predicting cell membranes in high-performance accuracy. In this study, a novel feature-based machine learning technique to identify the membrane cell. The proposed objective model was to enhance the accuracy of the classification. Each sequence chain of protein features has

mapped into a vector space. And, the multiclass membrane to recognize, the better performing multiclass classifier was chosen. Patterns match and similarity were calculated by using the standard test conducted on high dimensional multiclass protein data.

## II. MATERIALS AND METHOD

### A. Data Sets

The protein data bank has manually annotated proteins collected from Swiss-Prot PDB [14],[16],[17]. In this study, 560459 protein was obtained from a form data source. Datasets are further preprocessed for identifying a non-membrane and membrane protein correctly [19]. Here, 62029 membrane proteins are captured. For finding its types which is in eight classes, are: (i) GPI-anchored, (ii) lipid chain-anchored, (iii) multipass transmembrane, (iv) peripheral, (v) Type-1, (vi) Type-2, (vii) Type-3 and, (viii) Type-4. Further classification of the 62029 membrane proteins data sequence split into 43418 training and 18611 test samples. Table I have shown the sample details [2].

TABLE I. TOTAL SAMPLES IN THE DATASET

| Membrane protein (types) | No. of instances |
|--------------------------|------------------|
| GPI-anchored             | 651              |
| Lipid chain anchored     | 3032             |
| Multipass transmembrane  | 35480            |
| peripheral               | 17319            |
| Type-1                   | 2948             |
| Type-2                   | 2194             |
| Type-3                   | 211              |
| Type-4                   | 194              |
| Total                    | 62029            |

### B. Feature Extraction Methods

Feature selection is the main part of the machine learning process [4]. Specific knowledge is useful for identifying membrane types. Without knowing the sequence order, a sequence's composition loses the information and not used further evaluation. PseAAC (pseudo amino acid composition) to prevent the protein sequence order and pattern data. [29]. PseAAC has to generate ordered 50-dimensional vector space for each sequence data to be involved in computational proteomics [20], and sequence length generate 1 dimensional vector space each samples. In [30] suggest that it is feasible to predict membrane protein type when the features are derived directly from the amino acid sequence. A python-based toolkit iFeature integrates and calculating an extracting feature encode into specific properties of amino acid for generating 51 numerical descriptor value.

## III. PROPOSED METHODOLOGY

A practical method develops for predicting the function and structure of protein class from its discrete dimensional vector value. For doing this, the main steps are followed by

step 1. Collect protein benchmark data, step 2. Establish a well-built prediction algorithm and step 3. Valuable intrinsic relates as an emphasis for the membrane data samples that can match their desire object to predict.

This study is focusing on the 3rd step, a necessity. In this regard, various methods for formulating protein samples. Therefore, they can categorize into various representations as to the discrete value for sequential description. The flow diagram is shown in Fig. 1. In this study, a significant improvement as an order to,

- SVM (RBF) to expand functional parameters reflect in high-dimensional membrane cell descriptors protein and,
- Enhanced predicting results merits the use of further enriched training data samples and identify different types of membrane cell descriptors.
- Integrated features of PseAAC and sequence length are used for analysis to evaluate membrane. The learning model is based on kernel SVM for functional prediction and similarity matches in sequence to a query membrane.
- Unique multiclass are in eight batches—each sample descriptor value is 51D space for supporting multiple membrane types.
- Machine learning classifiers K nearest neighbors and Random Forest was added for simplifying the collective samples via computation of protein functions by multiple types.

Cross-validation is one method to overcome the class imbalance problem. Therefore, in this study, we use k-fold cross-validation. Membrane protein data consisting of N tuple has divided in k=10 folds (D1, D2, D3...D10), and if the N tuple is not divisible by k, then the last part is considered as a (k-1). Here in our estimation, learning using 10-fold cross-validation. A sequence of k = 10 runs is carried out with the decomposition and  $i^{th}$  = iteration, and  $D_i$  use as test data and other fold as training data. Thus, each tuple uses the same amount of time for training samples, and once for testing. The overall average of each iteration is estimated.

### A. PseAAC (Pseudo Amino Acid Composition)

Membrane protein information senses its molecular action. Its process is in a molecules system that allows organisms to endure these basic life processes—various inherited diseases caused by mutations and changes observed in a protein sequence result. The amino acid sequence is a pattern of 20 unique amino acid residues. As per chemical composition, amino acid 20 sets are further categorized into four groups: polar, nonpolar, positive charged, and negatively charged [12]. These are comparable and a varying side chain. Each amino acid has distinct chemical properties due to the different groups' side chains. The 20 amino acids composition computed as in eq. (1), [3].

$$P = [p_1, p_2 \dots p_{20}, p_{20+1} \dots p_{20+\lambda}]^T \quad (1)$$

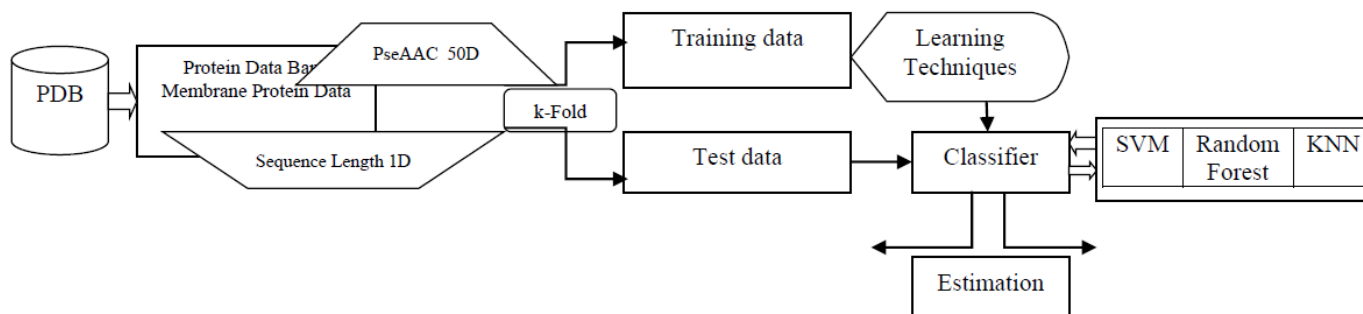


Fig. 1. Proposed Methodology.

The composition portions of amino acids are evaluated by using the mass of 20 amino acids, hydrophilic value and hydrophilic value.  $p_1, p_2, \dots, p_{20+\lambda}$  are calculated by eq. (2):

$$p_u = \begin{cases} \frac{f_u}{\sum_{i=1}^{20} f_i + \omega \sum_{k=1}^{\lambda} \tau_k}, & (1 \leq u \leq 20) \\ \frac{\omega \cdot \tau_{u-20}}{\sum_{i=1}^{20} f_i + \omega \sum_{k=1}^{\lambda} \tau_k}, & (20 + 1 \leq u \leq 20 + \lambda) \end{cases} \quad (2)$$

Integrating features descriptor value has produced a verity of amino acid patterns on regular occurrence in the protein sequence. The length of PseAAC depends on the descriptor value. This study uses a set of 30 amino acid composition. Thus the feature dimension from PseAAC is 50 (20+30=50D) descriptor vector space [2],[23].

### B. Sequence to Integer Encoding

This method is configured for a particular integer value (range from 1-20) with 20 amino acid residues made from the protein sequence. A protein sequence can be translated to an integer sequence by replacing each letter with a corresponding mapping integer value. The sum of residues in the sequence is proportional to its weight. For instance, in a protein sequence, AJKJLMLLK, L, is seen three times. The weight of L is then measured as  $3/9=0.33$ . Then, We use the following formula in eq. (3) to find the weight of a residue:

$$w_i = \frac{n_i}{L} \quad (3)$$

where,  $w_i$  is the weight of  $i^{\text{th}}$  residue,  $n_i$  is the number of occurrence of  $i^{\text{th}}$  residue in the protein sequence and  $L$  is the length of the protein sequence. A weighted total volume of each residue represents the required protein sequence and is performed by measuring each residue's weight. The numerical value encoded is then found as follows.

$$SEQ_{\text{encoded}} = k_1 \cdot w_1 + k_2 \cdot w_2 + \dots + k_{20} w_{20} \quad (4)$$

Where  $k_i$  is the  $i^{\text{th}}$  residue's mapped integer value and  $w_i$  is the corresponding weight of the residue got from equation (4). the resultant values gives one dimensional data for protein sequences. Where, weight factor is  $\omega$  (set to 0.05) and  $\tau_k$  is the  $k^{\text{th}}$  tier correlation factor that represents all correlation order of the  $k^{\text{th}}$ -most continuous residues.

### C. Classification Algorithm

Feature-based classification algorithm mapping the input data samples into the desired class, model build to predict class labels for unseen samples, the main part of machine learning applied in all fields are in bioinformatics and data

science. These were training techniques used to train most 70% data samples, and classifier testing test data sets 30%, respectively. This study classifier model based on SVM, KNN, and RF (Random Forest) used to classify the membrane features pseudo amino acid composition and sequence length descriptors into eight types.

1) *Support Vector Machine (SVM)*: In the bioinformatics data source, protein information has generally gathered in an amino acid sequence. However, a knowledge-based learning system dealing with homogeneous and heterogeneous datasets still needed some basic models based upon classification and clustering techniques. To implement the classifier support vector machine trendy and powerful for predicting protein structure and function. SVM (support vector machine) classification techniques have been used for dual-mode separation as a binary or multiclass. SVM outline draw hyperplane, which separates the decision surface data into two different class [31]. The use of the SVM learning model in high dimensional datasets creates a multiclass problem, so resolving that need to build a modified classification technique. In this study, unique features are based on a novel classifier model design on predicting membrane protein of achieving high accuracy for multiclass in the high-dimensional protein data source.

SVM transforms the given data first into a large vector space and then draws the maximum hyperplane margin to separate non-linear datasets, represented in Fig. 2. The functions for Radial Basic Function (RBF) in the SVM algorithm were used: Step 1. The built a feature vector from the input sequence. It can represent classes based on PseAAC and Sequence length properties. Step 2. RBF kernel selects to predict function while training eq. (5-16). Step 3 Selection of the prime parameter during training kernel function fit data to get maximum accuracy eq. (17-22)[2]. When the SVMs are using for the classification, the known set  $(\{+1, -1\})$ , marked training data is segregated by a hyperplane, which is as far as possible distant from positive negative samples. [See "optimal separating plane" (OSH) in Fig. 2]. The test data 'plot' then defines the positive or negative OSH for the high-dimensional sphere. The kernel model enables SVMs to work in combination with the nonlinear mapping into a function space to classify membrane protein types. For these problems, SVM is not linearly detachable. The SVM's optimal separating hyperplane within functional space is a nonlinear decision limit within the input space.



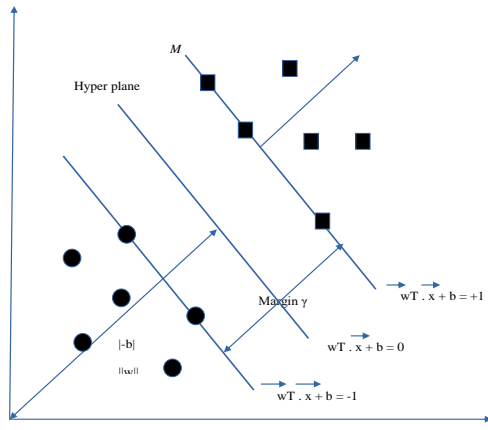


Fig. 2. Hyperplane and Margin Description. Samples of Class -1 and Class +1 are Represented Respectively by the Circular Dots and Square Dots.

### Linear SVM classification

where  $\vec{w} = (w_1, w_2, \dots, w_n)^T$  is a vector of  $n$  elements. Consider, the training data of two groups of  $n$  instance  $(\vec{x}_1, y_1), (\vec{x}_2, y_2), \dots, (\vec{x}_n, y_n)$ ,  $i = 1, 2, \dots, n$ , on each instances,  $y$  ( $i=1..n$ ),  $i = 1, 2, \dots, n$ , where specified a weight vector  $\vec{w}$  and bias  $b$ , and  $\vec{x}_i \in R^N$  is an  $N$  dimensional space, and  $y_i \in \{-1, +1\}$  is the class index.

$$\vec{w}^T \cdot \vec{x}_i + b \geq 1, y_i = +1, \quad (5)$$

$$\vec{w}^T \cdot \vec{x}_i + b \leq -1, y_i = -1, \quad (6)$$

The vector of  $n$  elements is where  $\vec{w} = (w_1, w_2, \dots, w_n)^T$ . Uniformities (1), (2) can be fused into one.

$$y_i(\vec{w}^T \cdot \vec{x}_i + b) \geq 1, \quad i = 1, 2, \dots, n. \quad (7)$$

For each training group, there are a number of hyperplanes. SVM's classification aims to create an optimum weight  $\vec{w}_0$  and an optimal bias  $b_0$  to achieve the maximum margin between the training data and the chosen hyperplane. The defined hyperplane by  $(\vec{w}_0)$  and  $b_0$  is optimal separating hyperplane. Any hyperplane can be represented as equated in eq. (8).

$$\vec{w}^T \cdot \vec{x}_i + b = 0 \quad (8)$$

and the difference between the two margins is in eq. (9)

$$\gamma(\vec{w}, b) = \frac{\min_{\{\vec{x}|y=+1\}} \vec{x}^T \cdot \vec{w}}{\|\vec{w}\|} - \frac{\max_{\{\vec{x}|y=-1\}} \vec{x}^T \cdot \vec{w}}{\|\vec{w}\|}. \quad (9)$$

The optimum separating hyperplane is being identified by raising the distance above or reducing the norm of  $\|\vec{w}\|$  by trying to restrict discrimination eq. (7), and

$$\gamma_{max} = \gamma(\vec{w}_0, b_0) = \frac{2}{\|\vec{w}_0\|}. \quad (10)$$

The following Lagrangean saddle point provides solutions to the above problems with optimization

$$L(\vec{w}, b, \alpha) = \frac{1}{2} \vec{w}^T \cdot \vec{w} - \sum_{i=1}^n \alpha_i [y_i(\vec{w}^T \cdot \vec{x}_i + b) - 1], \quad (11)$$

where  $\alpha \geq 0$  are Lagrange multipliers. To solve the quadratic programming problem, the gradient of  $L(\vec{w})$  to

$L(\vec{w}, b, \alpha)$  disappears in respect of  $\vec{w}$  and  $b$ , which gives a calculation of the following terms:

$$\frac{\delta L}{\delta \vec{w}} \Big|_{\vec{w}=\vec{w}_0} = 0, \text{ and } \frac{\delta L}{\delta b} \Big|_{\vec{w}=\vec{w}_0} = 0$$

$$\vec{w}_0 = \sum_{i=1}^n \alpha_i y_i \vec{x}_i, \quad (12)$$

$$\sum_{i=1}^n \alpha_i y_i = 0. \quad (13)$$

Via replacement of Eqs. (12, and 13) into (11), Maxing the following expression becomes the quadratic programming (QP) problem:

$$L(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j (\vec{x}_i^T \cdot \vec{x}_j) \quad (14)$$

in the constraints  $\sum \alpha_i y_i = 0$  and  $\alpha_i \geq 0, i = 1, 2, \dots, n$ .

Non-zero  $\alpha_i$  coefficients are among Eq. (14) solutions at the two optimal margins, and is known as vectors support (SV). The bias  $b_0$  can be estimated accordingly:

$$b_0 = -\frac{1}{2} \left( \min_{\{\vec{x}_i|y_i=+1\}} \vec{w}_0^T \cdot \vec{x}_i + \max_{\{\vec{x}_i|y_i=-1\}} \vec{w}_0^T \cdot \vec{x}_i \right) \quad (15)$$

The decision function that divides the two groups can be written as after evaluating the support vector and bias

$$f(\vec{x}) = \text{sign}[\sum_{i=1}^n \alpha_i y_i \vec{x}_i^T \cdot \vec{x} + b_0] = \text{sign}[\sum_{SV} \alpha_i y_i \vec{x}_i^T \cdot \vec{x} + b_0] \quad (16)$$

### Non-linear SVM classification

Since membrane protein types are typically nonlinear, these problems have been implemented in the SVM [30]. In the input space  $X$ , the original training data  $\vec{x}$  are translated into a high-dimensional F-function through the operator kernel Mercer  $K$  [34], in which the optimum separating hyperplane is formed. The set of classifiers will be converted into the form in mathematical terms.

$$f(\vec{x}) = \text{sign}[\sum_{i \in \{SV\}} \alpha_i y_i K(\vec{x}_i, \vec{x}) + b_0], \quad (17)$$

Where  $K$  is a symmetric positive function that fulfills the conditions of Mercer.

$$K(\vec{x}, \vec{y}) = \sum_{m=1}^{\infty} \alpha_m \phi(\vec{x}^T) \cdot \phi(\vec{y}), \quad \alpha_m \geq 0$$

$$\iint K(\vec{x}, \vec{y}) g(\vec{x}) g(\vec{y}) d\vec{x} d\vec{y} > 0, \int g^2(\vec{x}) d\vec{x} < \infty \quad (18)$$

The kernel is a valid internal product in the input field

$$K(\vec{x}, \vec{y}) = \phi(\vec{x}^T) \cdot \phi(\vec{y}). \quad (19)$$

The dual Lagrangian in the F space, given in Eq. [14].

$$L(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K(\vec{x}_i, \vec{x}_j) - \lambda \sum_{i=1}^n \alpha_i y_i \quad (20)$$

subject to  $\sum_{i=1}^n \alpha_i y_i = 0$  and  $\alpha \geq 0, i = 1, 2, \dots, n$

and the decision function is

$$f(\vec{x}) = \text{sign}[\sum_{i \in \{SV\}} \alpha_i y_i K(\vec{x}_i, \vec{x}) + b_0], \quad (21)$$

Where

$$b_0 = -\frac{1}{2} \left\{ \min_{\{\vec{x}_i | y_i = +1\}} \left( \sum_{j \in \{SV\}} \alpha_j y_j K(\vec{x}_i, \vec{x}_j) \right) + \max_{\{\vec{x}_i | y_i = -1\}} \left( \sum_{j \in \{SV\}} \alpha_j y_j K(\vec{x}_i, \vec{x}_j) \right) \right\} \quad (22)$$

SVM has employed a variety of candidate kernel functions, including poly-nominal  $K(\vec{x}, \vec{y}) = (1 + \vec{x} \cdot \vec{y})^d$ , Gaussian RBF  $K(\vec{x}, \vec{y}) = \exp\left(-\frac{\|\vec{x}-\vec{y}\|^2}{2\sigma^2}\right)$ , exponential RBF  $K(\vec{x}, \vec{y}) = \exp\left(-\frac{\|\vec{x}-\vec{y}\|}{2\sigma^2}\right)$ , and their kernel summing combinations of kernel coefficients products [33]. The Gaussian RBF kernel function is employed in this work to predict membrane protein types.

2) *K-Nearest Neighbor (KNN)*: K- nearest neighbor classifier, input data based on instance-based learner, into its feature space. KNN is based on the neighbor set that will be found near k object. KNN locate on majority voting among the k-data samples. Which store all value of the training data and wait till new data arrived to be classified on similarity measures or as a pattern matching techniques [2]. K-nearest finding based on Euclidean distance eq. (23). To classify membrane proteins, predicting the functional types of membrane proteins is indispensable [24]. Therefore similarity measures formula as the Euclidean distance ( $E_{Dis}$ ) phrase between two points( $y_1, y_2$ ) [27].

$$E_{Dis}(y_1, y_2) = \sum_{r=1}^N \sqrt{(y_1 r - y_2 r)^2} \quad (23)$$

The next steps are to generalize K-nearest neighbor classifier innovations. Metric distance and functions are measured to measure the distance between characteristics. The k-parameter must be designed for training data.

#### IV. RESULTS AND DISCUSSION

In this study, two different types of feature extraction techniques, namely PseAAC and sequence to integer encoding, are used, giving a feature vector of 62029 instances in a row and 51-dimension in a column classified by the proposed model, were 43418×51 training and 18611×51 test samples are implemented. For getting best accuracy, various type of classifiers such as SVM(Support vector machine), KNN (K Nearest neighbor), RF (Random Forest), classifiers are used and based on the result obtained from them, and the model is built [2],[16],[22],[32].

##### A. Accuracy

The number of instances rightfully predicted out of total number of instances in eq. (24).

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} \quad (24)$$

where  $TP$  is total number of true positive,  $FN$  is total number of false-positives,  $TN$  is total number of true negatives and  $FP$  is total number of false positives [3],[4]. The overall accuracy as eq. (24) of different classifiers are shown in Table II.

TABLE II. OVERALL ACCURACY OF DIFFERENT CLASSIFIERS

| Classifier    | Accuracy |
|---------------|----------|
| Random Forest | 89.38    |
| KNN           | 93.24    |
| SVM (RBF)     | 85.86    |

##### B. Specificity

Specificity of classifiers are good as the true negatives are correctly identified as calculated by eq. (25).

$$Specificity = \frac{TN}{TN+FP} \quad (25)$$

The specificity of classifiers is shown in Table III. The specificity range is 85% to 99% because TNR (true negative rate) is good. Specificity results noticed that wrongly classified samples are significantly less in KNN classifier [2].

##### C. Sensitivity

The classifier can correctly predict in eq. (26) the true positives shown in Table IV.

$$Sensitivity = \frac{TP}{TP+FN} \quad (26)$$

##### D. F-measure

Every model design to handle various kinds of multiclass problem to look at the accuracy of that model as the number of samples corrects predicted and misclassified from all prediction. Confusion Matrix gives detailed information about the failure in predictions for an unseen dataset sample. The F1 measures mathematically computed in eq. (27-28) recorded precision and recall balance values.

$$Precision = \frac{TP}{TP+FP} \quad (27)$$

$$F_{measure} = 2 \times \frac{Precision \times Recall}{Precision + Recall} \quad (28)$$

TABLE III. SPECIFICITY

| Types      | Specificity of classifiers |      |      |
|------------|----------------------------|------|------|
|            | Random Forest              | KNN  | SVM  |
| GPI        | 0.88                       | 0.99 | 0.99 |
| Lipid      | 0.93                       | 0.95 | 0.97 |
| Multi-pass | 0.95                       | 0.96 | 0.96 |
| Peripheral | 0.97                       | 0.96 | 0.96 |
| Type1      | 0.94                       | 0.97 | 0.98 |
| Type2      | 0.94                       | 0.95 | 0.85 |
| Type3      | 0.96                       | 0.95 | 0.93 |
| Type4      | 0.93                       | 0.93 | 0.94 |

TABLE IV. SENSITIVITY

| Types      | Sensitivity of classifiers |      |      |
|------------|----------------------------|------|------|
|            | Random Forest              | KNN  | SVM  |
| GPI        | 0.06                       | 0.25 | 0.05 |
| LIPID      | 0.53                       | 0.65 | 0.47 |
| MULTI-PASS | 0.97                       | 0.96 | 0.99 |
| PERIPHERAL | 0.93                       | 0.95 | 0.82 |
| TYPE1      | 0.59                       | 0.64 | 0.35 |
| TYPE2      | 0.43                       | 0.63 | 0.51 |
| TYPE3      | 0.49                       | 0.64 | 0.53 |
| TYPE4      | 0.15                       | 0.35 | 0.26 |

F1 balanced accuracy are used as a better metrics for a multi class imbalanced dataset classification task. F1-measure of various classifiers are shown in Table V.

#### E. Mathew's Correlation Coefficient (MCC)

Standard measure in machine learning MCC was suggested in 1975 by Brain W. Matthews [28]. Mathew's correlation coefficient is balanced in binary classifications into true and false positives and negatives classes [25]. It found a degree of correlation in the predicted level. It returns a value between -1 and +1. +1 represents a perfect prediction, and -1 represents the entire disqualifying range between predicting and observation eq. (29), shown in Table VI. If D datasets and N is the total number of the outcome of true and false positives and negatives views from a single instance, the Matthews correlation coefficient best such measures in larger dataset achieves a high proportion of correct predictions from the confusion matrix [11].

$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP+FP)(TP+FN)(TN+FP)(TN+FN)}} \quad (29)$$

From confusion matrix observation, is found that all classifiers perform well on multiclass datasets, KNN measure better as compared to other classifiers on the various parameter such as precision, recall, specificity, accuracy and F1-measure in MCC value.

Prediction result indicated that shown in Fig. 3, the proposed method achieved high prediction accuracy for the independent datasets. The different classifier prediction performance measures in confusion matrix results are represented in Table VII, VIII, and IX. Statistical Problem handed through machine learning is known as a confusion matrix. The proposed learning model, field error defined in the matrix table, also describes the classifier's efficiency to testing data samples, the actual value visualizing true identity on an algorithm. The confusion matrix makes easy identification of confusion between classes or mislabeled class of others in performance on various scales.

#### F. Confusion Matrix Results

The classifiers output was examined using independent tests [7]. The ensemble classifiers such as Random Forest 89.38% value, SVM 85.86%, and KNN value is improved, i.e., a maximum of 93.24%. Membrane Protein is a multilabel dataset. The classification results of models are shown in the confusion matrix, the total number of passable similarity matching with other multi-class functions. The confusion matrix is a critical way to summarize machine learning classifiers' performance, like SVM, RF (Random Forest), and

KNN classifiers. This Square matrix consists of based on features PseAAC and Sequence Length encoding. There are 62029 rows (43418 training rows and 18611 test rows in datasets) in total protein sequence and 51D descriptor size in columns. Moreover, this is listing the number of instances as absolute or relative actual class vs. predicted class ratio. The confusion matrix results demonstrate a major role in prediction identification in terms of accuracy, precision, recall, and F-1 score. SVM, KNN, and RF three learning techniques were analyzed based on outcome comparisons to find model performance. Parameter of the confusion matrix observed that the learning model KNN performs well in all eight membrane protein types. Overall, classifier performance observed a high boosting rate for large data training samples. Multipass large data sample observed F1-Score 95% in RF, 96% in KNN, and 89% in SVM where GPI class score poorly 7% in RF, 25% in KNN, and 5% in SVM. Peripheral and multipass transmembrane class are more sensitive in all classifiers, where GPI and Type-4 found a less sensitive class, with 99% GIP specificity found in KNN and SVM. In all classifiers observation found, the integrated 51D features of protein sequences and different patterns length, KNN classifier, provide better performance for membrane protein types, as shown in Fig. 4.

TABLE V. F1 SCORE

| Types      | F1 Score      |      |      |
|------------|---------------|------|------|
|            | Random Forest | KNN  | SVM  |
| GPI        | 0.07          | 0.25 | 0.05 |
| LIPID      | 0.60          | 0.65 | 0.56 |
| MULTI-PASS | 0.95          | 0.96 | 0.91 |
| PERIPHERAL | 0.91          | 0.92 | 0.89 |
| TYPE1      | 0.68          | 0.69 | 0.51 |
| TYPE2      | 0.61          | 0.69 | 0.67 |
| TYPE3      | 0.66          | 0.74 | 0.70 |
| TYPE4      | 0.27          | 0.49 | 0.28 |

TABLE VI. MCC VALUE OF DIFFERENT CLASSIFIERS

| Classifier    | MCC Value |
|---------------|-----------|
| Random Forest | 80.6      |
| KNN           | 85.4      |
| SVM(RBF)      | 82.2      |

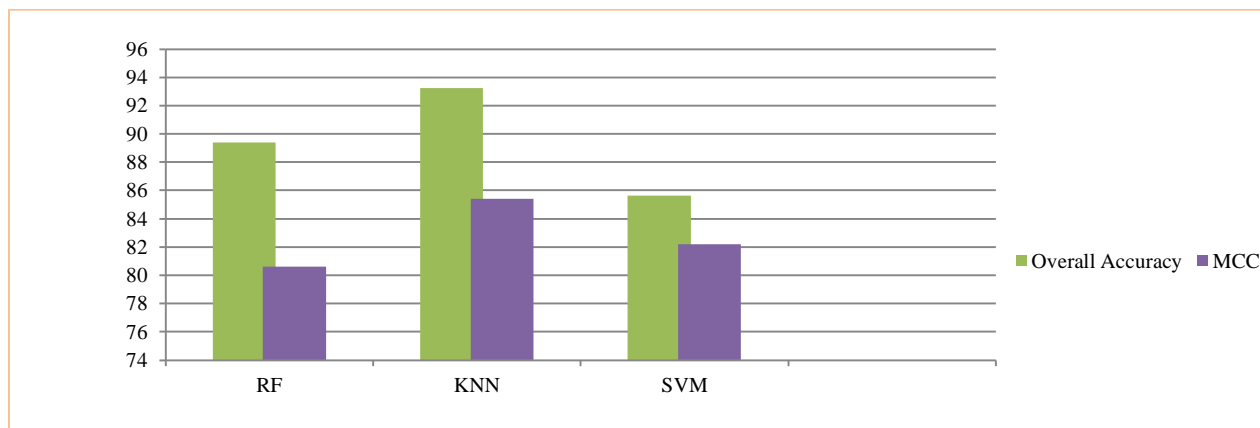


Fig. 3. Overall Accuracy and MCC Performance Scale in Bar-Chart.

TABLE VII. RESULT OF SVM RBF CONFUSION MATRIX; OVERALL ACCURACY: 85.86018485678972

|            | Gpi | Lipid | Multi-pass | Peripheral | Type1 | Type2 | Type3 | Type4 |
|------------|-----|-------|------------|------------|-------|-------|-------|-------|
| Gpi        | 9   | 139   | 37         | 1          | 2     | 0     | 0     | 0     |
| Lipid      | 158 | 418   | 262        | 44         | 1     | 2     | 0     | 0     |
| Multi-pass | 0   | 10    | 10554      | 39         | 5     | 5     | 0     | 0     |
| Peripheral | 0   | 19    | 901        | 4294       | 2     | 2     | 0     | 0     |
| Type1      | 1   | 4     | 570        | 20         | 317   | 2     | 0     | 0     |
| Type2      | 0   | 4     | 291        | 32         | 2     | 340   | 0     | 0     |
| Type3      | 0   | 2     | 15         | 11         | 0     | 0     | 32    | 0     |
| Type4      | 0   | 0     | 38         | 14         | 0     | 0     | 0     | 10    |

TABLE VIII. RESULT OF KNN CLASSIFIER CONFUSION MATRIX; OVERALL ACCURACY: 93.24188725885324

|            | Gpi | Lipid | Multi-pass | Peripheral | Type1 | Type2 | Type3 | Type4 |
|------------|-----|-------|------------|------------|-------|-------|-------|-------|
| Gpi        | 53  | 123   | 9          | 9          | 14    | 3     | 0     | 0     |
| Lipid      | 111 | 586   | 41         | 121        | 33    | 14    | 0     | 1     |
| Multipass  | 14  | 62    | 10184      | 262        | 62    | 42    | 3     | 1     |
| Peripheral | 13  | 56    | 119        | 4962       | 34    | 41    | 2     | 1     |
| Type1      | 21  | 24    | 145        | 99         | 545   | 14    | 0     | 1     |
| Type2      | 3   | 21    | 71         | 127        | 28    | 422   | 1     | 0     |
| Type3      | 0   | 4     | 4          | 7          | 4     | 2     | 38    | 0     |
| Type4      | 2   | 7     | 1          | 19         | 1     | 4     | 0     | 18    |

TABLE IX. RESULT OF RANDOM FOREST CLASSIFIER CONFUSION MATRIX; OVERALL ACCURACY: 89.38606749422322

|            | Gpi | Lipid | Multi-pass | Peripheral | Type1 | Type2 | Type3 | Type4 |
|------------|-----|-------|------------|------------|-------|-------|-------|-------|
| Gpi        | 13  | 152   | 26         | 7          | 9     | 4     | 0     | 0     |
| Lipid      | 132 | 502   | 95         | 153        | 13    | 12    | 0     | 0     |
| Multi-pass | 1   | 41    | 10355      | 193        | 36    | 4     | 0     | 0     |
| Peripheral | 3   | 38    | 294        | 4863       | 25    | 5     | 0     | 0     |
| Type1      | 3   | 20    | 299        | 32         | 489   | 6     | 0     | 0     |
| Type2      | 1   | 21    | 186        | 139        | 20    | 303   | 0     | 0     |
| Type3      | 0   | 4     | 13         | 9          | 4     | 0     | 29    | 0     |
| Type4      | 0   | 0     | 11         | 30         | 2     | 0     | 0     | 8     |

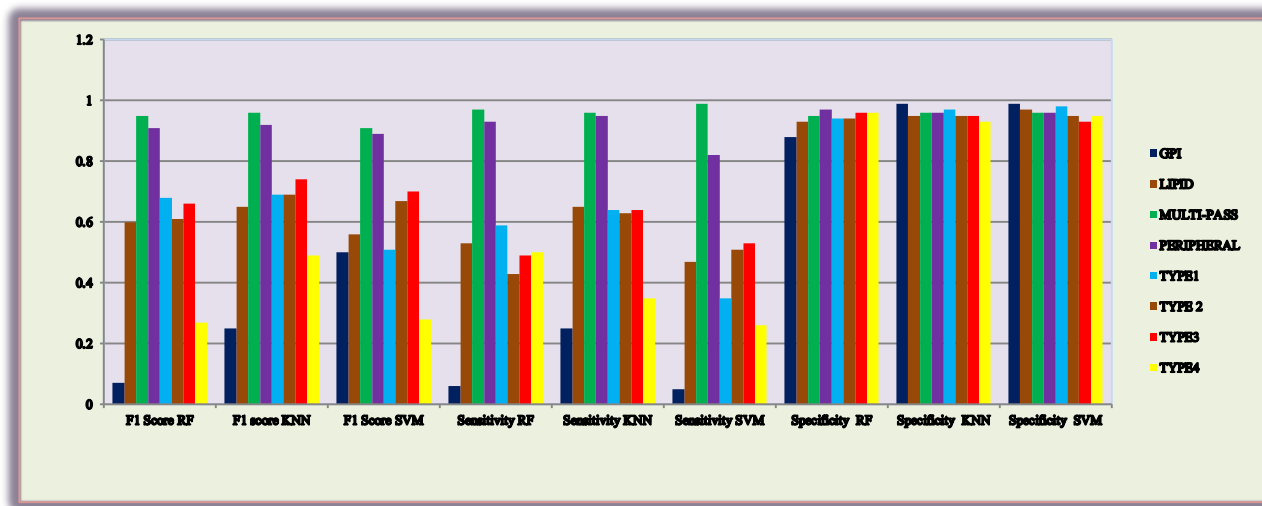


Fig. 4. F-1 Score, Sensitivity and Specificity Bar Chart of Membrane Protein Types.

## V. CONCLUSION

The proposed model objective is to score proper functions based on PseAAC and Sequence length of 51 descriptor features. This study confirmed a large sample size and fine-tuning techniques enforcement provides to build superior models that allow integrations of variant feature levels. In KNN, learning strategies based on the nearest neighbor's weight vector exploit the overall membrane protein types in a biological cell network to find the correct eight types of membrane protein. Prediction based on 51D feature vectors is used to learn three classifiers Random forest, K-nearest neighbors, and SVM. Python programming is supported by many machine learning techniques potent today. Python library provides many functions to learn about the Specify-Compile-Fit workflow that will be easy to make predictions. It can build simple necessary tools for various learning methods and generate predictions with them. Real classification results show that the proposed model achieves the desired goal significantly.

## REFERENCES

- [1] Ali, F., Hayat, M., "Classification of membrane protein types using voting feature interval in combination with Chou's pseudo amino acid composition". *J. Theor. Biol.* 384 78-83,2015.
- [2] Anjna J.Deen, Manasi Gyanchandani, "Improved Machine Learning using Adaptive Boosting algorithm in Membrane Protein Prediction", *International Journal of Innovative Technology and Exploring Engineering* Vol.8(12), page 3131-3137,2019.
- [3] Cai, Y.D., Chou, K.C., "Predicting membrane protein type by functional domain composition and pseudo amino acid composition". *J. Theor. Biol.* 238, 395-400,2006.
- [4] Cai, Y.D., Ricardo, P.W., Jen, C.H., Chou, K.C., "Application of SVM to predict membrane protein types". *J. Theor. Biol.* 226 (4), 373-376,2004.
- [5] Chen, W., Ding, H., Feng, P., "iACP: a sequence-based tool for identifying anti-cancer peptides", *Oncotarget* 7, 16895-16909,2016.
- [6] Chen, W., Lin, H., "Pseudo nucleotide composition or PseKNC: an effective formulation for analyzing genomic sequences". *Mol. Biosyst.* 11, 2620-2634,2015.
- [7] Chen, Y.K., Li, K.B., "Predicting membrane protein types by incorporating protein topology, domains, signal peptides, and

physiochemical properties into the general form of Chou's pseudo amino acid composition". *J. Theor. Biol.* 318, 1-12,2013.

- [8] Elisabeth P carpenter, Konsatinos Beis, Alexander D , So Iwata., Overcoming the challenges of membrane protein crystallography. *Curr Opin Struct Biol.* PMC ID. 18(5), 581-586,2008.
- [9] [9] Chen Z, Zhao P, Li F, Leier A, Marquez LogoTT, Ang Y, Webb GI, Smith AI, Daly RJ, Chou KC, Song J." iFeature: a python package and web server for feature extraction and selection from protein and peptide sequence". *Bioinformatics* Volume 34 issue 14,15 page2499-2502, 2018.
- [10] Qiao Ning, Zhiqiang Ma, Xiaewi Zhao. "dformKNN -PseAAC detecting formylation site from protein sequence using K-nearest neighbor algorithm via chou's 5-step rule and pseudo component", *Journal of Theoretical Biology.*470,43-49,2019.
- [11] Xiao-Sheng, Run-Jing Zhan. "clustering based subset ensemble learning method for imbalance data", *proceeding 13, ICMLC* ;35-39,2013.
- [12] Marco Punta, Lucy R. Forrest, Henry Bigelow, Andrew Kernytsky, Jinfeng Liu, and BurkhardRost. "membrane protein prediction methods" *NIH Public access* PMC ; 41(4): 460-474,2007.
- [13] Cheng, X., Zhao, S.G., Xiao, X., "iATC-mISF: a multi-label classifier for predicting the classes of anatomical therapeutic chemicals". *Bioinformatics* 33, 341-346,2017.
- [14] Chou, K.C., "Insights from modelling three-dimensional structures of the human potassium and sodium channels". *J. Proteome Res.* 3, 856-861,2004.
- [15] Chou, K.C., "Impacts of bioinformatics to medicinal chemistry". *Med. Chem.* 11, 218-234,2015.
- [16] Chou, K.C., Elrod, D.W., "Prediction of membrane protein types and sub cellular locations". *Proteins Struct. Funct. Bioinf.* 34 (1), 137-153,1999.
- [17] Chou, K.C., Shen, H.B., "MemType-2L: a web server for predicting membrane proteins and their types by incorporating evolution information through Pse-PSSM". *Biochem. Biophys. Res. Commun.* 360 (2), 339-345,2007.
- [18] Gao, Q.B., Ye, X.F., Jin, Z.C., He, J., "Improving discrimination of outer membrane proteins by fusing different forms of pseudo amino acid composition". *J. Anal. Biochem.* 398, 52-59,2010.
- [19] Golmohammadi, K.S.K., Crowley, L., Reformat, M., "Classification of cell membrane proteins". *Front, Convergence Biosci. Inf. Technol.* 153-158,2007.
- [20] Golmohammadi, S.K., Kurgan, L., Crowley, B. Reformat, M., "Amino acid sequence based method for prediction of cell membrane protein types". *Int. J. Hybrid Inf. Technol.* 1 (1), 95-109,2008.
- [21] Hayat M., Khan, A., "Predicting membrane protein types by fusing composite protein sequence features into pseudo amino acid composition". *J. Theor. Biol.* 262, 10-17,2011.

- [22] Hayat M., Khan, A., Yeasin, M., "Prediction of membrane proteins using split amino acid and ensemble classification". *Amino Acids* 42 (6), 2447-2460,2012.
- [23] Wang, S.Q., Yang, J., Chou, K.C., "Using stacked generalization to predict membrane protein types based on pseudo-amino acid composition". *J. Theor. Biol.* 242 (4), 941-946,2006.
- [24] Wan, S., M.W., Kung, S.Y., "Mem-mEN: predicting multi-functional types of membrane proteins by interpretable elastic nets". *IEEE/ACM Trans. Comput. Biol. Bioinf.* Doi: 10.1109/TCBB.2015. 2474407,2015.
- [25] Jia, J., Zhang, L., Liu, Z., Psumo-CD: "predicting sumoylation sites in proteins with covariance discriminate algorithm by incorporating sequence-coupled effects into general PseAAC", *Bioinformatics* 32, 3133-3141,2016.
- [26] Cheng, X., Zhao, S.G., Xiao, X., "iATC-mHyb: a hybrid multi-label classifier for predicting the classification of anatomical therapeutic chemicals". *Oncotarget* doi:10.18632/oncotarget.17028,2017.
- [27] Shen, H., Chou, K.C., "Using optimized evidence-theoretic K-nearest neighbour classifier and pseudo-amino acid composition to predict membrane protein types". *Biochem. Biophys. Res. Commun.* 334 (1), 288-292,2005.
- [28] B.W. Matthews, "Comparison of the predicted and observed secondary structure of T4 phage lysozyme", *Protein Structure.* 405 (2),pp.442-451,1975.
- [29] Mahdavi, A., Jahandideh, S., "Application of density similarities to predict membrane protein types based on pseudo-amino acid composition". *J. Theor. Biol.* 276, 132-137,2011.
- [30] Nanni, L., Lumini, A., "An ensemble of support vector machines for predicting the membrane protein type directly from the amino acid sequence". *Amino Acids* 35 (3), 573-580,2008.
- [31] Wang, M., Yang, J., Liu, G.P., Xu, Z.J., Chou, K.C., "Weighted-support vector machines for predicting membrane protein types based on pseudo-amino acid composition". *Protein Eng. Des. Sel.* 17 (6), 509-516,2004.
- [32] Shen, H.S., Chou, K.C., "Using ensemble classifier identify membrane protein types". *Amino Acids* 32, 483-488,2007.
- [33] N. Cristianini, J. Shawe-Taylor, *An Introduction to Support Vector Machines*, Cambridge University, Cambridge, 2000.
- [34] V.N.Vapnik, *The Nature of Statistical Learning Theory*, second ed., Springer, New York, 1999.

# Real Time Implementation and Comparison of ESP8266 vs. MSP430F2618 QoS Characteristics for Embedded and IoT Applications

Krishnaveni Kommuri<sup>1</sup>

Research Scholar, Dept. of ECE Koneru Lakshmaiah  
Education Foundation, Vaddeswaram, AP

Venkata Ratnam Kolluru<sup>2</sup>

Associate Professor, Dept. of ECM Koneru Lakshmaiah  
Education Foundation, Vaddeswaram, AP

**Abstract**—This research article proposes a novel Smart Communication Platform (SCP) to improve the Quality of Service (QoS) parameters in real time by using MSP430F2618. A static network has been implemented with narrow band Internet of Things (IoT) architecture which contains 10 nodes. SCP performs tracking of environmental parameters like Temperature, Humidity, Pressure, Proximity and light. A prototype has been developed by using Open source Red hat Linux 14.4 version and programmed in Embedded C.MSP430F2618 has been configured as master and slave nodes, the output is observed in a serial monitor and Gateway as well. The QoS parameters of MSP430F2618 and ESP8266 are compared in terms of power. The power consumption improvements of QoS (Quality of Service) analysis results are around 1.01mW has been seen with the experimental setup. These empirical results are much useful for wireless sensor network and IoT applications.

**Keywords**—Communication; ESP8266; gateway; Internet of Things (IoT); MSP430; power; sensor

## I. INTRODUCTION

The Embedded systems and Internet of Things (IoT) are the network of tiny, intelligent devices, actuators and other tools build together. These systems are mostly wireless and powered by batteries to ease the deployment. The popular version of the Embedded Systems development life cycle is designed from waterfall model. The basic feature of this model is managerial control and Modulate. This has distinct goals at every development step. Different scheduling mechanisms can be introduced at every stage of module and interface. This model is suitable for linear process models takes several steps. Those are 1) Requirement gathering 2) Pre-Design 3) Design review 4) Design implementation 5) Design final app Overview. These can also be put it another terms as Design, Coding, Testing, Maintenance. It is been seen that each phase is also connected back to earlier stage where it helps needed verification.

## II. LITERATURE REVIEW

Recent advancements in IC technologies, Communication devices, Sensors places an immediate impact on embedded systems and IoT [1]. Hardware and Software co design technique like, design flow approach with prototyping principles are being used in [2] to build SCP implementation. Embedded system requires immense knowledge in the areas of Electronic devices like C-motes, CPU, GPU, and Gateway [3] to check

the functionality of the proposed prototype. In the recent embedded system research domains, processor integrate new and relevant devices by adopting code motion techniques and independent APIs [4]. A prototype is designed to classify the systems, such as one to one node, one to multi node, multi to one node and multi to multi nodes via WINGS gateway and stores the data in cloud. Traditional data processing and measurement methods [5] are being slowly replaced by single wire and wireless communication technology where it is useful lot of starter embedded system application.

To reduce the duration of data transmission, and to improve the network lifetime in [6] like PRIMS algorithm, proposed a distance-based transmission rate selection and maximum emission rate (MER) determination. The web based remote dynamic data collection and Storage monitoring system has been proposed in [7] for data transfer and system characteristics are analyzed for centralized management. The power optimization techniques such as resource consolidation, virtualization, selective connectivity, and proportional computing are considered in [8] to improve the proposed experimental set up. The Network Time Protocol (NTP) is used in [9] for time synchronization, allowing a more flexible network by avoiding the system placement needs. By using these devices one can get data in the form of (i) analog, (ii) digital values (iii) channels of various communication blocks, (iv) Software development tools and (v) Real Time Operating System (RTOS) parameters can be utilized [10]. Interfacing of an advanced micro controller like MSP430 with Wi-Fi makes the user to perform interactive operations. A practical approach [11] on MSP435 based TI CC3220SL results the pervasive computing aspects. Synchronized sampling control method [12] gives better ways to manage wireless node state information to improve flexibility and agility during load balancing, fail over, and life cycle operations, and designing VNFs to allow for their transparent migration across central and edge clouds. Inter-networking arrangement [13] gives the integration of individual modules into systems of architecture, often in network protocol which build the Green Network. Design of structures require environments and modelling methods [14] to understand the operational types and processes affecting the system. Popular patterns for VNFs are set up in future study directions. Future trend universal computing work explores to exploit utility computing and the Internet of Things in [15-17] to move ubicomp systems where computing is made to appear anytime and everywhere. Today's research

turned in various paths such as Power Optimizations [18-19], Wireless sensor networks [20], Embedded systems [21], and reconfigurable antennas with IoT applications [22-25].

This platform gives 1) Design and development of hardware setup with soft real time values and test results. 2) Supports experimental digital/analog interfacing with wireless communication. 3) Design of new device set ups for data transmission and reception techniques useful for Edge triggering applications. 4) Experimental results for QoS Analysis.

This article has been divided into various sections. Section II briefed about Literature survey. Section III deals with motivations of current Heterogeneous systems. Section IV discusses about proposed algorithm and implementation steps of MSPF2618. Section V presents results and discussions for QoS comparison analysis. The paper is ended with conclusion and future scope in Section VI.

### III. MOTIVATION TOWARDS HETEROGENEOUS EMBEDDED SYSTEMS

Central Processing Unit (CPU) and Graphical Processing Unit (GPU) has employed in various applications like monitoring, data analytics etc. Master simulator design useful to help developers of slave devices or other devices how to test and simulate. Slave simulator as shown in Fig. 1(a) send the data from the foreign procedure calls. Assign the slave ID, address, size can read and write the registers. There are several data formats available with word order swapping, such as float, double, and long to program. These Processing units (PUs) shows unique features and strengths to conquer high performance computing [2]. Various techniques have applied to develop heterogeneous computing with peta scale range fused with design of CPU-GPU chips as shown in Fig. 2. Partitioning enables PUs performance and energy efficiency for heterogeneous computing techniques (HCTs) [4]. Computing has explained in below sections.

#### A. Design Strengths of PUs

Multi core GPU keeps Instruction set arises in CPU and few tens of PUs integrated on single chip though it is built with several different architectures. Single thread, large size caches with minimum latency, at high frequency rates leads the high through puts. CPU are clearly used for critical latency applications while GPU used for large throughput-critical applications. Thus, heterogeneous system used in number of applications with high performance context.

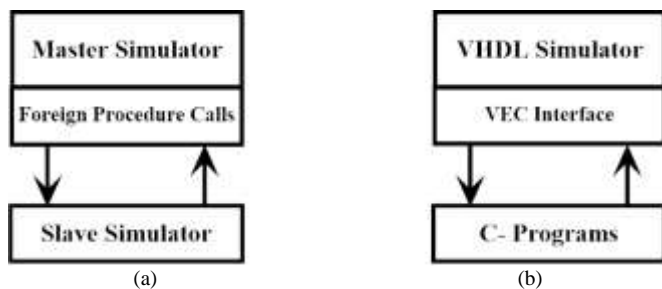


Fig. 1. Master Slave Co-Simulation (a) Hardware Architecture (b) Software Architecture in Embedded System and IoT Platform.

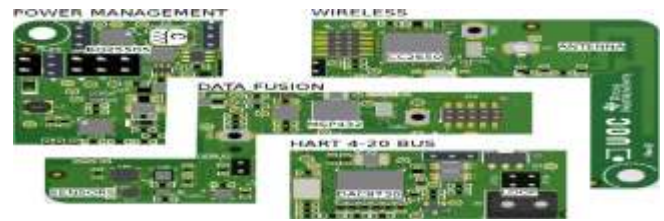


Fig. 2. Schematic of Embedded Chip Fabrication.

#### B. Embedded Algorithmic Features of PUs

For many applications of data transfers relishes execution time or case GPU cores does not allow uninterrupted execution and branch mitigations etc. In compared to GPUs CPUs own better performance in single applications at different stages which leaches the time factors of delivery of application.

#### C. Resource use Improvements

To conclude the resource utilization's to both CPU and GPU are over featured though the use sends as low. Sometimes CPU stays idle if it sends the control to GPU and in reverse GPU bandwidth memory decreased.

#### D. Heuristic Algorithms

By the design of Wolf [2], this kind of algorithm have a pool of heterogeneous elements of tasks and provides communication among it. Communication links, PEs and Task graphs acts as input to the algorithm. There are major aims and minor goals like to meet specified rate for execution, minimize the total cost.

Many optimization problems [1] get the solutions from approximations of Heuristic algorithms. The best possible solutions are in maximum or minimum solutions to objective function. The solution is a function used to evaluate the objective function. Optimization problems are nothing but talk about several real-world issues. In all search algorithms solutions can be defined as search space and optimization algorithms, etc.

Dynamic programming branch and bound techniques effectively present in Heuristic practices, gives us time complexity and not completed tasks. Premature convergence is a significant drawback in Hill-Climbing algorithm because of it always fetches the nearest local optima of low quality. This can be targeted by, SIMULATED ANNEALING ALGORITHM: It resembles the same to Hill climbing but sometimes accepts results worse than the present scenarios where the fault tolerance will be acceptable and decreasing with time slices. TABU SEARCH: Continued the idea where neglect the local optimization by introducing several memory data structures. The "jump" instruction repeated in its loop acts as bottleneck, this is prohibited in Tabu Search. SWARM INTELLIGENCE: Introduced in 1989 by Gerardo Beni and Jing Wang. This is developed among self-organized, decentralized systems collective behavior cellular Robotic systems. The popular approaches are Particle Swarm Optimization (PSO) and Ant Colony Optimization (ACO). In ACO the design done on future ants to build improved solutions by showing the variations in the graph and changing its way. PSO it is entirely different that solutions produce point or surface in an n-



dimensional space. The main possible advantage is that local optima [7] can be resolved from impressively resistant algorithms in multi-application design aspects.

IV. PROPOSED ALGORITHM AND IMPLEMENTATION STEPS OF MSP430F2618 FOR QOS COMPARISON ANALYSIS

The rapid advancement in the areas of wireless networks, information technologies, sensor design and semiconductor has been put to the proliferation of Wireless Sensor Network (WSN). WSN is growing as a backbone technology for various applications such as agriculture, traffic control, natural disaster relief, health monitoring and control, home automation, environment and habitat monitoring, consumer and industrial applications, product quality monitoring, seismic sensing etc. A WSN features are Low power consumption, self-organization, fault tolerance, low cost, and longer standalone life.

To implement the C-Mote set up in [3] used Red hat Linux 14.4 with SMA Antenna, Ubisense sensor, USB power jacks, data acquisition card through serial window, with I/O expansions and processor controls. Coming to C-Mote its specifications it has USB interface to PC, 256kB RAM+1MB XIP FLASH provides application level security, External USB peripherals, 2 to 20 pin connectors, 4 wire JTAG pins, On board antenna and GNU debugger support as well.

UbiSense is the sensor board developed for testing and integrating various sensors with C-Mote, UbimoteHR. UbiSense has applications where the users can make use of sensor data by plugging in the module directly to the Mote. The Sensors on UbiSense are I2C compatible and user can avail the advantage in interfacing. C-Mote has Female SMA Antenna connector. Antennas configured for 2.4GHz 50Ohms.

SMA compatible should be connected to the SMA connector. IEEE 802.11 is a wireless local area network (WLAN) for implementing computer communication in the 900 MHz and 2.4, 3.6, 5, and 60 GHz frequency bands. Most of the applications like Smart phones, laptops, Office Networks, Homes, etc. were widely used by this Wireless Networks.

A. Working Prototype

In Tele communication, Internet, and Data Communication the application interfacing happens through wireless via several protocols for remote locations. Communication happens from small to large from one point to another point. Depending on the height of the antennas and other devices, the frequency and power level used, and the surrounding environment, communications signals can travel up to tens of miles to its designated location.

The proposed algorithm, shown in Fig. 3 starts with sensors, here Ubisense sensor used. This Ubisensor board gives the experimental set up environmental values like Temperature, Humidity, Proximity, Light Intensity and Barometric pressure, etc. C-Motes are connected as Master and Slave node assigned with node ids. It is easily understandable from the flow diagram, as shown in Fig. 3. The sensor outputs are taken from the minimum values. If not, the setup again initializes from the start. The values are checked in the Wings gateway as well as it is seen in every

slave node. QoS analysis improvised [9] by the values predicted from C-Mote Master and Slave nodes. Validation with respect to Power are considered and observations were notified else algorithm repeated.

It is arranged C-Mote nodes as per the proposed setup shown in Fig. 4. The nodes are categorized as Master node which is shown as light grey color marked as M and Slave nodes which is shown as dark grey color marked as S. The Fig. 4(a) shows the point to point arrangement of C-Motes one as Master and another as Slave Node. Fig. 4(b) shows the point to multi point arrangement of C-Motes one as Master and two as Slave Node. Fig. 4(c) shows the point to multi point arrangement of C-Motes one as Master and nodes as Slave Node. In Fig. 4(c) Multi point communication made from Master node and Slave node and vice versa.

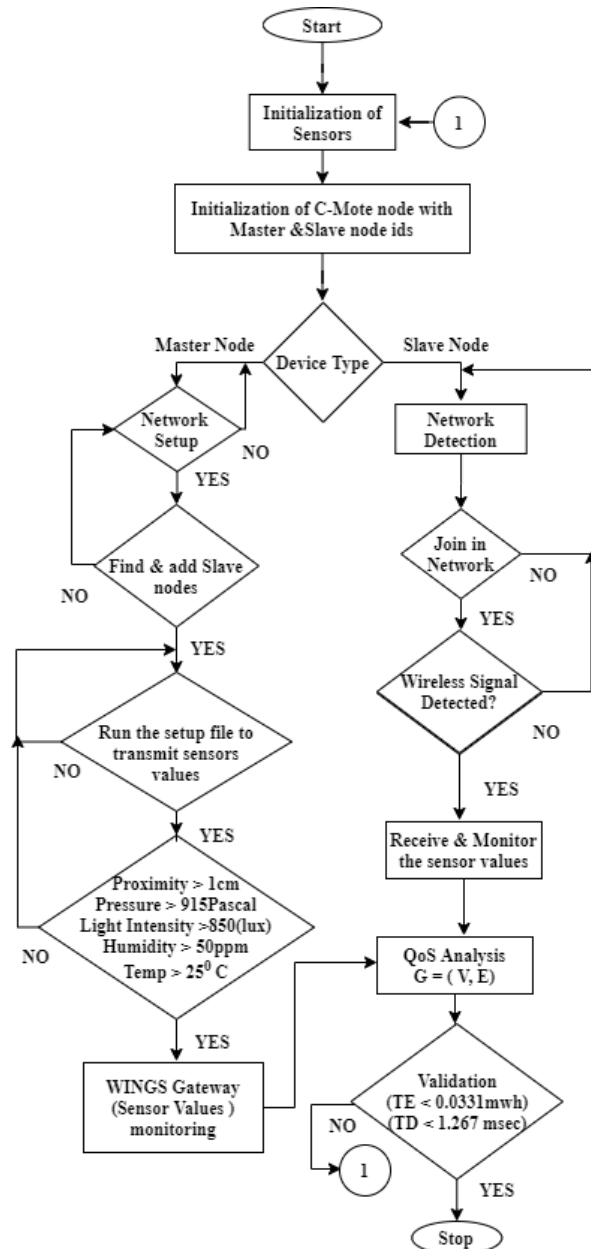


Fig. 3. Flowchart of MSP430F2618 QoS Analysis.

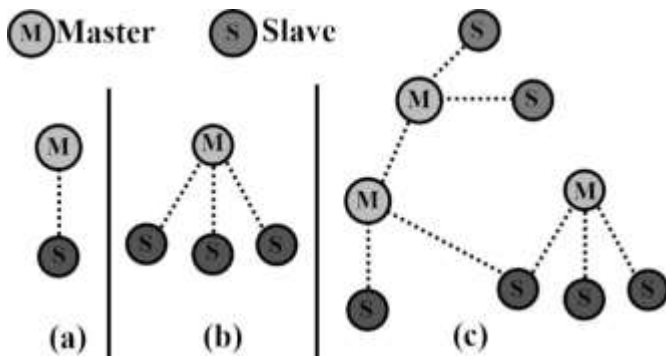


Fig. 4. Proposed Setting up of C-Mote Nodes in Various Combinations.

1) Wireless Point to Point Communication Setup (a) First C-Mote libraries are taken from Red Hat packages and is installed and configured from GCC, GNU and LIB, etc. With the help instruction set, built the needful configuration to make the nodes as master as well as slave. The setup powered with 5V, connected the UBISENSE sensor board to C-Mote and compiled the program. Master C-Mote node send the sensor data to Slave C-Mote node. Thus, set up point to point communication as shown in Fig. 5(a).

2) Wireless Point to Multi Point communication Setup (b) Here C-Motes, one is the transmitter as master Node and three are receivers as slave nodes setup is shown in Fig. 5(b). Ran the program to set up Point to Multi point communication among Master and slave C-Mote nodes. Transmitter reads the physical parameter value from ubisense connected to the device and transmits. Receiver receives the data packet which includes the measured physical parameter transmitted by the transmitter. Receiver parses the packet and prints the data received on Hyper Terminal. This way Point-to-Multi point communication platform can be build [14] carried out via a distinct type of one-to-many connection.

3) Wireless Multi Point to Single/Multi Point node communication Setup (c) The Multi point-to-Point network [3] communication has been configured to make communication between multiple remote user terminals and central hub which in turn makes and reduces the technical needs for remote locations. Data transmissions takes place from Master nodes to Slave nodes.

C-Mote is provided with a 10-pin connector which allows for the direct plug-in of UbiSense as shown in Fig. 5. It is a sensor board with temperature and relative humidity, light intensity, barometric pressure, proximity sensing and buzzer. The sensors communicate to the micro controller through I2C protocol. The USCI module B0 is incorporated for I2C functionality from micro controller viewpoint. The A1 module of USCI is used for UART with the required baud rate. Buzzer requires a PWM control signal used for alarm generation. This value is displayed on the Hyper Terminal as shown in Fig. 6. The outputs in the serial window checked and saw.

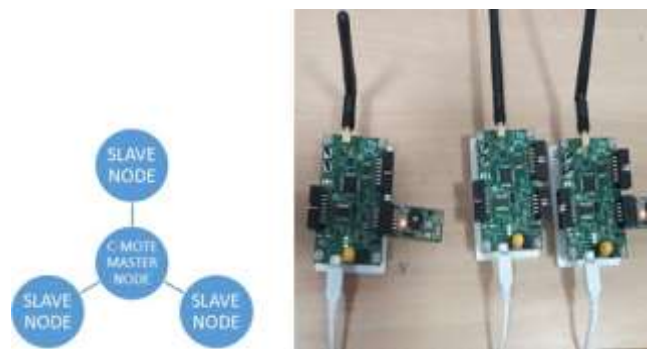


Fig. 5. C-Mote 1 Master Node and 2 Slave Nodes Setup.

```

File Edit View Search Terminal Help
Proximity: 609
Light Intensity: 786
Humidity: 51.63

-----UbiSense Data-----
Temperature: 30.2
Pressure: 916.5
Proximity: 78
Light Intensity: 856
Humidity: 51.62

-----UbiSense Data-----
Temperature: 30.2
Pressure: 916.8
Proximity: 22
Light Intensity: 863
Humidity: 51.62

-----UbiSense Data-----
Temperature: 30.2
Pressure: 916.20
Proximity: 85
Light Intensity: 858
Humidity: 51.64

-----UbiSense Data-----

```

Fig. 6. Experimental Results in Serial Window.

System design supports multi point-to-multi point communication with efficiency and reliability. A protocol management mechanism for the multi-point-to-multi point communication protocol is [6] proved from the node configuration. A user joining the system network can change the topology of the multi-point-to-multi point communication network as per the need.

In the proposed SCP, the data is collected through mounts of sensor nodes distributed in the field of sensors and supported by multi-hop wireless communication to users

towards application-specific, energy-constraint and uncertain topology. MIS (maximum independent sets) made up of cluster heads is achieved by clustering algorithms, based on which MCDS (minimally connected dominant sets) and MST (smallest spanning tree) are developed with improved Prim's algorithm and created. In addition to that, the spanning tree maintenance and update algorithm is also provided. Simulation analysis eventually shows in Fig. 8 that the proposed algorithm is successful.

A gateway [5] output terminal shown in Fig. 7 is a node on a network that serves as an entrance to another network. In enterprises, the gateway is the computer that routes the traffic from a workstation to the outside network that is serving the Web pages. In homes, the gateway is the ISP that connects the user to the internet. The wings gateway tool has totally 10 channels meet so total communication between nodes would be predicted. The open source gate which is configured the 10 maxima of nodes.

The TX nodes and Rx nodes IDs were stored and mapped to the configured networks. Data rate is high and good in accuracy was seen. It keeps the monitoring the nodes data at the input and at the output. In the hyper terminal saw the parameters readings as shown in Fig. 8.

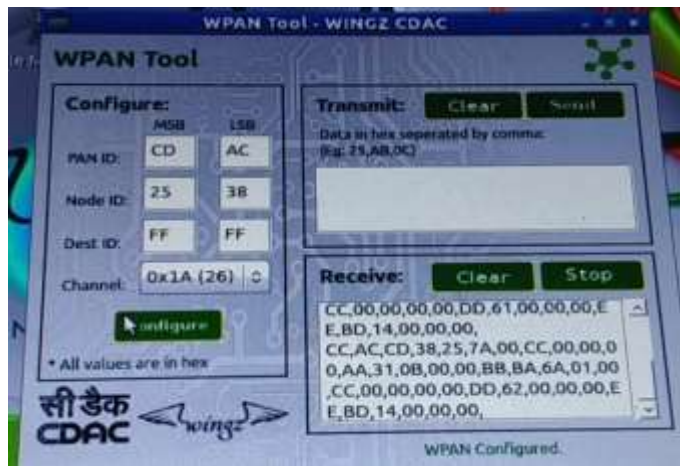


Fig. 7. Gateway Output C-Mote set up for 10No.Channels.

## V. RESULTS AND DISCUSSIONS

The results are seen from the serial window in the IDE. With the serial println functional libraries in Embedded board [10] for the environmental values are tracked for the built network. Initially it is shown for Temperature, Humidity, Light Intensity, Barometric pressure, etc. These values were chronically repeated till the node is active. The node data transfer happens from all the modes like point to point, point to multi point, multi-point to point, multi-point to multi point configurations.

All the series of sensor values are taken into the course of time (24 Hrs.) is displayed in Fig. 9 and predicted the values that affects the energy optimization methods. Improved QoS calculations QoS calculations results were shown in Fig. 9 with respect to ESP8266 and MSP430 boards.

Ubisensor prompts the values of Temperature, Humidity, Pressure, Light Intensity continuously measured over period. The nodes which are configured are static nodes [11]. If nodes are aligned dynamically causes various energy, vertices, and location change in values also. The sensor used for this prototype had a standby current resulting in power used up by both the sensor and MSP430 board during sleep mode. For the temperature sensor with the shutdown feature, the standby current parameters could be saved. As future predictions this would result in more power savings and replacing the sensor with a no-standby current would be a choice for the future to reduce power usage. The UBISENSE sensor would be high while taking the reading and then turned to low within the code before it goes into sleep mode. This leads the power optimization [12], and only the standby current would be used.

As the graph from Fig. 8 mentioned initially plotted with the point to point communication. It is seen that the energy values are almost similar except in the first 10min and 50 mins of time. There the value started from 0.3 mwh when compared to earlier result.

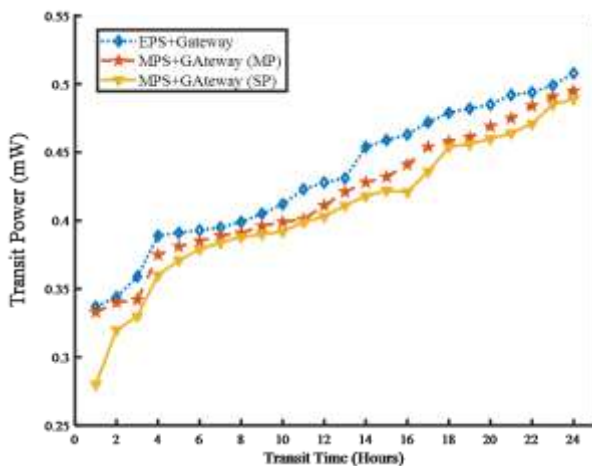


Fig. 8. Comparison Results for P-P P-M M-P M-M.

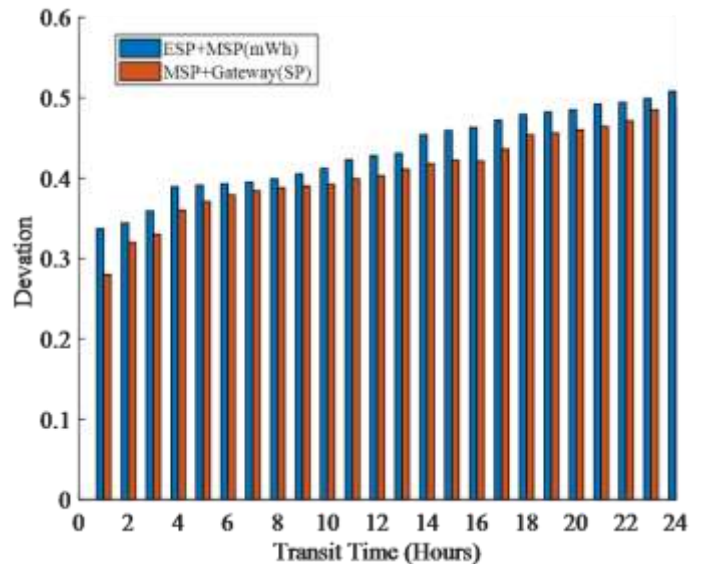


Fig. 9. Improved QoS Analysis.

The minimal spanning distance too calculated according to the theory of Prims algorithm which finds an edge of the least possible weight that connects any two trees in the forest. It is a greedy algorithm in graph theory as it finds a smallest spanning tree for a connected weighted graph adding increasing cost arcs at each step. The Greedy Choice is to pick the smallest weight edge that does not because a cycle in the MST constructed so far. own in Table I with the combination of MSP430+Gateway. This arrangement produces better results for communication applications.

In the Table II its represented with deviation of power in with respect to ESP8266 and MSP430. The comparison between both boards were compared and drawn the plots as shown in Fig. 9. The blue bar in the chart shows ESP+MSP board as the experimental set up for SCP. The red bar shows.

In the WINGZs gateway the values were seen from the Master and Slave nodes. Over 10 channels configured, 1 as master node and being still all as slave nodes and made the sensor data transmission. Nodes are assigned with static ids to build the green energy [13] network model. The values of the output sensor are measured for the software compatibility in Hexadecimal number system and converted into sensor values.

As the graph shown in Fig. 8 mentioned initially plotted with the point to point communication. It is seen that the energy values are predicted compared to earlier result [15]. By using Cooja simulator these nodes can be configured as a Zone based node, foreign nodes, local nodes, etc. The predictions or outputs are based on static allocation node positions.

From the observations of Multi point access Transit Power per node started from 0.333 mw to extend up to 0.5 mw as improvised algorithm experimental set up for SCP. In the Table II, it is shown that there is deviation.

TABLE I. SINGLE AND MULTI-POINT COMMUNICATION SETUP RESULTS

| Transmission Time in Hrs. | ESP+Gateway Power in mW. | MSP+Gateway Power M-M in mW. | MSP+Gateway Power S-P in mW. |
|---------------------------|--------------------------|------------------------------|------------------------------|
| 1                         | 0.337                    | 0.28                         | 0.28                         |
| 2                         | 0.334                    | 0.28                         | 0.32                         |
| 4                         | 0.389                    | 0.375                        | 0.36                         |
| 6                         | 0.393                    | 0.379                        | 0.379                        |
| 8                         | 0.399                    | 0.388                        | 0.388                        |
| 10                        | 0.412                    | 0.399                        | 0.392                        |
| 12                        | 0.428                    | 0.411                        | 0.400                        |
| 14                        | 0.454                    | 0.418                        | 0.418                        |
| 16                        | 0.463                    | 0.440                        | 0.428                        |
| 18                        | 0.479                    | 0.458                        | 0.454                        |
| 20                        | 0.485                    | 0.460                        | 0.46                         |
| 22                        | 0.494                    | 0.484                        | 0.471                        |
| 24                        | 0.508                    | 0.489                        | 0.489                        |

TABLE II. COMPARISON OUTPUTS OF ESP8266 AND MSP430F2618

| Transmission Time in Hrs. | ESP8266 Power in mW. | ESP8266 Power in mW. | Difference Power in mW. |
|---------------------------|----------------------|----------------------|-------------------------|
| 1                         | 0.337                | 0.28                 | 1.04                    |
| 2                         | 0.334                | 0.28                 | 1.05                    |
| 6                         | 0.393                | 0.379                | 1.04                    |
| 8                         | 0.399                | 0.388                | 1.03                    |
| 14                        | 0.454                | 0.418                | 1.09                    |
| 20                        | 0.485                | 0.460                | 1.05                    |
| 24                        | 0.508                | 0.489                | 1.04                    |

The scenario that we have considered is for unidirectional communications from [22] the sensor master node to slave node, which fits well with a sensor gathering readings or a smart button triggering an alert. Bidirectional communications [23] have significant implications on energy use because the reception circuitry must be left on in listening mode. For Gateway, bi-directional communication requires the device to be attached to the access point, which requires it to be active to receive beacon frames. MSP430 has a light sleep mode that has a timer to switch the central processing unit and radio circuitry off between beacons to save power, waking the chip up before the next beacon. However, while this offers significant reductions over keeping the chip active, the overall power usage stays in the 0.5–1-mA range, which is clearly far too high for long-term battery operation. Delays can be tolerated, for example, for updating configuration values, data can be sent to the sensor as part of the acknowledgment when the sensor sends data to the server, or the sensor can periodically poll the server even if there is no data to be sent.

## VI. CONCLUSIONS

In the last few years wireless sensor networks and IoT have drawn the attention of the research community, driven by a wealth of theoretical and practical challenges. This progressive research in WSN and IoT explored various new applications enabled by larger scale networks of sensor nodes capable of sensing information from the environment, process the sensed data and transmits it to the remote location. WSN and IoT [24-25] are mostly used in, low bandwidth and delay tolerant, applications ranging from civil and military to environmental and healthcare monitoring. The output values can be processed and mapped for several Environmental measure applications. Based on the results from this article, when combined with a low-power processor such as the MSP430, is power efficient for use in an IoT device. The results illustrated using ESP-8266 efficient if it is to be used for a short period. When used for static defined networks, the MSP430 coupled with the gateway offers effective monitoring and data handling results. The power consumption results were then carried out using static IP to prove. The use of the processor MSP430 showed in Fig. 9. an increased power saving compared to the ESP8266. This configuration can be used for IoT devices Academic experiments, data analysis, Data mining sources setups. Hence concluding that suggesting this kind of set ups for Industrial, Environmental, Cold storage's, Health Monitoring applications and many more. As

the future scope this concept can be extended for even more node set up pool and much more WSN network data transfer predictions via several spanning tree techniques for any one dedicated application platform.

#### ACKNOWLEDGMENTS

This research work was carried and supported by DST-FIST grant number FST/ETI-410/2016(C) sponsored, Internet of Things Excellence Centre, KLEF. I would also like to show my gratitude to Dr Venkata Ratnam Kolluru, Assoc.Prof., KLEF for assistance and comments that greatly improved the manuscript, although any errors are our own and should not tarnish the reputations of these esteemed persons.

#### REFERENCES

- [1] Darshana Thomas, Ross McPherson, Greig Paul, and James Irvine, "Optimizing Power Consumption of Wi-Fi for IoT Devices", IEEE Consumer Electronics Magazine, Volume: 5, Issue: 4, Oct. 2016.
- [2] Jorgen Staunstrup, Wayne Wolf, Daniel D. Gajski, Jianwen Zhu, Rainer Doerner, "Hardware Software Co-Design Principles and Practice", 1997.
- [3] Sparsh Mittal, Jeffrey S. Vetter, "A survey of CPU-GPU heterogeneous computing techniques", ACM Computing Surveys. Volume 47 Issue 4, Article No. 69, pp 1-35, July 2015.
- [4] Krishnaveni.Kommuri, K. Venkata Ratnam, Geetha Prathyusha, P. Gopi Krishna, "Development of real time environment monitoring system using with MSP430", International Journal of Engineering and Technology, vol. 7, Issue no.2.8, pp. 72-76, 2018.
- [5] Gabriel Martins Dias, Boris Bellalta and Simon Oechsner, "Using Data Prediction Techniques to Reduce Data Transmissions in the IoT", IEEE, 978-1-5090-4130-5/16/2016.
- [6] K. Sharma and T. Suryakanthi, "Optimizing Power Consumption of Wi-Fi for IoT Devices", International Conference on Green Computing and Internet of Things, pp. 1586-1593, 2016.
- [7] Junying Yuan, Huiru Cao\*, Choujun Zhan and Lin Wang, "A method of determining maximum transmission rate in wireless sensor network", Int. J. Autonomous and Adaptive Communications Systems, Vol. 11, No. 1, 2018.
- [8] Sven Jager, Tino Ungelded, Ralph Maschotta, and Armin Zimmermann, "Model-Based QoS Evaluation and Validation for Embedded Wireless Sensor Networks", IEEE Systems journal, 2014.
- [9] Niturkar Priyanka, Shinde, "Design and development of ethernet control system for embedded web server using ARM processor", International Journal Of Engineering And Computer Science Volume- 3 Issue -3, PP 4096-4099 March 2014.
- [10] P. Gopi Krishna, K. Sreenivasa Ravi, K Hari Kishore, K KrishnaVeni, K. N. Siva Rao, R. D. Prasad, "Design and development of bi-directional IoT gateway using ZigBee and Wi-Fi technologies with MQTT protocol", International Journal of Engineering and Technology, vol. 7, Issue no.2.8, pp. 125-129, 2018.
- [11] Wenxuan Yao, Haoyang Lu, Micah J. Till, Wei Gao, Yilu Liu, "Synchronized Wireless Measurement of High Voltage Power System Frequency Using Mobile Embedded Systems", IEEE Transactions on Industrial Electronics, Vol No. 65, Issue No. 3, pp 2775-2784, 2017.
- [12] Aruna Prem Bianzino, Claude Chaudet, Dario Rossi, Jean-Louis Rougier, "A Survey of Green Networking Research", IEEE Communications surveys, Tutorials, Vol. 14, NO. 1, FIRST QUARTER 2016.
- [13] Junying Yuan, Huiru Cao\*, Choujun Zhan and Lin Wang, "Towards a Virtual Network Function Research Agenda: A Systematic Literature Review of VNF Design Considerations", Journal of Network and Computer Applications. doi.org/10.1016/j.jnca.2019.102417.
- [14] Caceres, R., Friday, "Ubicomp systems at 20: Progress, opportunities, and challenges", IEEE Pervasive Computing, Cross Ref Google Scholar, Vol.No. 11, Issue no.1, pp.1421, 2012.
- [15] Hashim, Y., Idzha, A.H.M. and Jabbar, W.A., 2018. The Design and Implementation of a Wireless Flood Monitoring System. Journal of Telecommunication, Electronic and Computer Engineering (JTEC), 10(3-2), pp.7-11.
- [16] Jabbar, W.A., Shang, H.K., Hamid, S.N., Almohammed, A.A., Ramli, R.M. and Ali, M.A., 2019. IoT-BBMS: Internet of Things-Based Baby Monitoring System for Smart Cradle. IEEE Access, 7, pp.93791-93805.
- [17] Jabbar, W.A., Kian, T.K., Ramli, R.M., Zubir, S.N., Zamrizaman, N.S., Balfaah, M., Shepelev, V. and Alharbi, S., 2019. Design and fabrication of smart home with internet of things enabled automation system. IEEE Access, 7, pp.144059-144074.
- [18] W. Xiao, M.G.J. Lind, W.G. Dunford and A. Capel, "Real-time identification of optimal operating points in photovoltaic power systems", IEEE Trans. on Ind. Elec., vol.53, no.4, pp.1017-1026, 2006.
- [19] M. N. M. Hussain, A. M. Omar and A. A. A. Samat, "Identification of multiple input-single output (miso) model for MPPT of photovoltaic system", IEEE International Conference on Control System Computing and Engineering (ICCSCE-2011), Malaysia, pp. 49-53, Nov 25-27, 2011.
- [20] Srikanth, Nandoori, and Muktyala Sivaganga Prasad, "Energy efficient trust node-based routing protocol (EETRP) to maximize the lifetime of wireless sensor networks in Plateaus", International Journal of Online and Biomedical Engineering (iJOE) 15.06 (2019): 113-130.
- [21] Meghana, M., et al. "Design and Development of Real-Time Water Quality Monitoring System." 2019 Global Conference for Advancement in Technology (GCAT). IEEE, 2019.
- [22] Allam, V.K., Madhav, B.T.P., Anilkumar, T. and Maloji, S., 2019., "A Novel Reconfigurable Bandpass Filtering Antenna for IoT Communication Applications.", Progress in Electromagnetics Research, 96, pp.13-26.
- [23] Krishna, A.V., Madhav, B.T.P., Avinash, R., Koukab, 2019. "A novel h-shaped reconfigurable patch antenna for IoT and wireless applications" (2019), International Journal of Innovative Technology and Exploring Engineering, 8 (7), pp. 1757-1764.
- [24] Vamseekrishna, A. and Madhav, B.T.P., 2018. "A frequency reconfigurable antenna with Bluetooth, Wi-Fi, and WLAN notch band characteristics.", International Journal of Engineering and Technology, 7(2.7), pp.127-130.
- [25] Vamseekrishna, A., Madhav, B.T.P., Anilkumar, T. and Reddy, L.S.S., 2019. "An IoT Controlled Octahedron Frequency Reconfigurable Multi-band Antenna for Microwave Sensing Applications.", IEEE Sensors Letters, 3(10), pp.1-4.

# A Critical Analysis of IS Governance Frameworks: A Metamodel of the Integrated use of CobiT Framework

Lamia MOUDOUBAH<sup>1</sup>, Abir EL YAMAMI<sup>2</sup>, Mansouri KHALIFA<sup>3</sup>, Mohammed QBADOU<sup>4</sup>

SSDIA Laboratory, ENSET Mohammedia  
Hassan II University of Casablanca  
Mohammedia, Morocco

**Abstract**—Information Systems Governance (ISG) is an essential component of corporate governance. It refers to the implementation of the means of decision-making. A considerable number of studies on information systems governance (ISG) have been published. Nevertheless, there is a need to conceptualize and model this theoretical context. The aim of this paper is to provide a study of frameworks that integrates this domain as well as to bring a modeling of the concepts that structure the framework of this domain and a profound and clear understanding of the IS process, IS governance has been studied as a concept. The results demonstrated that the adoption of the COBIT repository in the organization could amplify its efforts. This input therefore enables the organization to capitalize on and build up knowledge in the field of IS governance, and to propose models for delivering an integrated, business-aligned IS.

**Keywords**—Information Systems Governance (ISG); IS process; business process; COBIT

## I. INTRODUCTION

This paper focuses on the area of information systems governance. It corresponds to the implementation of the ways and means by which stakeholders can ensure that their concerns are taken into account in the operation of the information system (IS).

According to [1] IS management thus aims to define the objectives assigned to the information system and to plan, define and implement the processes related to IS lifecycle management.

These activities are based on the control and measurement of the performance of these processes with respect to the objectives underlying the use of the IS [1]. The object of IS governance is therefore the Information System [2]. The mission of an IS is to make the main activities of the organization generate more added value. It takes advantage of computer technologies (memorization, communication, calculation, transformation, and presentation) to establish a network of coordination between the organization's activities as well as a network of cooperation between the organization's actors.

In this paper, authors address a twofold question in order to answer, on the one hand, the choice of good practice frameworks for IS governance, and on the other hand, the research gaps in the formalization and conceptualization of the IS object that is the IS process.

This work presented as follows. In the first section, authors briefly present a repository of good IS governance practices. In the second section, authors present the proposed model for the conceptualization of ISG, by explaining the concept of ISG, clarifying the perimeters of ISG and modeling the ISG process. In the third section, authors defined the place of COBIT in the ISG, and then they proposed a model for the conceptualization of COBIT. Finally, Discussion of this work to sum up with a conclusion.

## II. THEORETICAL BACKGROUND AND MOTIVATION

### A. Benchmarks of Good Practices of ISG

According to [4], Standards and benchmarks of good practices in ISG is relatively little studied in the academic literature. However, the last few years have been marked by an increase in the number of these good practices, each coming from a professional community with its own issues and its own culture. The professional literature offers all kinds of books, catalogues and guides with comments on the use and fields of application of good practices [5], [6], [7], [8]. The reading of these documents shows a context rich in knowledge about the content and orientations of these standards.

According to [9], the notions of "standard" and "benchmark of good practice" are only two sides of the same coin. Their common denominator lies in their willingness to serve as a model or reference system recognized by a competent body and disseminated to a wide public. The authors retain the following characteristics relating to these two concepts in the table (Table I) [9]:

TABLE I. CHARACTERISTICS RETAINED FOR THE CONCEPT OF GOOD PRACTICES [9]

| Concept    | Features   |
|------------|--|
| Standards  | <ul style="list-style-type: none"><li>- A document established by consensus and approved by a formal standards body.</li><li>- Provides rules, or characteristics, for activities or their results.</li><li>- Defines an optimal requirement level to be achieved.</li><li>- Is a public statement because of its official origin?</li></ul> |
| References | <ul style="list-style-type: none"><li>- Document established and approved by a profession.</li><li>- Contains a set of recommendations.</li></ul>  |

### B. Existing Standards and Benchmarks

Existing standards and benchmarks, considered as operational solutions, can be summarized from the following list, and depending on their use, can be divided into several domains in the table below (Table II):

- Information System Development:

**CMMI** (Capability Maturity Model Integration) and maturity levels. It is a model for evaluating processes during the design of software or applications [3].

**UML** (Unified Modeling Language), a unified modeling language. It is a development tool allowing modeling a problem in a standard way. It is the reference in terms of object modeling [3].

**SPICE** (Software Process Improvement and Capability determination). Standard for software process evaluation, synthesis of software process evaluation and improvement approaches. Essentially, it includes an implementation guide for the evaluation of software development projects [3].

- Information System Management:

**ITIL** (Information Technology Infrastructure Library) offers a structured library of best practices for a better management of the Information System [3].

**Norme BS 15000**: Guide to good practice for supply and service management. It is associated, for its implementation, with ITIL recommendations [3].

- Management and organization of the Information System:

**COBIT** (Common Objectives for Business Information Technology). This method was developed by ISACA (Information Systems Audit and Control Association) about ten years ago [3].

- Project Management:

**PRINCE 2**: Projects IN Controlled Environments is a structured project management and certification method that focuses on three points: project organization, management and control [3].

**PMBOK**: Project Management Body of Knowledge. It is the reference document for project management. It describes knowledge and methods applicable to the majority of projects, whether IT or not, on which there is a consensus on their value and usefulness [3].

**PPM**: Project & Portfolio Management. Management of projects so that they can be considered as portfolios. A strategy allows organizations to align their IT application development projects and resources with business objectives by putting in place indicators to monitor these projects [3].

- Information System Security:

**ISO 27001**: This standard allows companies to validate the security practices they adopt for their Information System [3].

**ISO 15408/16949**: IT security management, common criteria. They define the procedures and standard technical measures to be considered in the life cycle of a software product [3].

- Company management and quality:

**COSO** (Committee Of Sponsoring Organizations): is to manage business risks [3].

**ISO 20000** and organization certification: this standard defines the needs of service management within the framework of the Information System. It defines the main processes for the efficient provision of these services [3].

**ISO 9001**: quality assurance model used for the certification of quality management systems [3].

**ISO 10006**: This standard provides guidance on the application of quality management to projects as part of project management processes [3].

**eSCM (e-Sourcing Capability Model)**: it is a repository presenting good practices in the client/provider relationship in the context of outsourcing services [3].

### C. Objectives of these Methods

These main references are complementary Associated; they bring value to the processes of the Information System and a fortiori to the whole organization, based on four main objectives [3]:

1) The implementation of good practices in the management of the services provided by the Information System.

2) The establishment of a development strategy for these processes including indicators related to budgets and projects.

3) The guarantee of a good organization (management, supervision) of the assets (hardware, software) and technologies implemented.

4) The alignment of the Information System with the strategy of the company on its core business, the requirements of regulations related to professional particularities.

TABLE II. RANKING OF THE MAIN REPOSITORIES IN TERMS OF USAGE AND BY ISD DOMAIN (SOURCE CIGREF)

| Order | Name of the Repository       | ISD Domain                                |
|-------|------------------------------|---|
| 1     | ITIL                         | Production                                |
| 2     | ISO 27001                    | Security                                  |
| 3     | Nomenclature RH du CIGREF    | Competency management                     |
| 4     | COBIT                        | Governance                                |
| 5     | CMMI                         | Development                               |
| 6     | PMBOK                        | Project Management                        |
| 7     | ISO 9001                     | Quality Management                        |
| 8     | Benchmarking of CIGREF costs | Costtracking                              |
| 9     | TOGAF                        | Customer-supplier relationship management |
| 10    | PRINCE 2                     | Project Management                        |
| 11    | eSCM                         | Customer-supplier relationship management |

### D. COBIT

The COBIT model (Control Objectives for Information and related Technology) presented as a model for governance and control in information technology [3].

Created by ITGI (IS governance Institute) and ISACA (Information Systems Audit and Control Association), COBIT has been adopted by many international companies. However, because of its concept, it is preferably implemented in large companies [3].

Indeed, it is mainly aimed at managers and auditors who may be involved in providing a methodology for [3]:

- 1) Corporate management. This framework helps them to control investments in order to better manage risks and meet their obligations to investors and shareholders.
- 2) The IT managers in charge of managing the Information System and the services provided.
- 3) The auditors, as they can make recommendations to management on the internal control of Information Systems.

The fact that this standard is intended for large companies does not prevent the implementation of processes adapted for small companies.

The methodology presented may contain improvement ideas for the governance of their Information System. CobiT is a set of recommendations and processes for evaluating IS resources. It is intended to guide practitioners in the implementation of internal controls.

CobiT was developed in 1994 (and published in 1996) by ISACA (Information Systems Audit and Control Association). ISACA has been represented in France since 1982 by AFAI (Association Française de l'Audit et du Conseil Informatiques).

CobiT is a control framework that aims to help management to manage risks (security, reliability, compliance) and investments [1].

### III. PROPOSED METAMODEL OF ISG

As noted earlier, information systems governance (ISG) is a goal-driven project management activity that is driven by the execution of a process. This observation allows us to consider a representation of ISG as a whole made up of a product, describing the system of concepts that underlies ISG, and a process that aims to change the context of ISG [1].

In addition, any system can be directed and controlled provided that it can define (i) the devices for measuring whether the objectives assigned to it are being achieved and, if not, (ii) the levers (variables) of action for correcting deviations [1].

Governance is therefore first, and foremost a matter of making decisions in the face of uncertainty. The mediation of decisions to be taken and the resulting actions is mediated by a decision-maker driven by the desire to move towards the target assigned to the project or project portfolio [1].

In this section, authors will present a conceptual model of GSI, the objective of which is to describe the conceptual system underlying GSI. This work is done to overcome the inadequacy in the conceptualization of IS management and to build an IS of governance. Proposed model is based on observation and analysis of the literature.

#### A. ISG Concept

IS Governance can be reduced to a simple approach based on good practices inspired by standards and reference frameworks. However, it leads to ambiguity, misunderstanding regarding the notion of IS Governance and the respective roles of management on the one hand, and service governance on the other.

In order to understand the place of IS Best Practice Standards and Reference Frameworks in IS Governance processes, authors believe it is necessary to clarify the meaning and scope of IS Governance. This will allow them to understand a posteriori the actual role of the IS Best Practice Standards and Reference Frameworks, in relation to IS management and IS governance.

According to [9], [10], [11] and [12], the concept of IS governance often referred to as IS governance in specialized language, is a relatively new concept, emerging from several disciplines, including the social and information sciences.

Following an analysis of the work as a whole, authors deduce that there is a lack of agreement between these different authors ([13], [14], [15], [16], [16], [17], [18], [19], [20], [21], [22], [36], [37], [38], and [39]) on the definition of the concept of IS governance.

According to [23], there are several definitions of the concept of IS governance on the web. In order to understand the exact meaning of the concept of IS governance; the author proposes to return to the notion of corporate governance, often referred to "corporate governance". According to CIGREF (2002), the transposition of IS governance from corporate governance presupposes a good understanding of the principle of separation between "owners" and "managers". The implementation of this principle at the IS level presupposes the existence of a control body independent of the ISD, responsible for reducing the gaps between the decisions taken by those in charge of the IS (managers) and the interests of the owners (business and functional departments) [24].

Following the example of these excerpts, it is worth noting that IS governance, as a subset of the principles transposed from corporate governance, aims to strengthen the overall consistency of IS decisions with the interests of stakeholders.

There is a considerable gap between actual IS governance practices in companies and theoretical approaches, according to [12] and [14]. This is mainly due to the common confusion in the professional community between the respective roles of management and governance (inferred in [23]).

#### B. Perimeters of IT Service

IS Governance stems "from initiatives for strategic alignment with the expectations of managers and the business processes from which the principles of business governance result", according to [25]. This should not be confused with two closely related sub-domains, namely, IS Governance and infrastructure governance.

To clarify the scope of IS Governance, and according to ITGI [26], "IS governance is the responsibility of the board of directors and executive management. It is an integral part of



enterprise governance and consists of the leadership and organizational structures and processes that ensure that the organization's IT sustains and extends the organization's strategies and objectives" [26]. Consequently, the IS Governance body as a supervisory body exogenous to the IS function must define under the responsibility of the supervisory bodies.

The framework and processes that support the company's strategy while respecting the objectives of corporate governance [27].

All this analysis by the authors cited above leads us to conclude that ISG is based on the implementation and management of a set of processes that are modelled on the objectives of corporate governance. Normally, these processes are intended to support the objectives relating to the following areas:

- S.A: The strategic alignment of the IS with the business;
- R.M: Risk management;
- V.C: Value creation;
- R.M: Resource management;
- P.M: Performance management;

This study pushes us towards the conceptualization of a model, which models ISG as a concept (Fig. 1):

### C. ISG Process

IS governance is based on a set of processes that make it possible to control that the objectives assigned to the IS are properly considered and to react if necessary.

[21] Proposes to consider the IS processes that are essential for IS management around a control process (reporting) and an action process for decision-making. It is in line with the idea developed earlier in [28], which recommends six steps for aligning business and IT. They mainly concern identification of objectives, understanding of alignment links, analysis (in-fine, measurement and control) and prioritization of gaps, specification and choice of actions to be taken.

The IS processes that the authors consider are thus linked to the achievement of IS quality by a control mechanism based on the generic Deming approach of the PDCA (Plan, Do, Check, Act) [28].

The PROCESS SI facet allows this aspect to be represented. The values associated with this facet measure the degree of control of these processes based on the principle that an IT PROCESS is at least documented. The identification of metrics, indicators and control rules allows decision making on the audit process: the process is then steered. An evaluative process is a process under control whose evolution has been considered and which is representative of mature governance.

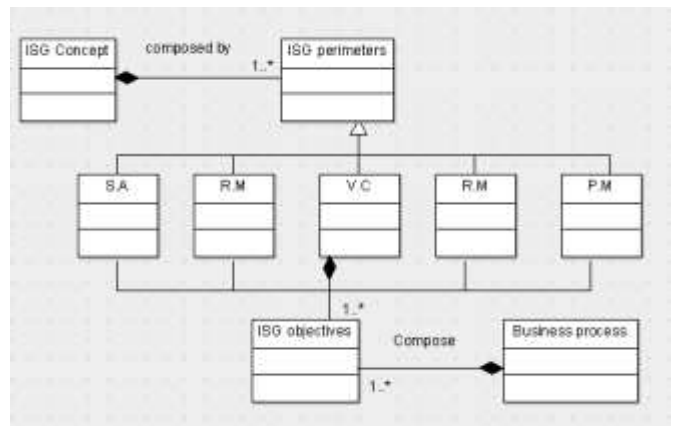


Fig. 1. Metamodel of ISG Concept.

Do not confuse IS processes with business processes; IS processes are essential for IS management around a control process (reporting) and an action process for decision-making [28]. Thus, the IS PROCESSES that are essential within the framework of good governance are those dedicated to audit, control and reporting according to [29].

While the business process is defined in [30] as "a structured and measured framework of activities designed to produce a specific output for a customer or market. This implies focusing on how work is done within an organization, rather than focusing on the product.

A process is therefore a precise order of activities across time and space, with a beginning and an end, clearly defined inputs and outputs: a structure of action." [30].

The typologies of business processes are defined in several ways in the previous works; authors will clarify typologies of business processes by quoting:

RUMMLER's article [31]: According to his approach, he distinguishes primary processes, which are in direct contact with the customer and directly generate value, from supporting processes. The support processes are invisible from the customer's point of view and are functional: they concern accounting, recruitment or technical support. The primary processes concern activities and operations dedicated to procurement, production and sales.

ALONSO's article [32]: His approach is based on the nature of the business process. It distinguishes four types of processes:

- Productive: The process is repeatable and implements the primary processes of the company.
- Administrative: The process is bureaucratic and is governed by clearly established rules.
- Collaborative: The process is characterized by important interactions between actors. This is the case, for example, with steering committee processes.
- Ad-hoc: The process is defined on the fly during its execution. It is a process that is not planned, it is often linked to exceptions.

Authors' study leads us to conceive the IS governance process across the domains or in other words the perimeters of IS management, from which the objectives of value creation derive from the strategic alignment of the IS with the business while risk management derives from the control and accountability policies in the company. The whole is supported by resources, and managed with the aim of achieving the desired performance.

This conclusion gives rise to the following Metamodel (Fig. 2).

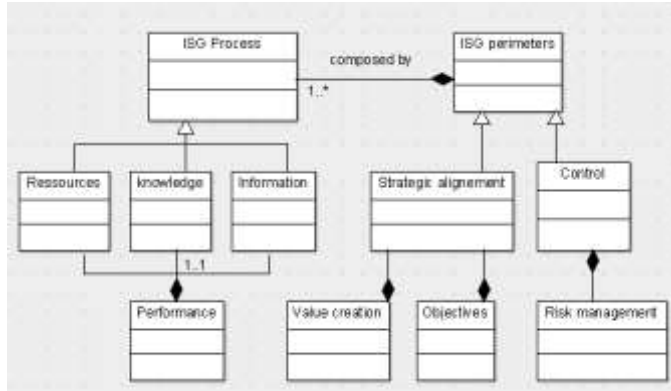


Fig. 2. ISG Process Metamodel.

#### IV. COBIT AT CORE OF ISG

Performance is at the heart of ISG concerns. It is the result of mastering the maturity of business and IT processes. Also the application of methods oriented by process maturity such as COBIT [34], [33].

Authors' thesis topic is about the COBIT repository, so after this study, researchers of this paper will focus on COBIT. In fact, authors' paper don't underestimate the value of the other standards. However, I take COBIT because it indicates the main lines to follow, the main axes to have a good ISG. For example, for the "Plan and organize" axis, COBIT tells you that you need to define a strategic IT plan aligned with the company's strategy, then for "acquire and implement" that you need to put in place solutions, infrastructure and processes that are consistent with this plan. Then that you need to define service levels, ensure a level of security to manage risks, train employees, etc. and finally that you need to ensure effective control of IT processes to guarantee a level of reliability, security, compliance and confidentiality. All this is based on strategic alignment: aligning this entire cycle with the company's objectives.

##### A. COBIT Proposed Metamodel

The CobiT repository is structured by components on which a conceptualization process will be applied. In this part, researcher's paper describe these components and propose a conceptual model to show the concepts of CobiT.

**CobiT** refers to four Generic Process Areas. Each contains the processes audited by the CobiT approach and refers to a stage of the governance cycle: Plan and Organize, Acquire and Implement, Deliver and Support, and Monitor and Evaluate.

In total, CobiT **includes 34 processes** (COBIT Process) that meet five IS governance requirements (Domain of ISG).

A process is audited according to information criteria (Information Criterion) against a set of control objectives (Control Objective). It is analyzed according to its level of maturity, which is representative of its effectiveness and efficiency.

According to CobiT a process uses resources in terms of skills, information, applications and infrastructure (IT Resource), and requires input and output information elements (Element, Input, Output).

A process organizes Activities during which actors intervene in accordance with their functions and responsibilities (Role). CobiT proposes a RACI grid (Responsible, Accountable, Consulted, and Informed) which allows visualizing the responsibilities of each person in relation to the activities. For a particular activity, an ISD can be responsible (R), accountable (A), consulted (C) or simply informed (I) [1].

The means of control proposed in CobiT meet control objectives. They implement a set of metrics allowing judging the achievement of the control objective. **A control objective** is defined in relation to the business goals and IT goals which are the objectives that stakeholders set for themselves within the framework of IS management processes.

In general, CobiT processes meet a set of **28 goals (ButCOBIT)**. Indicators (COBIT Indicator) measure the level of achievement of the goals.

This analysis led us to apply a conceptualization process, and to describe the whole study of the COBIT product in the following metamodel (Fig. 3):

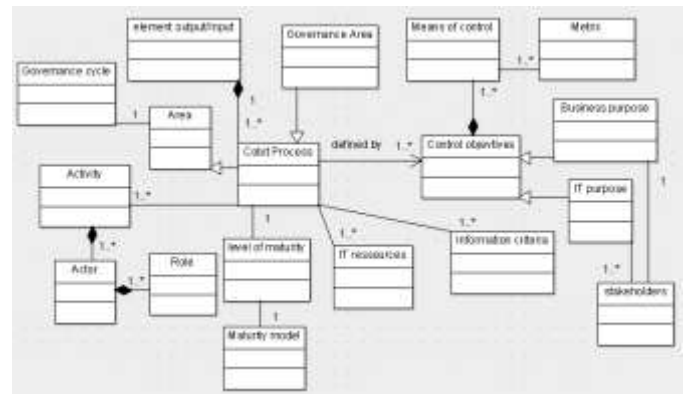


Fig. 3. COBIT Proposed Metamodel.

#### V. OSTERLE PRINCIPLES

In order to differentiate scientific research from solutions designed by practitioners, Osterle [35] indicates that scientific research must be marked by abstraction, originality, justification and benefit.

1) Abstraction: This paper clarifies the notions that characterize the field of ISG and proposes a metamodel to determine the place of COBIT in the conceptualization of ISG.

- 2) Originality: The proposed metamodel is not present in the body of knowledge of the domain.
- 3) Rationale: The proposed method for evaluating the model must justify the model.
- 4) Advantage: The COBIT framework allows a better conceptualization of the ISG and guarantees a better IT management for the company that adopts it.

## VI. DISCUSSION

ISG includes the entire management system (processes, procedures, organization) used to steer IT. This concern is an expression of the desire to ensure corporate governance.

There are a large number of repositories that reflect the best practices, developed over the years. This may come as a surprise. The reality is that each of them starts from a particular concern: safety, quality, services offered to customers, auditing, project development, etc. [33].

This is unavoidable for each function to recognize itself in its own practices. At the same time, the question arises of setting up a single, global framework for the IT department that meets all expectations [33].

CobiT positions itself as both an audit reference and a governance reference. In terms of governance, it is immediately in line with the company's business lines and strategy. Beyond this positioning, CobiT is designed, developed and continuously improved to federate all IT-related repositories.

As a repository for information systems governance, the scope of CobiT goes beyond the scope of information systems management to encompass all the stakeholders in the company's information systems.

Indeed, implementing the ISG processes is not an easy task, as its definition and concepts are not clear. In this context, this work aims to provide a global approach for the conceptualization of the ISG and a benchmark of good practices in this field.

Even though the number of researches dealing with the conceptualization of the ISG is increasing, there is no study that models the concept of ISG in a way that identifies the interesting role of the CobiT at the heart of this field.

It is therefore mandatory to build a shared representation of ISG concepts and to show how these concepts are structured within the CobiT framework.

The objective is to strengthen the professional literature by providing a machine-readable document for the ISG domain model. Then to the scientific literature that is interested in improving information systems governance frameworks by improving the understanding of the CobiT architecture.

Similarly, the main objective of the proposed metamodel, is to represent the ISG domain concepts, their properties, and relationships, to build a shared representation of ISG concepts between researchers and practitioners, to show how these concepts are reinforced by the CobiT framework, to make ISG knowledge reusable in similar IS engineering and

management situations and to support the creation of new ISG models.

## VII. CONCLUSION

In this article, authors have proposed a framework for the analysis of information systems governance (ISG), starting with a study of information systems standards and repositories, showing the link of these standards and repositories with the ISG. Then proceeding to the conceptualization of the ISG by proposing metamodel, then the ISG process, and finally the conceptualization of COBIT in order to highlight the need for research on the globality of the ISG. This work contributes, confirms and proposes a plus on the subject of IS governance.

### REFERENCES

- [1] Bruno Claudepierre, "Conceptualisation de la Gouvernance des Systèmes d'Information : Structure et Démarche pour la Construction des Systèmes d'Information de Gouvernance," Paris, 2012.
- [2] Cigref, "Gouvernance du système d'information," 2002.
- [3] SUPINFO International University, "principauxreferentielsgouvernanc esysteme information," 2016.
- [4] Randa Ben Romdhane, "Les effets de la multiplicité des normes et des référentiels de bonnes pratiques : le cas de la Direction des Systèmes d'Information," Conservatoire national des arts et métiers - CNAM, Français 2015.
- [5] G.Teneau, J.Ahanda, "Guide commenté des normes et référentiels," Editions d'Organisation, 2011.
- [6] CIGREF, "Les emplois\_métiers du SI dans les grandes entreprises," 2009.
- [7] G. T. B. Axel Hochstein, "Service Oriented IT Management: Benefit, Cost and Success Factors," in European Conference on Information Systems, 2005.
- [8] C. & C.-S. A. Pollard, "Justifications, Strategies, and Critical Success Factors in Successful ITIL Implementations in U.S. and Australian Companies: An Exploratory Study," Information Systems Management, 2009.
- [9] J. Iden, "Setting the stage for a successful ITIL Adoption; A delphi study of IT experts in the Norwegian armed forces," Information systems management, 2010.
- [10] Cigref, "gouvernance du SI problématiques et démarches," 2001-2002.
- [11] Bounfour et Epinette, "Valeur et Performance des SI: une nouvelle approche du capital immatériel," 2006.
- [12] Peter David Weill, Jeanne W. Ross, "IS governance: How Top Performers Manage IT Decision Rights for Superior Results", Reviewed by Lester P. Diamond, U.S. Government Accountability Office, USA. International Journal of Electronic Government Research, 1(4), 63-67, 2005.
- [13] Young, "An introduction to IT service management," 2004.
- [14] Denise Ko, Dieter Fink, "information technology governance: an evaluation of the theory practice gap," Corporate Governance, 2010.
- [15] M. B. T. a. K. H. Schermann, "Explicating Design Theories with Conceptual Models: Towards a Theoretical Role of Reference Models," 2009.
- [16] Candida C. Peterson., "Theory of mind development in oral deaf children with cochlear implants or conventional hearing aids," 2004.
- [17] Mårten Simonsson and Pontus Johnson KTH, "Assessment of IS governance - A Prioritization of Cobit -", Royal Institute of Technology Osquidasväg 12, 7 tr, S-100 44 Stockholm, Sweden.
- [18] Pontus Johnson & Mathias Ekstedt, "The Effect of IS governance Maturity on IS governance Performance," journal information systems management. Vol. 27, 2010.
- [19] Amrik S. Sohal & Paul Fitzpatrick, "IS governance and management in large Australian organisations," International Journal of Production Economics, Vol. 75, Issues 1-2, 2002, pp.97-112.

- [20] Peter Weill, Jeanne W. Ross, "IS governance: How Top Performers Manage IT Decision Rights for Superior Results," 2004.
- [21] Wim Van Grembergen, Steven De Haes, "Structures, Processes and Relational Mechanisms for IS governance," *Strategies for Information Technology Governance*, pp. 36, 2004.
- [22] Willson, Phyl, Pollard, Carol, "Exploring IS governance in Theory and Practice in a Large Multi-National Organisation in Australia," *Journal Information Systems Management*, vol.26, 2009.
- [23] Eric Fimbel, "Pour un système d'information synchrone," *L'Expansion Management Review*, pp.114-129, 2007.
- [24] Cigref, "Comment le contrôleur de gestion peut-il assister le DSI," 2001-2002.
- [25] Georges Epinette, "Valeur et performance des SI Une nouvelle approche du capital immatériel de l'entreprise," 2006.
- [26] H. Simon, "The Science of the Artificial," MIT Press, 19.
- [27] S.D.Haes ; W.V.Grembergen, "Analysing the Relationship between IS governance and Business/IT Alignment Maturity," *Proceedings of the 41st Annual Hawaii International Conference on System Sciences (HICSS 2008)*.
- [28] Luftman, J., Papp, R., & Brier, T. "Enablers and Inhibitors of Business-IT Alignment. *Communications of the Association for Information Systems*," 1999.
- [29] Riadh Manita, "La qualité du processus d'audit : une étude empirique sur le marché financier tunisien. La place de la dimension européenne dans la Comptabilité Contrôle Audit", France, 2009.
- [30] Thomas H. Davenport, "Process Innovation: Reengineering Work Through Information Technology," Harvard Business Press, 1993.
- [31] P. Brache, A. Rummier, "Invited reaction: Performance improvement: A methodology for practitioners," 1995.
- [32] CIGREF, "Referentiels\_de\_la\_DSI\_CIGREF\_2009.pdf," 2009.
- [33] Dominique Moisand, Fabrice Garnier de Labareyre, "COBIT pour une meilleure gouvernance des systèmes d'information," 2009.
- [34] T. Ravinchandran, C. Lertwongsatien, "Effect of information systems resources and capabilities on firm performance: A resource-based perspective," *Journal of Management Information Systems*, pp.237-276, 2005.
- [35] H.Österle, J.Becker, U.Frank, T.Hess, D.Karagiannis, H.Krcmar, "Memorandum on Design-Oriented Information Systems Research," *European Journal of Information Systems*, EJIS, pp. 7-10, 2011.
- [36] A. El Yamami, K. Mansouri, M.Qbadou, "Multi-objective IT project selection model for improving SME strategy deployment," *International Journal of Electrical and Computer Engineering* 8 (2), 1102, 2018.
- [37] K Benmoussa, M Laaziri, S Khouilji, KM Larbi, A El Yamami, "Enhanced model for ergonomic evaluation of information systems: application to scientific research information system," *International Journal of Electrical and Computer Engineering* 9 (1), 683, 2019.
- [38] Moudoubah, L., Mansouri, K., & Qbadou, M., "Towards an Ontological Analysis of the Alignment Problems of Fields in the Architecture of an Information System," 2020.
- [39] Moudoubah, L., Yamami, A.E., Mansouri, K., & Qbadou, M., "Towards the implementation of an ontology based on COBIT framework (CobitOnylogy)," *1st International Conference on Smart Systems and Data Science (ICSSD)*, 1-6, 2019.

# Prioritization of Software Functional Requirements from Developers Perspective

Muhammad Yaseen<sup>1</sup>, Aida Mustapha<sup>2</sup>, Noraini Ibrahim<sup>3</sup>

Faculty of Computer Science and Information Technology  
Universiti Tun Hussein Onn Malaysia, Parit Raya  
86400 Batu Pahat, Johor, Malaysia

**Abstract**—Prioritizing software requirements is important and difficult task during requirements management phase of requirements engineering. To ensure timely delivery of project, software developers have to prioritize functional requirements. The importance of prioritization increases when size of requirements is big. Software for large enterprises like the Enterprise Resource Planning (ERP) systems are more likely to be developed by a team of software developers where large size requirements are distributed in parallel team members. However, requirements are dependent on each other, therefore development of pre-requisite requirements must be carefully timed and should be implemented first. Therefore, assigning importance and priority to some requirements over others is necessary so that requirements can be available on time to developers. This paper proposes a prioritization approach for functional requirements on the basis of their importance during implementation. The design of research method consists of Analytical Hierarchical Process (AHP) technique based on spanning trees. Through spanning trees, dependent requirements were linked in hierarchical structure and then AHP were applied. As a result of prioritization, requirements were distributed in such a way that dependency among requirements of developers were kept minimum as much as possible so that waiting time of requirements for their pre-requisite were reduced. With reduced effect of dependency in requirements of parallel developers, timely delivery of software projects can be assured.

**Keywords**—Requirements prioritization; Functional Requirements (FRs); directed graph; spanning tree (ST); Analytical Hierarchical Process (AHP)

## I. INTRODUCTION

Requirements Engineering (RE) is a systematic way of collecting software requirements [1][2][3]. There are different types of software requirements [4][5][6]; Business Requirements (BRs) that deal with benefits of implementing requirements, Process Requirements (PRs) that deal with time and cost issues during development, Functional Requirements (FRs) that deal with the actual functionalities of the software, and finally Non-Functional Requirements (NFRs) that deal with requirements such as usability, security, and performance. The collected FRs need proper management in determining issues such as which requirements should be given higher priority, which team member will implement a particular requirement, when the requirements is expected to be delivered, how will the requirements be integrated and other concerns related to requirements management [7][8].

Requirements Prioritization (RP) is a task in RE that focuses on giving priority or ordering a group of requirements [9][10]. Techniques such as cost-value ranking, attribute goal-oriented, and value-oriented approaches work well for BRs in combination with high level FRs [11][12]. FRs are prioritized either from client's perspective or developer's perspective [13][14]. FRs from client's perspective are normally high level requirements that are also known as user requirements (URs). Techniques like the Analytical Hierarchical Process (AHP), binary trees, Genetic Algorithm (GA) are more suitable to prioritize FRs from user perspective [15][16][17]. Meanwhile, techniques like Quality Function Deployment (QFD) and contextual preference-based technique are suggested for prioritizing NFRs [18][19]. Although most of the techniques like AHP work well for small size requirements, they are not scalable and suitable to apply on large requirements. While machine learning techniques and intelligent based techniques such as Artificial Neural Networks (ANN) and SNIPR are suitable for prioritizing large-sized FRs, but they are not suitable techniques to prioritize FRs from developer's perspective where requirements are distributed in parallel development team [20][21][22].

As FRs are not isolated but inter-related so prioritization of FRs is necessary especially when parallel team members are assigned to implement the entire requirements. Giving importance and priority to some requirements over the others is necessary so that pre-requisite requirements can be available on time for other requirements. According to [23], successful projects are not only those that meet all their FRs and NFRs but timely delivery of these requirements is also necessary. Most of big size software's fail to deliver in time, thus proper management and prioritization of FRs from developer's perspective is necessary for successful implementation and delivery of any software project [24].

Although the current prioritization techniques are able to prioritize FRs from user perspective effectively in selecting particular modules or requirements, the same techniques are not either capable or applied to prioritize FRs from developer's perspective when it involves the internal structure and dependency of one requirement on others. Another problem is that most techniques are suitable for prioritizing small-sized requirements but not scalable for large set of requirements. Therefore, a new prioritization is needed for focusing on prioritizing FRs from developer perspective

within the setting of large size requirements especially in parallel developing projects.

Technique like AHP can be applied with pre-defined prioritization rules to FRs but it is not scalable for big-sized requirements. However, we can use technique such as AHP that pairwise compare requirements to prioritize requirements from developer's perspective.

To address this gap, this research work proposes a new approach to prioritize FRs using AHP but based on spanning trees, called the SAHP. The proposed prioritization approach will then be evaluated on FRs of ODOO ERP as case study. Finally, this paper will also investigate the scalability of SAHP in ERP systems by comparing time complexity of the SAHP with existing AHP. The remaining of this paper proceeds as follows. Section 2 presents preliminary studies related to AHP. Section 3 presents the proposed AHP based on Spanning Trees called the SAHP. Section 4 reports evaluation of prioritization experiments using requirements of ERP system. Section 5 presents efficient distribution of requirements in parallel team members and finally Section 6 concludes with some indication for future work.

## II. BACKGROUND STUDY

Analytical Hierarchical Process (AHP) is an organized decision-making method that is intended to compute complex multi-criteria decision problems. AHP is technique that is also applied efficiently in many other fields such as biology and social sciences for prioritization. In fact, AHP is the utmost frequently discussed prioritization technique within decision making in requirements engineering. AHP is led by comparing all possible pairs of hierarchically categorized entities such as requirements as well as stakeholders for obtaining comparative priorities for all objects [15].

Research in [25] revealed that AHP is capable of improving total time of calculations for pairwise comparisons of the requirements by using eigenvalues and matrix evaluation. The research also proposed Consistency Index (CI) to remove errors like inconsistency. Basically the requirements are arranged in groups called bins in the form of hierarchy. This form of prioritization although be helpful in those cases where requirements are not too much and we need to prioritize with the help of AHP. Number of comparisons will be less as compared to traditional AHP but still it fails large set of requirements.

According to [26], although we assign priorities to FRs, we can also assign priorities on the basis of PRs. The work discussed prioritization of PRs by considering both local priority and perspective priority and proposed the Correlation-Based Priority Assessment (CBPA) that prioritizes requirements from different stakeholder perspectives while to highlight the key issues among them. Two types of requirements were considered (1) from business point of view and (2) from management point of view. Increased profit, lead in competition, reduced cost of development, reduced time to development are business-oriented process requirements while maintaining a project within budget, on schedule, high customer satisfaction, increase productivity are management-oriented process requirements that are considered and

prioritized in the research work by author. The relationship between different requirements, its prioritization and impact are discussed in the paper in the form of matrix. Apart from PRs, prioritization of requirements from multiple stakeholder's point of view is also discussed. High priority requirement needs more attention and leads to project success [26].

Apart from fully AHP-based solutions to prioritization of requirements, intelligent-based solution has also been proposed for prioritization of requirements collected from stakeholders by applying machine learning techniques to first group similar requirements, and then apply Artificial Neural Networks (ANN) for further prioritization. Finally, AHP was applied at the end for final comparisons. In first step, before clustering, stakeholders are requested to prepare requirements, then on the basis of profiles of stakeholders and through expert opinions using ANN, requirements can be prioritized [22].

Along with stakeholder preferences, it is also necessary to have prioritization which can handle dependencies in between requirements from user perspective. DRANK is an automated algorithm was presented to perform comparisons based on the importance of dependent requirements and compared the results with AHP and other techniques. Experiments proved that this technique is more efficient and scalable for large size URs [27].

Though many authors have used AHP and tried to reduce number of comparisons from different perspectives, AHP are still unable to cater prioritization of FRs during an active implementation software life cycle. Existing AHP implementation needs user input for pairwise comparison of requirements, while we need this process to be automatic i.e. to take input from its internal structure rather than user. The purpose of this study is to reduce this research gap to prioritize FRs from developer's perspective.

## III. PROPOSED AHP BASED ON SPANNING TREE (SAHP)

This section proposes spanning trees based approach to represent FRs and then prioritized with AHP. Spanning tree represents hierarchal order and dependencies of all inter-related requirements. From spanning tree, one can easily pairwise compare requirements with AHP. FRs collected from any sources using appropriate elicitation technique and must be specified in the form of Software Requirement Specification (SRS). In this research, the FRs are represented as alphabets R1, R2, ..., Rn and are enclosed in circles as nodes.

### A. Spanning Trees

In graph theory, a spanning tree is a subset of graph. A graph  $G = (V; E)$  consists of finite set of vertices  $V$  and finite set of edges  $E$ . Edge is something that connects two vertices. Graphs are useful for the representation of any kind of data in particular sequence [28][29]. This research uses directed acyclic graphs (DAG) rather than cyclic graphs. Requirements are represented as vertices and arrows in the graph indicates the dependency of a requirement on another requirement. The requirement generates arrow and points to another requirement indicating that it is necessary or required for

another requirement. For example,  $R1 \leftarrow R2$  indicates that R1 is depended on R2 or R2 is required for the completion of R1. Given the requirements collected, Fig. 1 shows the graphical representation of requirements through DAG. Cycles in requirements are not possible because if one requirement is needed for the implementation of other requirement than opposite is not possible e.g. if R1 is required for R2 and R2 is required for R3 than it is not possible that R3 will be required for R2 and R1. Graph based approach is also used in one of our previous research study to related FRs [30].

**Spanning trees** are special graph that have several important properties. First, if T is a spanning tree of graph G, then T must span G, meaning T must contain every vertex in G. Second, T must be a sub graph of G. In other words, every

edge that is in T must also appear in G. Third, if every edge in T also exists in G, then G is identical to T [31]. Spanning trees can be formed simply either by performing breadth-first search (BFS) or depth-first search (DFS) or it can be formed directly from adjacency matrix. Because spanning trees use graph-based search algorithms that are only dependent on the number of vertices in the graph, the algorithms are considerably fast [32][33]. The general properties of spanning trees are as follows.

The resulting spanning trees from graph of Fig. 1 are shown in Fig. 2. From a spanning tree, one can easily see the need of particular requirement in relation to other requirements.

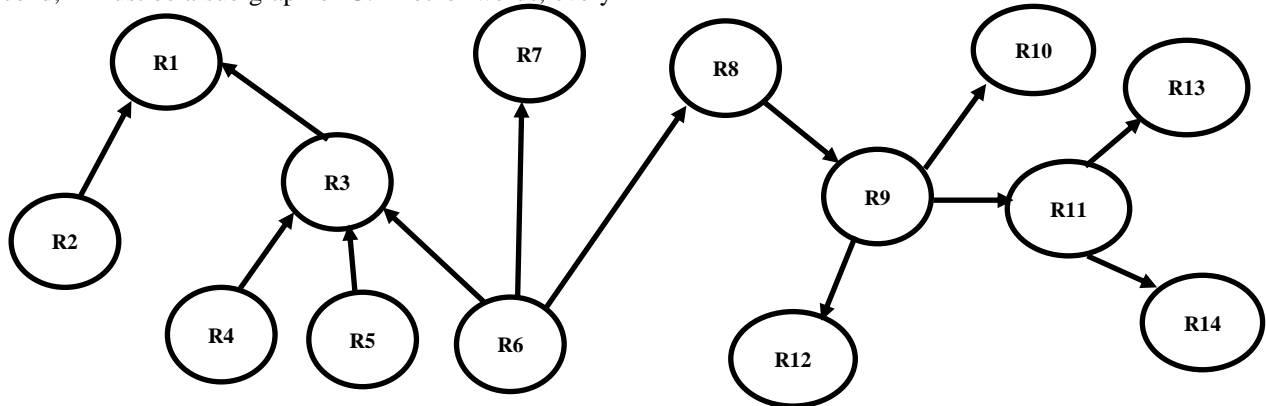


Fig. 1. Graph Connecting Requirements for Making Spanning Tree from Graphs.

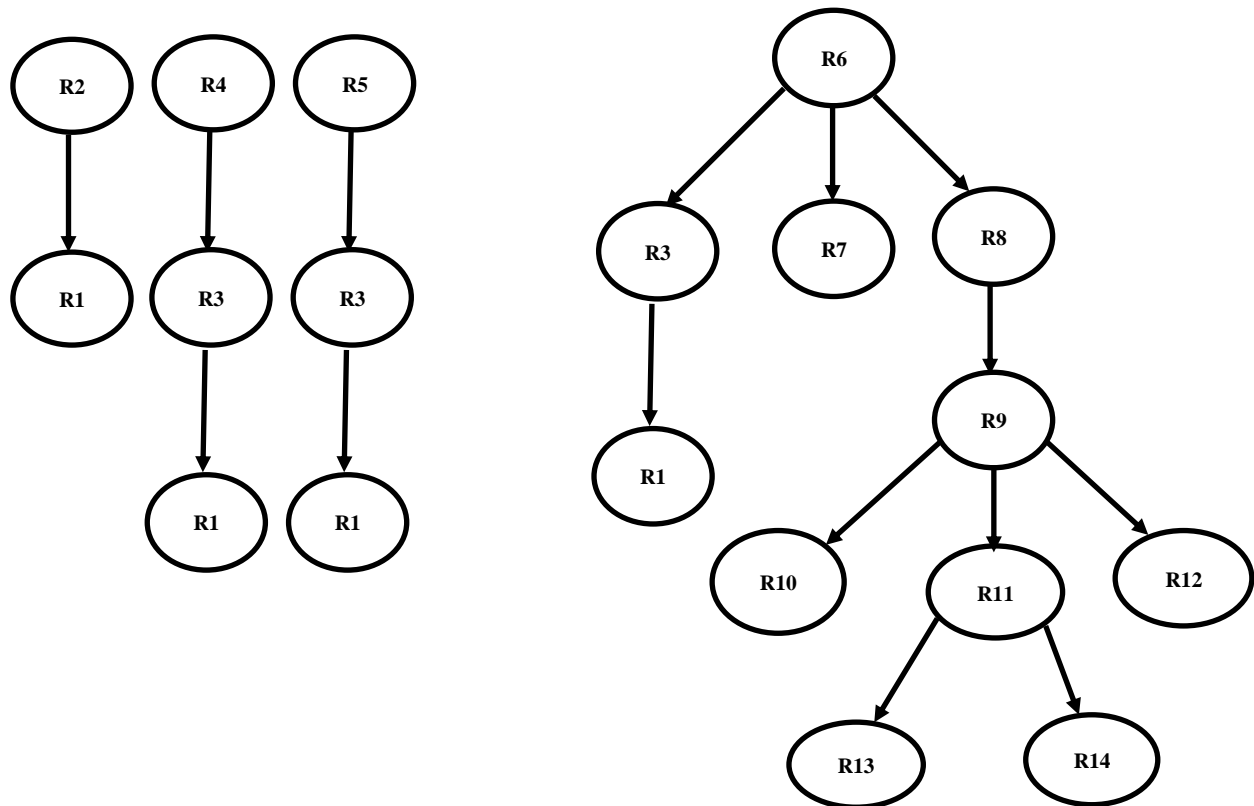


Fig. 2. Tree 1, Tree 2, Tree 3, Tree 4, Respectively.

B. Analytical Hierarchical Process (AHP)

Spanning trees will show the relationship of requirement with other requirements. As shown in Fig. 2, a finite number of spanning trees will be produced from directed graph. Next, AHP will be applied to individual trees or combination of many trees that have common requirements. The main idea is that while applying AHP to spanning tree, only depended requirements will be compared, hence resulting in optimal prioritization in a reduced time. For example, consider the spanning tree with starting node R6 shown in Fig. 2, R6 will be compared with R3, R1 and R7 as it is required for all these requirements. However, R7 will be not compared with R1 or R3 as there is no direct relation with these requirements. In this case, when R6 is compared with R3 or any other requirement, then there is no need to compare between R3 with R6. Requirements that are not depended can be considered as equal during comparison and assigned with value 1. This means with the help of spanning tree, the number of comparisons can be greatly reduced. AHP can be applied to either every spanning tree individually or combination of two or more trees if they have some requirements in common. We have five spanning trees as given in Fig. 2. AHP will be applied to first four spanning trees combined as they are related by some common requirements. First, apply AHP to Tree 5 starting with root R8 and then apply AHP to combined four trees. Table I shows requirements of Tree 5 for comparison and calculation.

From Table I, we can see that we can put value either 1 or greater than 1 while comparing any two requirements. We can

only put 1 or greater value where 1 represents equal priority requirements and value greater than 1 represents those requirements that have not equal priorities. For instance, we can use values such as 2, 3, 4, ..., n for requirements that are not equal in priorities. If we increase the value, the difference in both requirements will be increased. The value 2 is taken for requirement that is needed for other requirement. For instance, if R1 is required for R2 and R2 is required for R3, then we will put 2 for R1 against R2 and will put 4 for R1 against R3. The value 1 is taken for requirements that have either equal priority or not related and 2 against those requirements that need this particular requirement as well and value 1/2 for the reverse case. In this case, as R8 is required for R9, therefore the value is 2 against R9 for R8. Priority value for each requirement against other requirements is shown in Table I e.g. priority of R9 against R8 is 0.5 which means priority of R8 is double as compare to R9. For independent requirements like R10 and R12, we put value 1 because these requirements have no relation.

Next, the task is to calculate normalized values for each requirement by dividing the values of each column value in Table I by column sum. Column sum for each column is shown in Table II. For example, the value 1 in the first row and the first column will be divided by 2.5, which comes to 0.4. Consequently, normalized values for each requirement are shown in Table II. The column sum2 represents the averaging over normalized values for each row. The same process is then repeated for the combined four trees together and the values obtained in shown in Table III.

TABLE I. PAIRWISE COMPARISON FOR TREE 5

|     | R8    | R9    | R10    | R11    | R12    | R13    | R14    |
|-----|-------|-------|--------|--------|--------|--------|--------|
| R8  | 1.000 | 2.000 | 4.000  | 4.000  | 4.000  | 8.000  | 8.000  |
| R9  | 0.500 | 1.000 | 2.000  | 2.000  | 2.000  | 4.000  | 4.000  |
| R10 | 0.250 | 0.500 | 1.000  | 1.000  | 1.000  | 1.000  | 1.000  |
| R11 | 0.250 | 0.500 | 1.000  | 1.000  | 1.000  | 2.000  | 2.000  |
| R12 | 0.250 | 0.500 | 1.000  | 1.000  | 1.000  | 1.000  | 1.000  |
| R13 | 0.125 | 0.250 | 1.000  | 0.500  | 1.000  | 1.000  | 1.000  |
| R14 | 0.125 | 0.250 | 1.000  | 0.500  | 1.000  | 1.000  | 1.000  |
| Sum | 2.500 | 5.000 | 11.000 | 10.000 | 11.000 | 18.000 | 18.000 |

TABLE II. NORMALIZATION AND AVERAGING AND FOR TREE 5

|     | R8    | R9    | R10   | R11   | R12   | R13   | R14   | Sum2/<br>priority | Out of 1<br>(x = sum/7) | Z= (x/2) |
|-----|-------|-------|-------|-------|-------|-------|-------|-------------------|-------------------------|----------|
| R8  | 0.400 | 0.400 | 0.360 | 0.400 | 0.360 | 0.440 | 0.440 | 2.800             | 0.400                   | 0.200    |
| R9  | 0.200 | 0.200 | 0.180 | 0.200 | 0.180 | 0.220 | 0.220 | 1.400             | 0.200                   | 0.100    |
| R10 | 0.100 | 0.100 | 0.090 | 0.100 | 0.090 | 0.055 | 0.055 | 0.600             | 0.090                   | 0.045    |
| R11 | 0.100 | 0.100 | 0.090 | 0.100 | 0.090 | 0.110 | 0.110 | 0.700             | 0.100                   | 0.050    |
| R12 | 0.100 | 0.100 | 0.090 | 0.100 | 0.090 | 0.055 | 0.055 | 0.600             | 0.090                   | 0.045    |
| R13 | 0.050 | 0.050 | 0.090 | 0.050 | 0.090 | 0.055 | 0.055 | 0.440             | 0.060                   | 0.030    |
| R14 | 0.050 | 0.050 | 0.090 | 0.050 | 0.090 | 0.055 | 0.055 | 0.440             | 0.060                   | 0.030    |



TABLE III. CALCULATING PRIORITIES OF TREE 1 TO TREE 4 (COMBINED)

|    | R1    | R2    | R3    | R4    | R5    | R6    | R7    | Sum2/<br>priority | Out of 1<br>(y = sum/7) | Z = (y/2) |
|----|-------|-------|-------|-------|-------|-------|-------|-------------------|-------------------------|-----------|
| R1 | 0.055 | 0.076 | 0.050 | 0.040 | 0.040 | 0.043 | 0.125 | 0.430             | 0.060                   | 0.030     |
| R2 | 0.110 | 0.153 | 0.105 | 0.170 | 0.170 | 0.173 | 0.125 | 1.000             | 0.140                   | 0.070     |
| R3 | 0.110 | 0.153 | 0.105 | 0.086 | 0.086 | 0.086 | 0.125 | 0.751             | 0.110                   | 0.055     |
| R4 | 0.220 | 0.153 | 0.210 | 0.170 | 0.170 | 0.173 | 0.125 | 1.221             | 0.200                   | 0.100     |
| R5 | 0.220 | 0.153 | 0.210 | 0.170 | 0.170 | 0.173 | 0.125 | 1.221             | 0.200                   | 0.100     |
| R6 | 0.220 | 0.153 | 0.210 | 0.170 | 0.170 | 0.173 | 0.250 | 1.346             | 0.200                   | 0.100     |
| R7 | 0.055 | 0.153 | 0.105 | 0.170 | 0.170 | 0.086 | 0.125 | 0.864             | 0.123                   | 0.062     |

The column sum2 also shows the priority value of every requirement of the spanning tree, in particular, or combination of spanning trees. The sum of these sum2 values will equal to number of requirements i.e. 7. These values can be divided on number of requirements to find priority of requirements out of 1. For considering whole set of requirements i.e. In 14 requirements, priority value will be divided on 2 (2 is sum value for all requirements priorities). Column value z for Table II and Table III shows priority out of 14 requirements. Priority out of 14 is calculated. Similarly, for calculating priority of requirement in 100, value 100 is multiplied.

### C. Time Complexity of SAHP

Time complexity of AHP depends on total number of pairwise comparisons. With spanning tree, total number of comparisons are reduced because of limited number of relations. Either we consider combination of all spanning trees in one table or individual trees, the number of comparisons of dependent or related requirements will be always same (from adjacency matrix one can see how much relations exists). The number of comparisons in all cases will depend on how much relations of requirements in graph exist. In this example, as only 20 relations are possible, the total number of comparisons will equal to only 20. Therefore, in this way, number of necessary comparisons are reduced from  $n*(n-1) / 2$ , which was from 91 to only 20 in this example. This reduction in value shows the advantage of using spanning trees for related depended requirements only. Overall values and calculations during comparing requirements can be reduced by considering individual trees for prioritization as explained.

For given requirements set, maximum relations that can exists are equal to  $((n-1) + (n-2) + (n-3) + \dots + (n-n))$ , where n are total number of requirements. This is possible when all requirements are connected point to point in chain like structure such that one requirement is dependent on other requirement. The value of n will be decremented and will be added until it reaches to 0. In such case, total number of comparisons will become  $n*(n-1) / 2$  which is equal to number of comparisons of AHP. The minimum number of relations will be 0 in any requirements set. In such case, priority of all requirements will be consider to be equal i.e. 1. Fig. 3 shows number of comparisons of two techniques i.e. AHP without spanning trees by considering all requirements and AHP with spanning trees. Let's take 10 requirements. Minimum possible relations are 0 while maximum relations can be 45. Any number of relations can be possible between 0 and 45. The orange linear line of Fig. 3 shows that number of comparisons

in this proposed approach is directly proportional to number of relations. It is equal to 45 i.e. case of AHP where maximum relations exist. In small set of requirements where requirements are few in amount, this is possible that maximum relations exist (number of relations reaches number of requirements) such that each requirement is point to point connected with other requirement but we rarely can see such number of relations in large set of requirements like ERP.

From this discussion, it can be concluded that by comparing only the depended requirements through spanning tree, the number of comparisons and calculations can be greatly reduced. Therefore, although total comparisons of dependent requirements are same in all cases, but as the entire project, the number of comparisons and calculated normalized values are not same due to independent requirements.

### D. Requirements Priority

Priority is assigned to requirements on the basis of its position in spanning tree i.e. how much they are needed and dependent on other requirements. Requirements need can either increase breadth-wise or depth-wise. In either case, priority can increase but priority values in both cases can be different. Similarly, priority of requirement can decrease when requirements are dependent and wait for other requirements.

AHP can be applied for calculating priority of requirement on the basis of how much they are depended or required for other requirements. AHP is simple and accurate prioritizing technique that can find priority of requirements by comparing pairwise all requirements together. If requirement let say R1 is required for R2 and R2 is required for R3, then priority of R1 can be taken as double of R2 or it can be said that priority of R1 is two times as compared to R2 while R1 priority is 4 times as compare to R3. The following scenarios show different cases of requirements behavior as they change when applied with AHP.

**Scenario 1:** In this scenario priority of requirement is determined when its need for other requirements increases breadthwise. Breadthwise contain all requirements on same level with same priority. Two cases can be considered here, one with seven requirements and other with five requirements and calculate priorities.

**Case 1:** In this case, R1 is required for six other requirements with all requirements on same level with same priority as shown in Fig. 4.

Through AHP, we have calculated priority of R1 by comparing all seven requirements together which is equal to 1.75. R1 is considered to be double in priority as compare to individual requirements during pairwise comparison. The priority of all other requirements is shown in Table IV. Table IV summarizes priority values for all requirements.

**Case 2:** In this case, R1 is required for four other requirements with all requirements on same level or priority in Fig. 4. Now priority of R1 is reduced to 1.32 as shown in Table IV.

**Scenario 2:** In this scenario, requirements size increases depth wise. In case 01 of scenario 1, R1 is required for six

other requirements in depth wise structure such that one requirement is depended on other requirement as shown in Fig. 5. In case 02, number of requirements that need R1 are reduced from six to four. Priority of R1 in first case comes out 3.5 while in second case it is 2.21. The priority of R1 in 2<sup>nd</sup> case of scenario 2 (required for four requirements) is still greater than case 01 of scenario 1 (required for six requirements). This shows priority increases with greater ratio depth wise and this has advantage because in scenario 1, R1 is available to all requirements after implementation but in scenario 2, it is not available to all requirements e.g. R7 in scenario 2 can't be implemented when R6 is not available but in scenario 1, all requirements are dependent on R1.

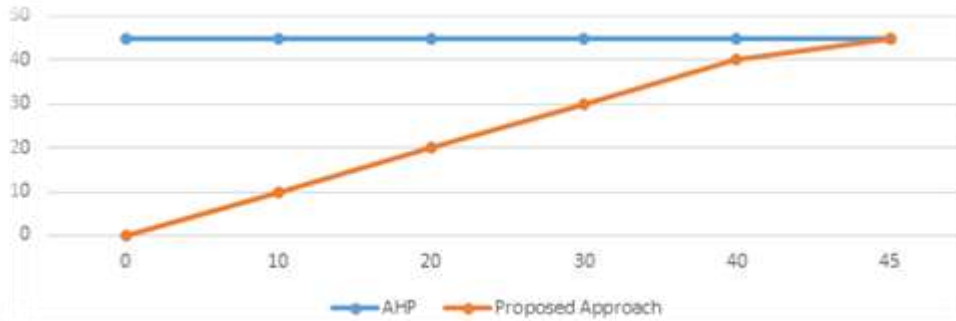


Fig. 3. Comparison of AHP and Proposed Approach.

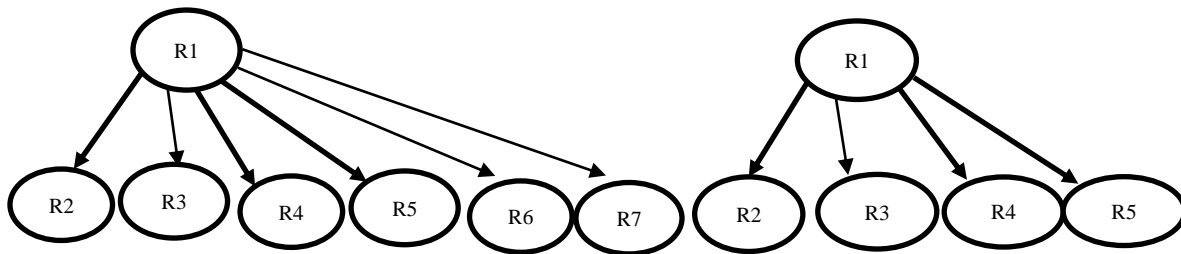


Fig. 4. Breadth-Wise Increase of Requirements.

TABLE IV. COMPARISON OF PRIORITY OF REQUIREMENTS AS RESULT OF AHP

| Requirements | Scenario 1 |        | Scenario 2 |        | Scenario 3 |
|--------------|------------|--------|------------|--------|------------|
|              | Case 1     | Case 2 | Case 1     | Case 2 |            |
| R1           | 1.750      | 1.32   | 3.500      | 2.210  | 1.0        |
| R2           | 0.875      | 0.68   | 1.750      | 1.245  | 0.5        |
| R3           | 0.875      | 0.68   | 0.875      | 0.758  | 0.5        |
| R4           | 0.875      | 0.68   | 0.437      | 0.515  | 0.5        |
| R5           | 0.875      | 0.08   | 0.210      | 0.400  | 0.5        |
| R6           | 0.875      | x      | 0.105      | 0.900  | 2.0        |
| R7           | 0.875      | x      | 0.500      | 0.900  | 2.0        |

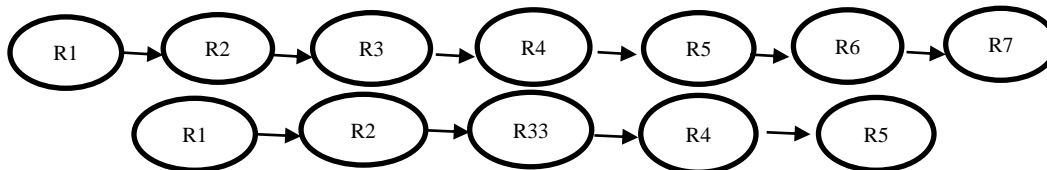


Fig. 5. Depth-Wise Increase of Requirements.

**Scenario 3:** Priority of requirement decreases when its dependency on other requirements increase. The reason is that during comparison against other requirements, sum of values are reciprocal of 1. Fig. 6 shows the priority of R1 against R6 and R7 will be equal to  $\frac{1}{2}$ . The sum of reciprocal values will reduce the priority of requirement. Priority of R1 is now 1, which is minimum as compared to all cases. Priority of other requirements are shown in Table IV.

From values given in above Table IV, it can be concluded that requirement priority is associated with its increasing size but the ratio in which it increases depth wise is greater than breadth wise and it should be increase with high ratio in depth wise as compare to breadthwise because in breadthwise, the

pre-requisite requirement is available for all requirements and the delay is not too much as compare to the case of depth wise where pre-requisite requirement is not available for all requirements and by delaying this requirement can delay the implementation of its requirements more in case of parallel developing project.

Similarly, if number of pre-requisite requirements and number of requirements for which particular requirement is needed are equal then priority of requirement will be equal. For example, in Fig. 7, the number of backward and forward requirements for R1 are equal, in all cases priority of requirement will be equal. With AHP, we have calculated priority of R1 that is 0.84 for all cases of Fig. 7.

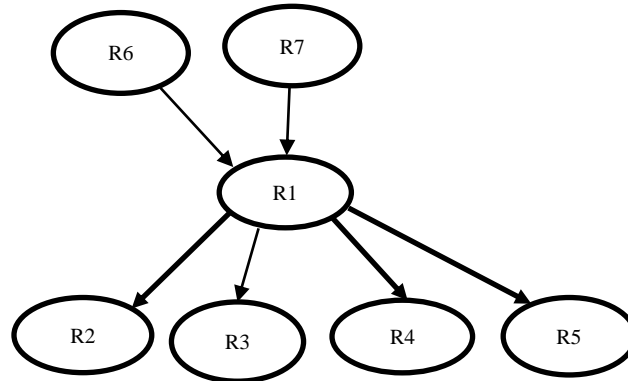


Fig. 6. Number of Pre-Requisite Connected with Requirement.

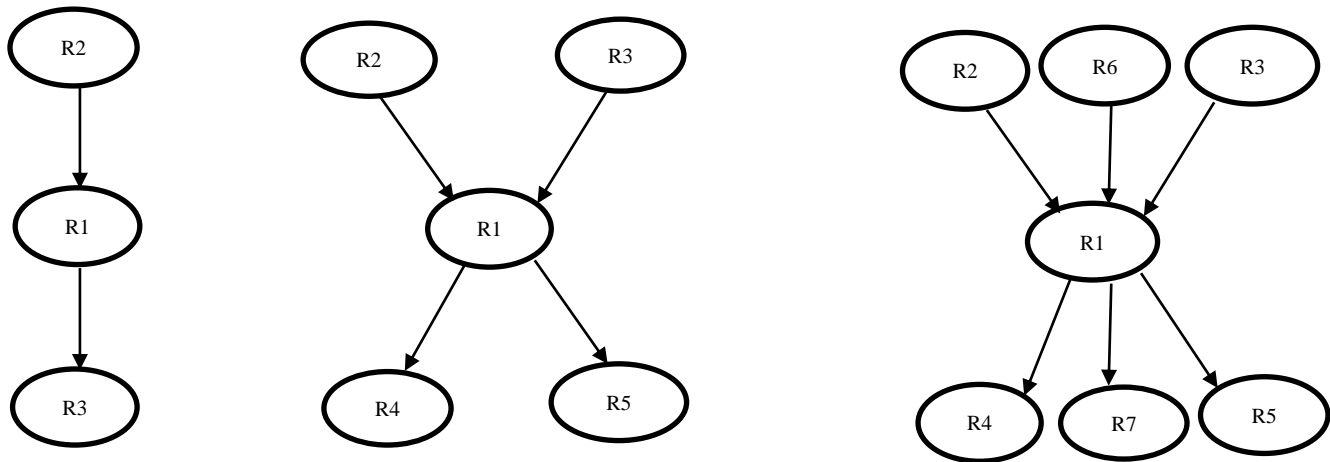


Fig. 7. Distribution of Requirements with Same Ratio.

#### IV. VALIDATION OF SAHP ON ODOO ERP

SAHP was evaluated on requirements of On Demand Open Object (ODOO). ODOO is open source ERP software system that is used by millions of users for managing hundreds of possible enterprises and their resources. In many of research studies, authors used different modules of ODOO ERP [34][35][36]. In ERP system, all modules are integrated which shows that all the requirements should be inter-related. Modules of ERP are highest level URs that are further comprised of low level FRs. With spanning tree, we can relate FRs that can belong to any module. Module is just high level abstraction to which requirements of same nature belong e.g.

customer and supplier creation are FRs that belong to HR module while customer sale and supplier sale are FRs that belong to sale management module. With spanning tree, we can relate these FRs that belong to different modules. Thus spanning tree does not show abstraction or high level representation of requirements because it relates only different requirements that belong to particular module. Selection of particular modules have impact on priority of their FRs. This means variations in selecting different modules by users have impact on FRs structure. The suggested prioritization approach will be applied on the FRs of ODOO to prioritize them. The modules of ERP consists of 96 FRs for this study as shown in Table V.

TABLE V. REQUIREMENTS OF ODOO ERP FOR HR MODULE

| Notation | Requirement                      | Module No. | Required For   | Tree    | Notation | Requirement               | Module No. | Required For                       | Tree       |
|----------|----------------------------------|------------|--|---------|----------|---------------------------|------------|------------------------------------|------------|
| R1       | employee creation                | 1          | R81, R25, R23, R67, R2, R4, R10, R11, R12, R17, R18, R20, R21, R22, R7, R9, R8 | T1      | R69      | sale return view          | 3          |                                    | T10        |
| R2       | public information's of employee | 1          |  | T1      | R42      | purchase                  | 4          | R51, R59                           | T4, T5, T6 |
| R3       | employee personal info           | 1          |  |         | R59      | purchase view             | 4          |                                    | T4, T5, T6 |
| R4       | contact info                     | 1          |  | T4      | R60      | purchase return           | 4          | R68                                | T4         |
| R5       | job position                     | 1          |  | T2, T3  | R68      | purchase return view      | 4          |                                    | T4         |
| R6       | department                       | 1          | R5, R61, R67   | T2, T3  | R34      | product                   | 5          | R42, R60, R66, R35, R70, R71, R90, |            |
| R7       | job information's                | 1          |  |         | R66      | stock ledgers             | 5          |                                    | T4         |
| R8       | manager                          | 1          | R5, R24, R67   |         | R70      | product transfer in       | 5          |                                    | T4         |
| R9       | coach                            | 1          |  |         | R71      | product transfer out      | 5          |                                    | T4         |
| R10      | contract information's           | 1          |  | T1      | R56      | company                   | 5          |                                    |            |
| R11      | contract reference information's | 1          |  | T1      | R90      | manufacturing orders      | 5          |                                    | T4         |
| R12      | salary generation                | 1          | R21  | T1, T18 | R24      | project management        | 6          | R26, R27, R28, R29                 | T3         |
| R22      | hr expenses                      | 1          | R23  | T1      | R25      | add team members          | 6          |                                    | T1         |
| R23      | hr expenses detail               | 1          |  | T1      | R26      | extra information's       | 6          |                                    | T3         |
| R33      | customer detail                  | 1          | R73, R55, R36, R35, R61, R64, R39  | T10     | R27      | project stages            | 6          |                                    | T3         |
| R37      | sales persons                    | 1          | R58, R63, R35  |         | R28      | view current task         | 6          |                                    | T3         |
| R41      | supplier detail                  | 1          | R44, R65, R72, R42, R52, R60   |         | R29      | create a task             | 6          | R31                                | T3         |
| R43      | sales man                        | 1          | R42, R44   | T5      | R30      | extra information's       | 6          |                                    |            |
| R57      | region                           | 1          | R58  | R8      | R31      | tasks stages              | 6          |                                    | T3         |
| R58      | area                             | 1          | R35  | R7, R8  | R93      | directories for documents | 7          |                                    | T11        |
| R80      | job position in recruitment      | 1          |  | T4      | R94      | documents history         | 7          | R96                                | T15        |
| R81      | job                              | 1          |  | T1, T2  | R95      | documents attachments     | 7          | R96                                | T14        |
| R82      | appraisal form                   | 1          |  |         | R91      | fleet management          | 8          | R92                                | T11        |
| R83      | create a job position            | 1          |  |         | R92      | vehicle repairing         | 8          |                                    |            |
| R84      | recruitment form                 | 1          |  |         | R13      | salary rules              | 9          |                                    | T18        |

|     |                       |   |               |            |     |                      |    |               |                   |
|-----|-----------------------|---|---------------|------------|-----|----------------------|----|---------------|-------------------|
| R85 | job selection process | 1 |               |            | R14 | salary structure     | 9  | R12           | T16               |
| R86 | link tracker          | 1 |               |            | R15 | salary categories    | 9  | R12           | T17               |
| R87 | mass mailing          | 1 |               |            | R16 | registers            | 9  | R12, R13      | T18               |
| R88 | contact               | 1 |               |            | R21 | hr payroll process   | 9  |               | T1, T16, T17, T18 |
| R89 | business pipeline     | 1 |               |            | R17 | apply for leave      | 10 | R19, R20      | T1                |
| R38 | customer receipts     | 2 |               | T10        | R18 | allocation request   | 10 |               | T1                |
| R39 | customer payment      | 2 | R55, R38      | T10        | R19 | leave approval       | 10 |               | T1                |
| R40 | supplier receipts     | 2 |               | T12        | R20 | leave summary        | 10 |               | T1                |
| R44 | supplier refund       | 2 |               | T5, T6     | R46 | bank statement       | 11 | R47           | T9                |
| R45 | supplier payment      | 2 | R40           | T12        | R47 | bank detail          | 11 | R49, R50, R53 | T9                |
| R52 | supplier payment      | 2 |               |            | R48 | cash registers       | 11 |               |                   |
| R53 | journals accounts     | 2 | R54           | T9         | R49 | put money in         | 11 |               | T9                |
| R54 | chart of accounts     | 2 |               | T9, T10    | R50 | put money out        | 11 |               | T9                |
| R55 | analytic accounts     | 2 | R54           | T10        | R51 | profit and lost      | 11 |               | T4, T5, T6        |
| R63 | salesman ledgers      | 2 |               | T7         | R75 | compose message      | 12 |               |                   |
| R64 | customer ledgers      | 2 |               | T10        | R76 | message inbox        | 12 | R79           | T13               |
| R65 | supplier ledgers      | 2 |               | T6         | R77 | message draft        | 12 |               |                   |
| R67 | hr expense management | 2 |               | T1, T2, T3 | R78 | sent messages        | 12 |               |                   |
| R74 | balance sheet         | 2 |               |            | R79 | message searching    | 12 |               | T13               |
| R32 | customer invoice      | 3 | R36           |            | R72 | order to suppliers   | 13 |               | T6                |
| R35 | sale                  | 3 | R61, R62, R32 |            | R73 | order from customer  |    |               | T10               |
| R36 | customer refund       | 3 |               | T10        | R96 | documents attachment |    |               | T14, T15          |
| R61 | sale return           | 3 | R69           | T10        |     |                      |    |               |                   |
| R62 | sale view             | 3 |               |            |     |                      |    |               |                   |

### A. Results and Discussion

Results of prioritization of ODOO ERP requirements after applying suggested framework using AHP and spanning tree combination have been calculated. Requirements are prioritized by applying the same criteria discussed.

1) *Spanning trees*: As the result, 8 spanning trees are constructed (T1, T2 up to T18) while 19 requirements are independent requirements which are neither required nor dependent on other requirements. The root and the detail requirements are given in Table VI. Spanning trees are categorized into different groups which are made on the basis

of common requirements in different spanning trees. For example, in T1 and T2, the common requirement is R67. Similarly, R21 is common in T1, T16, T17 and T18. Six groups (A, B, C, D, E, and F) of different trees are made which are shown in Table VI.

2) *Applying AHP to spanning trees*: The column “priority” as shown in Table VII shows priority of requirements as a result of applying AHP on spanning tree. Priority of requirements in spanning trees are calculated. We have calculated priority of these requirements out of 100 as shown in Table VII.

TABLE VI. COMBINING REQUIREMENTS OF SPANNING TREES

| Group                      | Tree | Root                              | Requirements   | Efforts (Hours) |
|----------------------------|------|-----------------------------------|--|-----------------|
| A                          | T1   | R1                                | R81, R23, R25, R2, R4, R10, R11, R12, R17, R18, R19, R20, R22, R21, R67        | 720             |
|                            | T2   | R6                                | R5, R67, R81,  |                 |
|                            | T3   | R8                                | R5, R67, R24, R26, R27, R28, R29, R31  |                 |
|                            | T16  | R14                               | R21  |                 |
|                            | T17  | R15                               | R21  |                 |
|                            | T18  | R16                               | R12, R13, R21  |                 |
| B                          | T9   | R46                               | R47, R49, R50, R53, R54  | 1230            |
|                            | T8   | R57                               | R58  |                 |
|                            | T7   | R37                               | R58, R63, R35, R61, R62, R32, R36, R69   |                 |
|                            | T10  | R33                               | R73, R55, R54, R35, R61, R62, R32, R36, R69, R64, R38, R39,                    |                 |
|                            | T4   | R34                               | R42, R51, R59, R60, R66, R68, R70, R71, R80, R90, R35, R61, R62, R32, R36, R69 |                 |
|                            | T5   | R43                               | R42, R51, R59, R44   |                 |
| T6                         | R41  | R42, R51, R59, R44, R52, R60, R68 |  |                 |
| C                          | T11  | R92                               | R93  | 50              |
| D                          | T12  | R45                               | R40  | 50              |
| E                          | T13  | R76                               | R79  | 50              |
| F                          | T14  | R95                               | R96  | 80              |
|                            | T15  | R94                               | R96  |                 |
| Individual requirements    |      |                                   |  | 470             |
| Total efforts in man hours |      |                                   |  | 2650            |

TABLE VII. REQUIREMENTS PRIORITY OF ODOO

| Notation | Combined Priority<br>(Out of 100) | Separate Priority<br>(Out of 100) | Notation | Combined Priority<br>(Out of 100) | Separate Priority<br>(Out of 100) |
|----------|-----------------------------------|-----------------------------------|----------|-----------------------------------|-----------------------------------|
| R1       | 1.66                              | 2.22                              | R62      | 0.72                              | 0.84                              |
| R2       | 0.96                              | 0.96                              | R69      | 0.62                              | 0.79                              |
| R3       | 1.03                              | 1.03                              | R42      | 0.9                               | 0.91                              |
| R4       | 0.96                              | 0.96                              | R59      | 0.77                              | 0.75                              |
| R5       | 0.98                              | 0.82                              | R60      | 0.9                               | 0.90                              |
| R6       | 1.08                              | 1.20                              | R68      | 0.81                              | 0.80                              |
| R7       | 1.03                              | 1.03                              | R34      | 2.37                              | 2.90                              |
| R8       | 1.65                              | 2.68                              | R66      | 0.9                               | 0.90                              |
| R9       | 1.03                              | 1.03                              | R70      | <b>0.93</b>                       | 0.94                              |
| R10      | 0.96                              | 0.96                              | R71      | <b>0.9</b>                        | 0.90                              |
| R11      | 0.96                              | 0.96                              | R56      | 1.03                              | 1.03                              |
| R12      | 0.97                              | 0.88                              | R90      | 0.9                               | 0.90                              |
| R22      | 0.96                              | 0.96                              | R24      | 1.16                              | 1.34                              |
| R33      | 2.72                              | 3.10                              | R25      | 0.96                              | 0.96                              |
| R37      | 2.57                              | 2.70                              | R26      | 0.92                              | 0.73                              |
| R41      | 1.47                              | 2.056                             | R27      | 0.92                              | 0.73                              |
| R43      | 1.20                              | 1.414                             | R28      | 0.92                              | 0.73                              |
| R57      | 1.031                             | 1.045                             | R29      | 0.96                              | 0.78                              |
| R58      | 0.78                              | 0.79                              | R30      | 1.03                              | 1.03                              |
| R80      | 0.86                              | 0.79                              | R31      | 0.92                              | 0.61                              |
| R81      | 0.98                              | 0.92                              | R93      | 0.72                              | 0.72                              |
| R82      | 1.03                              | 1.03                              | R94      | 1.23                              | 1.23                              |
| R83      | 1.03                              | 1.03                              | R95      | 1.23                              | 1.23                              |

|     |      |      |     |       |      |
|-----|------|------|-----|-------|------|
| R84 | 1.03 | 1.03 | R91 | 1.03  | 1.03 |
| R85 | 1.03 | 1.03 | R92 | 1.37  | 1.37 |
| R86 | 1.03 | 1.03 | R13 | 0.97  | 0.90 |
| R87 | 1.03 | 1.03 | R14 | 1.01  | 1.07 |
| R88 | 1.03 | 1.03 | R15 | 1.01  | 1.07 |
| R89 | 1.03 | 1.03 | R16 | 1.1   | 1.44 |
| R38 | 0.87 | 0.81 | R21 | 0.90  | 0.88 |
| R39 | 1.10 | 1.12 | R17 | 1.02  | 1.06 |
| R40 | 0.72 | 1.03 | R18 | 0.96  | 0.96 |
| R44 | 0.87 | 0.85 | R19 | 0.95  | 0.88 |
| R45 | 1.37 | 1.37 | R20 | 0.96  | 0.88 |
| R52 | 0.90 | 0.90 | R46 | 1.6   | 2.63 |
| R53 | 1.08 | 0.79 | R47 | 1.15  | 1.17 |
| R54 | 0.70 | 0.41 | R48 | 1.03  | 1.03 |
| R55 | 0.90 | 0.81 | R49 | 0.89  | 0.60 |
| R63 | 0.84 | 0.85 | R50 | 0.89  | 0.60 |
| R64 | 0.93 | 0.87 | R51 | 0.77  | 0.75 |
| R65 | 0.91 | 0.90 | R75 | 1.03  | 1.03 |
| R67 | 0.92 | 0.83 | R76 | 1.37  | 1.37 |
| R74 | 1.03 | 1.03 | R77 | 1.03  | 1.03 |
| R32 | 0.77 | 0.88 | R78 | 1.03  | 1.03 |
| R35 | 1.34 | 1.47 | R79 | 0.72  | 0.72 |
| R36 | 0.62 | 0.79 | R72 | 0.91  | 0.90 |
| R61 | 0.77 | 0.64 | R73 | 0.93  | 0.88 |
| R23 | 0.96 | 0.96 | R96 | 0.618 | 0.62 |

### B. Time Estimation

Time estimation is time taken by particular requirement to complete its implementation. Every requirement consume certain amount of efforts on the basis of which time can be calculated. Many models are suggested by authors for calculating efforts and time estimation of requirements and projects. We applied USE CASE point (UCP) estimation technique which was simple in use and more appropriate for our requirements. The UCP estimation method was presented initially in 1993 by Karner estimates efforts in person-hours based on use cases that primarily specify FRs of a system [11][12]. Use cases are assumed to be developed from scratch, be sufficiently detailed and typically have less than 10-12 transactions. The method has previous been used in numerous industrial software development projects. There have been promising outcomes and the method was highly accurate than expert estimates in industrial trials.

UCP defines the functional scope of the system to be developed. Attributes of a use case model may therefore serve as measures of the size and complexity of the functionality of a system. After following all steps of USE case estimation technique, effort in hours for each requirement is calculated. After approximation, we have divided requirements into three categories as follows.

- First category contains requirements that take approximately 20 hours to complete its implementation. This is time just needed to implement requirement with functionalities. This time contain unit and integration testing time.
- Second category contain requirements that contain requirements that take approximately 30 hours to complete its implementation.
- Third category contain requirements contain requirements that take approximately 60 hours to complete its implementation.

Completion time of particular module will be sum of time taken by all its requirements. This time reduces when the project is to be developed by parallel team members. But total actual time can exceed calculated time in parallel development projects because requirements are interrelated to each other's and waiting time for particular requirements can cause delay in projects. The purpose of prioritization is to minimize the delay or waiting time.

### V. DISTRIBUTION OF REQUIREMENTS IN PARALLEL DEVELOPERS

From results of prioritization we can conclude that not only priority value and order of requirements is necessary for reducing delays and assuring timely delivery of project but

distribution of requirements in team members is also necessary. Total delivery time of project is equal to maximum time taken by any team member to implement all requirements. Distribution of requirements as shown in Table VIII are not uniform e.g. actual time estimation of requirements of A = 410 hours, B = 610 hours, C = 880 hours and D = 670 hours. Total delivery time of the project can exceed from 880 hours due to waiting time which is the maximum time of team member C but total time can't be less than 880 hours. This is because C is given those requirements which take more time in hours. Efficient distribution will be in that case where everyone is given requirements with same efforts. The generalized formula we can make for equal distribution is as follows.

$$\begin{aligned} \text{Total efforts (for any team member)} \\ = \text{Total efforts (man hours)} / 4. \end{aligned}$$

Where total efforts (man hours) = Total efforts (for all requirements starts from R1 to Rn).

From this formula, we will get average time for every team member which becomes 660 hours. If every team member gets no more than 660 hours than in ideal case total estimation time of delivery of project can be 660 hours which is reduced from 880 hours. This means further adjustment will be needed to reduce time estimation more and for this purpose some requirements of C can be assigned to A.

Along with equal distribution of requirements, we should reduce dependency among requirements of different team members as much as possible. Requirements of A that are required for C can be adjusted and can be assigned to C. Similarly, some requirements of C can be adjusted and implemented by A. From the spanning tree, one can easily

identify which requirements are dependent on each other, so dependent requirements can be assigned to same team members. In ideal case, distribution of requirements will be uniform and dependency between different team member requirements will be zero.

The best way to distribute requirements is thus assigning requirements of whole spanning tree to same team member. Requirements of spanning trees should be adjusted in such a way that every team member get requirements with equal weight of man hours. Team members can either implement big spanning tree requirements or requirements of many small spanning trees. If some trees requirements are distributed in more than one member than requirements should be prioritized in order to reduce the waiting time.

#### A. Combining and Splitting the Spanning Trees

If two or more than two trees have some common requirements than we can combine two trees and consider as one group. The reason is that common requirements are depended on requirements of more than one trees requirements and hence this dependency can increase waiting time and cause delays in parallel developing projects. Splitting process is taken when tree size is either big or difficult to assign all its requirements to single developer or sometimes small size trees are split to assure equal distribution of requirements. Table VIII shows how different trees are combined. Six groups were made as result of combining trees with common requirements. Total efforts in man hours for each group are also shown below. It is better to split tree at edge where two trees are combined for assigning requirements to different developers. For example, T9 and T10 are combined with R54, so the tree can be break here.

TABLE VIII. COMBINING AND SPLITTING OF REQUIREMENTS OF SPANNING TREES

| No of team members | Efforts per team (hours) with equal distribution | Splitting of trees  | Combining trees   | Time estimation for implementing requirements | Time completion with prioritization |
|--------------------|--|---|---|---|-------------------------------------|
| 01                 | 2650   | NIL   | All trees are considered  | 2650 hours                                    | 2650 hours                          |
| 02                 | 1325   | NIL   | Developer 1: [B + C + D]<br>Developer 2: [A + E + F + 470 individuals]  | 1230 hours                                    | 1230 hours                          |
| 03                 | 880  | Breaking of Group B:<br>1230 = 350 + 880  | Developer 1: [880 of B]<br>Developer 2: [350 of B + C + D + E + F + 300 individuals]<br>Developer 3: [A+ 170 individuals]   | 880 hours                                     | 880 hours                           |
| 04                 | 660  | Breaking of Group B: 1230 = 640 + 590<br>Breaking of Group A:<br>720 = 660 + 60 | Developer 1: [640 of B + 20 individuals]<br>Developer 2:[590 of B + C + 20 individuals]<br>Developer 3:[660 of A]<br>Developer 4: [60 of A + D + E + F + 430 individuals] | 670 hours                                     | 670 hours                           |



Common requirements can be assigned to any tree. Common requirements normally get low priority as they are dependent on other requirements. Similarly, T4 and T6 are combined with R42. In this case, we can break here (edge of R42) in order to equally distribute requirements. Common requirements can be adjusted with any tree requirements. But for equal distribution in terms of time efforts, especially in case where a single tree is quite large and needs to split, then we will split it. In such case, try to split an edge and assign those requirements that have significantly high priority difference from their parent requirements. It is better to split at edge where there exists quite big difference in priorities between two requirements. E.g. if tree T4 is to be split, there can many options, either to split edge at R60, R66, R35, R70, R71 or R90.

The difference in priorities between R34 and R35 is less as compare to other requirements because R35 is high priority requirement, so splitting at R35 can increase waiting time if R35 is assigned to different team member. Splitting at edge of low priority requirement and assigning it to other team member will decrease the effect of dependency and waiting time. For maintaining balance and equal distribution, more than one trees can be split e.g. T3 can be split along with T4 but at point where there exists quite difference in priority. Thus from values of SAHP, distributed priority can be determined requirements can be easily assigned to team members such that effect of dependency in requirements become low as much as possible.

**B. Distribution of Requirements**

Requirements will be distributed in such a way that there does not exist either relation between requirements of different team members or if relation exist, then requirements should be prioritized so that waiting time can be reduced and timely implementation of requirements can be assured. Few cases are considered for distribution of requirements as shown below.

1) *Distribution of requirements in 2 team members:* In distributing requirements based on efforts in man hours per team member, the value will be equal to 1325 hours i.e. half of

total 2650 hours. There is no need to split any tree or group of trees because different groups can be managed to produce total efforts of 1325 hours. We can assign requirements of groups B, C, D to one developer for implementation and groups A, E, F along with 470 individual requirements to second developer. In this way two different developers will get independent requirements with no relationship between any two requirements.

2) *Distribution of requirements on 3 team members:* In this case efforts per team member will be equal to 880 hours. While distributing requirements on three developers, it is must to split large tree or group of trees to assure equal distribution of requirements on developers. Requirements with total efforts of 350 hours were separated from group B. The separated requirements from any tree of group based on values of SAHP. Group B requirements after splitting will remain with efforts of 880 hours. In this way two sub groups are made. Sub-group with 350 hours can be adjusted with groups C, D, E, F and 300 hours of individual requirements to comprise total of 880 hours. Similarly, requirements of group A can be implemented along with remaining individual requirements i.e. 170 hours. In such way equal distribution of requirements can be assured. After distributing requirements, it is necessary to prioritize it to reduce waiting time and delays in project.

3) *Distribution of requirements in 4 team members:* To assure equal distribution of requirements, every team member will get requirements of 660 hours. For equal distribution, we can split group B into two subgroups with 640 and 590 hours. Similarly group A can be split into two subgroups with 660 and 60 hours. Splitting Group A were necessary as requirements of A were exceeded from 660 hours. Efficient distribution and prioritization of requirements reduces the effect of dependency between requirements and waiting time in parallel developing projects which results in timely delivery of projects. Separated requirements are shown in Table IX.

TABLE IX. SEPARATED REQUIREMENTS OF SPANNING TREES

| Number of Developers | A            | B  | C   | D   | E   | F   |
|----------------------|--------------|--|-----|-----|-----|-----|
| 2                    | nil          | nil  | nil | nil | nil | nil |
| 3                    | nil          | (R66, R70, R71, R80, R90, R42, R43, R44, R38, R52, R72)<br>OR<br>(R41, R44, R65, R72, R52, R42, R51, R59, R52, R60, R68) | nil | nil | nil | nil |
| 4                    | R4, R10, R11 | R66, R70, R71, R80, R90, R42, R43, R44, R38, R52, R72, R41, R64, R68, R65, R60   | nil | nil | nil | nil |

## VI. CONCLUSION

This paper proposed an approach for prioritizing FRs using AHP based on spanning trees. The proposed approach of SAHP has been presented in detail with evaluation on ODOO ERP system. The proposed framework is capable of prioritizing large-sized FRs while in active development cycle. As FRs are inter-related, so prioritization will help in easy arrangement of requirements. Similarly, apart from its implementation priority, in which the requirement is pre-requisite for other requirements, if we compare two requirements that are totally independent of each other, then deciding about which requirement is more important is a very important task. Importance of requirement was measured from how much it can reduce delay or waiting time.

Prioritizing and implementing important requirements decrease not only total estimation time but also decrease non-critical delay. Although non-critical delay does not increase estimation time of the project, it affects the waiting time of requirements. Another big problem that needs to be solved is that how much the proposed technique is scalable of handling and prioritizing large requirements size. Prioritizing large size requirements on the basis of its importance was solved using AHP and spanning tree in combination. AHP is used because it can solve dependency issues of requirements as it statistically compares pairwise for each and every requirement against other requirements. Requirements were represented with directed graph and spanning tree. From spanning tree, it became easy to decide about not only which requirement was necessary for other requirement but it became easy to compare all neighbor requirements that belong to same tree. AHP was applied to each tree separately and only depended requirements were scored value greater than 1. Priority of all other requirements during comparison were considered equal.

The results were obtained and were evaluated on parallel developing requirements of ODOO ERP. From different cases for prioritized and un-prioritized requirements, we showed that the proposed framework not only deal with big size requirements but reduce all possible delays in projects. We have shown that how spanning tree can help in equal and efficient distribution of requirements in parallel developing team members so that the effect of dependency and waiting time of requirements can be reduced. In future, we aim to do more industrial based experiments in order to validate framework on big projects and get feedbacks from industry.

## ACKNOWLEDGMENT

This paper is supported by Research Fund E15501, Research Management Centre, Universiti Tun Hussein Onn Malaysia.

## REFERENCES

- [1] M. Yaseen, S. Baseer, and S. Sherin, 'Critical Challenges for Requirement Implementation in Context of Global Software Development: A Systematic Literature Review', pp. 120–125, 2015.
- [2] M. Yaseen, M. Bacha, and Z. Ali, 'REVIEW PAPER COORDINATION AND COLLABORATION PRACTICES IN GLOBAL By', vol. 14, no. 2, 2020.
- [3] Z. Ali and M. Yaseen, 'Critical Challenges for Requirement Implementation in Global Software Development: A Systematic Literature Review Protocol with Preliminary Results', vol. 182, no. 48, pp. 17–23, 2019.
- [4] M. Yaseen, Z. Ali, and M. Humayoun, 'Requirements Management Model (RMM): A Proposed Model for Successful Delivery of Software Projects', *Int. J. Comput. Appl.*, vol. 178, no. 17, pp. 32–36, 2019.
- [5] A. U. Rahman, M. Yaseen, and Z. Ali, 'Identification of Practices for Proper Implementation of Requirements in Global Software Development: A Systematic Literature Review Protocol', vol. 177, no. 13, pp. 53–58, 2019.
- [6] Z. Ali, M. Yaseen, and S. Ahmed, 'Effective communication as critical success factor during requirement elicitation in global software development', vol. 8, no. 03, pp. 108–115, 2019.
- [7] M. Yaseen, S. . Baseer, S. . Ali, S. U. . Khan, and Abdullahh, 'Requirement implementation model (RIM) in the context of global software development', 2015 *Int. Conf. Inf. Commun. Technol. ICICT 2015*, 2015.
- [8] M. Yaseen and Z. Ali, 'Success Factors during Requirements Implementation in Global Software Development: A Systematic Literature Review', vol. 8, no. 3, pp. 56–68, 2019.
- [9] M. Yaseen, A. Mustapha, and N. Ibrahim, 'Minimizing Inter-Dependency Issues of Requirements in Parallel Developing Software Projects with AHP', vol. 8, no. Viii, 2019.
- [10] M. Yaseen, A. Mustapha, and N. Ibrahim, 'An Approach for Managing Large-Sized Software Requirements During Prioritization', 2018 *IEEE Conf. Open Syst.*, pp. 98–103, 2019.
- [11] N. Garg, M. Sadiq, and P. Agarwal, 'GOASREP: Goal Oriented Approach for Software Requirements Elicitation and Prioritization Using Analytic Hierarchy Process', pp. 281–287, 2017.
- [12] M. A. A. Elsood and H. A. Hefny, 'A Goal-Based Technique for Requirements Prioritization', 2014.
- [13] M. Yaseen, A. Mustapha, A. U. Rahman, S. Khan, and W. Kamal, 'Importance of Requirements Prioritization in Parallel Developing Software Projects', vol. 9, no. 2, pp. 171–179, 2020.
- [14] M. Yaseen, N. Ibrahim, and A. Mustapha, 'Requirements Prioritization and using Iteration Model for Successful Implementation of Requirements', *Int. J. Adv. Comput. Sci. Appl.*, vol. 10, no. 1, pp. 121–127, 2019.
- [15] R. Beg, R. P. Verma, and A. Joshi, 'Reduction in number of comparisons for requirement prioritization using B-Tree', no. March, pp. 6–7, 2009.
- [16] P. Tonella, A. Susi, and F. Palma, 'Interactive requirements prioritization using a genetic algorithm', *Inf. Softw. Technol.*, vol. 55, no. 1, pp. 173–187, 2013.
- [17] A. K. Massey, P. N. Otto, and A. I. Antón, 'Prioritizing Legal Requirements', vol. 1936, no. 111, 2010.
- [18] C. E. Otero, E. Dell, A. Qureshi, and L. D. Otero, 'A Quality-Based Requirement Prioritization Framework Using Binary Inputs', pp. 0–5, 2010.
- [19] F. Dalpiaz, 'Contextual Requirements Prioritization and Its Application to Smart Homes', vol. 1, pp. 94–109, 2017.
- [20] N. Setiani and T. Dirgahayu, 'Clustering Technique for Information Requirement Prioritization in Specific CMSSs', 2016.
- [21] A. Perini, A. Susi, and P. Avesani, 'A Machine Learning Approach to Software Requirements Prioritization', vol. 39, no. 4, pp. 445–461, 2013.
- [22] M. I. Babar, M. Ghazali, D. N. A. Jawawi, S. M. Shamsuddin, and N. Ibrahim, 'Knowledge-Based Systems PHandler: An expert system for a scalable software requirements prioritization process', *KNOWLEDGE-BASED Syst.*, 2015.
- [23] H. Taherdoost and A. Keshavarzsaleh, 'A Theoretical Review on IT Project Success / Failure Factors and Evaluating the Associated Risks', 4th *Int. Conf. Telecommun. Informatics, Sliema, Malta*, no. August, pp. 80–88, 2015.
- [24] R. Prioritization and U. Hierarchical, 'Requirements Prioritization Using Hierarchical Dependencies', pp. 459–464, 2018.
- [25] M. A. Iqbal, A. M. Zaidi, and S. Murtaza, 'A new requirement prioritization model for market driven products using analytical

- hierarchical process', DSDE 2010 - Int. Conf. Data Storage Data Eng., pp. 142–149, 2010.
- [26] X. Frank, Y. Sun, and C. Sekhar, 'Priority assessment of software process requirements from multiple perspectives', vol. 79, pp. 1649–1660, 2006.
- [27] F. Shao, R. Peng, H. Lai, and B. Wang, 'The Journal of Systems and Software DRank : A semi-automated requirements prioritization method based on preferences and dependencies', vol. 126, pp. 141–156, 2017.
- [28] M. Yaseen, I. Journal, M. Yaseen, A. Mustapha, M. A. Salamat, and N. Ibrahim, 'International Journal of Advanced Trends in Computer Science and Engineering Available Online at <http://www.warse.org/IJATCSE/static/pdf/file/ijatcse09912020.pdf> Prioritization of Software Functional Requirements : A Novel Approach using AHP and Spanning Tree', vol. 9, no. 1, 2020.
- [29] S. Ma, J. Li, C. Hu, X. Lin, and J. Huai, 'Big graph search : challenges and techniques', 2015.
- [30] M. Yaseen, A. Mustapha, S. Qureshi, A. Khan, and A. U. Rahman, 'A Graph Based Approach to Prioritization of Software Functional Requirements', vol. 9, no. 3, pp. 64–73, 2020.
- [31] S. Kapoor and H. Ramesh, 'Algorithmica An Algorithm for Enumerating All Spanning Trees of a Directed Graph 1', pp. 120–130, 2000.
- [32] M. Usman, D. Sakethi, R. Yuniarti, and A. Cucus, 'The Hybrid of Depth First Search Technique and Kruskal ' s Algorithm for Solving The Multiperiod Degree Constrained Minimum Spanning Tree', no. Icidm, pp. 0–3, 2015.
- [33] S. Dhingra, 'Finding Strongly Connected Components in a Social Network Graph', vol. 136, no. 7, pp. 1–5, 2016.
- [34] E. Reitsma, P. Hilletoft, and U. Mukhtar, 'Implementation of enterprise resource planning using Odoo module sales and CRM . Case study : PT Ecosains Hayati Implementation of enterprise resource planning using Odoo module sales and CRM . Case study : PT Ecosains Hayati', 2017.
- [35] M. Yaseen, A. Mustapha, and N. Ibrahim, 'Prioritization of Software Functional Requirements : Spanning Tree based Approach', vol. 10, no. 7, pp. 489–497, 2019.
- [36] M. Yaseen, A. Mustapha, N. Ibrahim, and U. Farooq, 'International Journal of Advanced Trends in Computer Science and Engineering Effective Requirement Elicitation Process using Developed Open Source Software Systems', vol. 9, no. 1, 2020.

# Understanding user Emotions Through Interaction with Persuasive Technology

Wan Nooraishya Wan Ahmad<sup>1</sup>  
Ahmad Rizal Ahmad Rodzuan<sup>3</sup>

Faculty of Computing and Informatics  
Universiti Malaysia Sabah, Labuan International Campus  
F.T. Labuan, Malaysia

Nazlena Mohamad Ali<sup>2</sup>

Institute of Industrial Revolution 4.0 (IIR4.0)  
Universiti Kebangsaan Malaysia  
Bangi, Selangor, Malaysia

**Abstract**—Emotions play a vital role in persuasion; thus, the use of persuasive applications should affect and appeal to the users' emotions. However, studies in persuasive technology have yet to discover what triggered the users' emotions. Therefore, the objectives of this study are to examine user emotions and to identify the factors that affect user emotions in using persuasive applications. This study is conducted in three stages; pre-interaction, during-interaction and post-interaction, employed a mixed-method approach using Geneva Emotions Wheel (GEW) and open-ended survey questions that analyzed using thematic analysis. The result shows that most of the emotions that users felt belong to high-control positive valence emotions that consist of interest, joy and pleasure. User, system and interaction are the three factors that triggered the emotions encompasses of elements such as Individual Awareness, Personality, Interface Design, Persuasive Function, Content Presentation, System Quality, Usability, and Tasks. The findings contribute to the body knowledge of Persuasive Technology, where the discovered factors and its elements are the antecedents that should be the concern in constructing an emotion-based trust design framework that could bring emotional impact to users to ensure a successful persuasion.

**Keywords**—Emotion; emotional states; interaction; persuasive technology; captology

## I. INTRODUCTION

Persuasive technology (PT) is increasingly being developed commercially and has been one of the first research area related to shaping human behavior. It is a technology with a purpose to shape and/or change people's attitudes or behavior towards an issue [1]. For example, a persuasive application is developed to help and assist a smoker in becoming a non-smoking person. The persuasion process is to be done without using coercion but using many strategies such as social influence, self-monitoring, and personalization that capable of triggering emotions in users [2]. The PT work is often concerned with the use of practical approaches to address particular behavioral problems [3], hence making the emotional effect that one feels while using PT has been overlooked. Since the persuasion strategies are capable of triggering emotions, it is essential to know what exactly the users feel when using persuasive technology and what makes them feel those emotions. Persuasive technology must bring an emotional impact to the users to ensure the success of the persuasion process.

Recently, studies on emotion have become increasingly important as the need to incorporate emotion into computer application design has become a significant focus of HCI (Park, Lee, & Kim, 2011). However, most studies on emotion in the field of HCI focus more on the emotional stimulation of users during interactions [4][5][6], although studies on emotion can be done at all levels of interactions between computers and humans [7]. Thus, studies on emotions involving all levels of interaction between humans and computers have yet to be discovered.

Hence, the objectives of this study are 1) to investigate the emotions of the user at three interaction stages; pre, during and post, and 2) to identify the factors that trigger user emotions in using persuasive technology. The next section of this paper will further describe the emotion and how it contributes to persuasion as well as persuasive technology. The methodology section explains the comprehensive methods used for data collection, while the data analysis section elaborates on the methods used to analyze the collected data. Next, the result and discussion are presented to explain the discovered scenario. Lastly, the paper ends with a conclusion.

## II. LITERATURE REVIEW

Persuasion involves an attempt to bring about a change in attitudes or behaviors as a result of providing information on a topic or issue without intimidation. Significantly, attitudes are not only relying on thoughts and beliefs but also feelings and emotions [8]. Emotions are one of the components that influence the user experience, which is a medium for users to understand how they feel about using a technology [9]. It is defined as a physiological state of arousal that is triggered by beliefs about something which consists of cognitive physiological, social and behavioral aspects [10]. Emotions are sensitive and are stimulated by situations due to certain circumstances, actions or objects last for a short period [11]. This distinguishes between emotion and mood because the mood lasts for an extended period, which can last for weeks or months, occurring without a specific target and is often separate from what causes the mood to be triggered [12]. According to [13], emotions are stimulated by the culture in which they are produced through actions or tasks and interactions.

---

This study is sponsored by the Fundamental Research Grant Scheme (FRGS) (FRGS/2/2014/ICT01/UKM/02/3) by the Ministry of Education.

Emotions can be categorized in a variety of ways. Discrete emotions, often associated with facial expressions, are described as basic emotions. A study conducted by [14] showed that each of the people recognizes and expresses basic emotions in much the same way as others. There are currently six basic emotions that are universally expressed and recognized throughout the region and culture, which are happy, angry, sad, scared, disgusted and shocked [15]. Dimensional emotions are an exciting category of emotions because they offer a way to describe and distinguish emotional states [16]. Compared to discrete emotions that have only a few emotions, dimensional emotions are subdivided into bipolar dimensions of pleasure-displeasure, arousal, and dominance-submissiveness. Dimensional emotion assesses human feelings in terms of valence from positive to negative and arousal [11]. Measurement of the dimension of arousal from active to passive explains whether or not human actions are affected by emotional states. Dimensional emotions are better able to deal with non-discrete emotions as well as variations in an emotional state over time. For example, the angry feelings are considered to be negative valence with high arousal. At the same time, sadness is associated with negative valence with low arousal, indicating that the same category of emotional valence will have different levels of arousal. In this study, the emotional state's data refers to the stimulus data resulting from the use of persuasive technology that is likely to consist of one or more combinations of emotions.

Emotions can then be modelled as a form of information processing and another set of inputs into cognitive processes [13] where emotion can be a form of internal signal that contributes to cognitive actions to change one's attitude and behaviour. Individuals with positive emotions are more likely to be motivated to change their behaviours [17]. The research showed that emotions and website design influence a person's behavioural shift as a website that uses cue representation to view information is beneficial for individuals with positive emotions. Thus, [18] believes that understanding user emotions is one of the critical components in creating an impactful application apart from putting trust and persuasion in the design. He emphasized that simply providing persuasion technology might at first be useful, but over time, it created a "mess of persuasion", reduced trust and was no longer sufficient. Therefore, designers need to develop persuasive tools based on a deep understanding of the user's emotions.

Thus, the above-reviewed papers show that there is a gap in understanding user's emotions, especially in the interaction with persuasive technology. It is important to know precisely what makes the user feel a particular emotion when using persuasive technology because this might become the antecedents that could bring emotional impact to improve the persuasion process.

### III. RESEARCH METHODOLOGY

The methodology process consists of five parts: (a) participants, (b) material, (c) study design, (d) measures, and (e) setting and task. We present each part as a subsection here.

#### A. Participants

The participants were recruited through advertisements via Facebook. Due to the use of persuasive health applications in the study, such participants have undergone an evaluation using a questionnaire [19] to evaluate their readiness to improve the state of performing physical activity. The volunteers are qualified for the study if only they are in the state of contemplation, preparation, action, and maintenance based on the assessment results. These results not only showed the readiness level of behavior change towards physical activities but also as an indicator of committing to the study. The study was conducted in 6 weeks, with 25 participants managed to complete all three stages of interaction. They consist of 10 males and 15 females. The participants were among the university students and employees. Among the participants, 18 of them had experience in using similar persuasive applications with less than six months and between 6 months to a year experience. The participants' age was in the range from 21 to 45 years old.

#### B. Material

The process of determining appropriate persuasive applications to be studied, ranging from online searching, screening and exclusion. A list of persuasive applications regarding health and environmental was established from Google search activity using different keywords. For example, keywords such as "top health applications in Malaysia" and "top health apps" are used to find persuasive applications on health; meanwhile, keywords such as "game for change" and "persuasive games" are used to search for persuasive applications on the environment.

Three levels of criteria with overall 11 operational variables from [20] were used to screen the 63 persuasive applications gathered from the Internet. The first level of screening is to find applications that fulfil the needed definition of a persuasive application using two operational variables, which are persuasive application and app platform. The second level of screening consists of seven operational variables that were used to find applications that fulfil the specific criteria of a needed persuasive application. The seven operational variables are the theme, type of app, target user, delivery, app availability, interactivity styles, and device collaboration. The third screening level is to ensure that the content of the listed persuasive applications is fit to all range of users regardless of country and education background and meet the local community demands. Two operational variables that were used in the third screening process are app content and the focus of the content.

The exclusion process runs simultaneously with the screening process, where finally, five persuasive applications are selected for this study. All five applications have been identified to qualify as persuasive technology, as mentioned by [1] and [2]. Three from the five applications are about health, and another two applications are on the environment. Both health and environment applications fall under different categories of persuasive applications. Health applications that consist of MyFitnessPal, MapMyFitness and Fitocracy are in tool category that shared the same goal to support and enhance user capabilities to achieve the desired target behavior.

Meanwhile, the selected environment applications fall under medium category shared a goal to showcase the relationship of cause and effect thru simulation consist of Stop Disaster and Pandemic 2 game. Table I describes all five persuasive applications.

### C. Study Design

A mix of between-subject and within-subject design is used to design participants for the experiment. Each participant has the opportunity to use both types of persuasive applications simultaneously within six weeks. Participants were assigned to six groups randomly so that each participant in a group can use a pair of different applications. However, the groups were limited to several participants at a time. Group 1 participants were given a pair of MyFitnessPal and Stop Disaster. Group 2 used MyFitnessPal and Pandemic 2, while participants in group 3 get to used MapMyFitness and Stop Disaster. Meanwhile, participants in group 4 get to used MapMyFitness and Pandemic 2, while group 5 participants get a chance to use Fitocracy and Stop Disaster. Group 6 participants were given a pair of Fitocracy and Pandemic 2.

### D. Measures

A questionnaire that consists of ratings and open-ended questions related to emotional states construct is used as the measurement instrument. For each interaction stage, participants were asked to pick five emotional states to construct which they felt and rated the intensity of the emotions as well as stating the cause that makes them felt the emotions. Geneva Emotion Wheel (GEW) [11] [21] was used to measure emotional states construct. GEW is chosen since it has been used in various fields to study the user's emotions related product and technology. Studies [22] and [23] employed GEW to study emotions towards virtual learning environments. A study by [24] evaluates user emotions related to coffee machines and alarm clock, while in [25] the GEW is used to assess user emotions towards a higher learning website.

GEW composed of 20 discrete emotions arranged parallel in the form of circles according to emotional groups divided into two dimensions of valence (positive and negative) and control (high-low). The division of valence and control dimensions have split the 20 discrete emotions into four groups; positive valence-high control (i.e., interested, amusement, pride, joy, pleasure), positive valence-low control (i.e., contentment, love, admiration, relief, compassion), negative valence-high control (i.e., anger, hate, contempt, disgust, fear), and negative valence-low control (i.e., disappointment, shame, regret, guilt, sadness). The emotional states construct was measured using a 5-point Likert scale, illustrates the intensity of the emotions from low intensity (towards the wheel center) to high intensity (towards the wheel circumference). Five out of 20 emotions from the GEW that were rated by users were listed as the most frequent or dominant emotions that users felt. For each five listed emotions, participants are required to outline the cause that makes them felt those emotional responses.

### E. Setting and Tasks

For this study, user emotional states and also the reasons that triggered the emotions were studied in three interaction stages; pre-interaction, during-interaction, and post-interaction. The pre-interaction stage is defined as an initial interaction with the persuasive applications that they are assigned. During-interaction is a stage where participants used the persuasive applications on their own without the help of others according to their needs and free time. Post-interaction is a stage where the participant has passed a specific period in using the persuasive applications.

The same set of questionnaires each to assess tool and medium type of persuasive applications were distributed in the lab at a specific seat. Upon arrival, participants were asked to fill up the demographic details in the questionnaire. Each interaction stages consist of two sessions of experiment and one relaxing session at the beginning of every session using Calm – a web application that provides scenic pictures and the soothing sound of nature to eliminate stress and stabilizing emotions. When participants tested different applications (AP-1: tool, AP-2: medium), the same experimental procedure was repeated. The following description describes further details of each procedure according to the interaction stage:

- Pre-interaction: It takes about 1 hour to complete two sessions of the experiment. Participants spend 5 minutes for relaxing as a form of control to emotions at the beginning of each session using Calm web application. Before the participants were permitted to have initial interaction with the applications, a demo on how to use the applications was showed to them using slide presentation. In the first session of the experiment that takes about 20 minutes, participants get to explore the AP-1 app based on the tasks given include answering a questionnaire after finishing the tasks. In the next twenty minutes, the participants are required to explore the AP-2 app based on the tasks give and answer a questionnaire after the tasks completed.

TABLE I. PERSUASIVE APPLICATIONS

| Type   | Application   | Description  |
|--------|---------------|--|
| Tool   | MyFitnessPal  | A tracking tool to help weight loss and remain physically fit based on calorie food intake and exercise.   |
|        | MapMyFitness  | A tracking tool to record food intake and exercises by calculating intake and burned calories.   |
|        | Fitocracy     | A social network for fitness that tracks users' fitness progress using gamification principles such as the quest, points and level up to encourage people.                                   |
| Medium | Stop Disaster | A simulation game on the prevention of natural disasters. The game educates users in preventing natural disasters from getting worse and reducing the cost of destruction from the disaster. |
|        | Pandemic 2    | A simulation game on disease spreading. The game teaches users how diseases spread throughout the world.   |

- **During-interaction:** Before the first session of the experiment begun, participants rested for 5 minutes by wearing a headphone to use the Calm web application. With about 20 minutes for each session, participants are required to perform a list of tasks using the AP-1 app and another 20 minutes using the AP-2 app. After the tasks completed, the participants are then required to answer the questionnaires provided, referring to the application they used.
- **Post-interaction:** The study at this stage was conducted in week six, where the lab experiment is no longer necessary. In the two previous stages of the experiment, participants spent 5 minutes' rest at the beginning of each session for emotional relaxation. Afterwards, participants need to answer a provided questionnaire each to evaluate the post-use of AP-1 and AP-2 apps. The participants received a token of appreciation after completed both questionnaires.

In general, participants are required to perform the tasks for AP-1 and AP-2 applications with a minimum addition of new tasks based on the interaction stage. Below are the given tasks:

- **AP-1:** For the pre-interaction stage, participants must first register a user account using either an email or Facebook for the health application (i.e., MyFitnessPal, MapMyFitness, Fitocracy) to which they have been assigned before are allowed setting up a user profile and target goal. Next, participants created physical and/or nutritional activities daily by using the available databases. To start with, participants must record diaries "yesterday" and "today." Finally, participants must see the feedback or advice provided based on the data entered. For the during-interaction stage, participants are generally required to perform the same tasks as the pre-interaction stage. Additional tasks have been established where participants are expected to assess specific functions in AP-1 applications, such as the Body Mass Index (BMI), Body Metabolic Rate (BMR), and other actions that can be taken to obtain feedback or progress reports from the system.
- **AP-2:** In the pre-interaction stage, first, for the Stop Disaster game, participants are required to choose the "easy" game level. For Pandemic 2, participants are required to choose the "virus" role play from the role options such as bacteria and parasites. Within the specified time and/or budget, participants played the game to accomplish a mission. Finally, participants viewed their achievements from the report provided by the system. In the during-interaction stage, for the Stop Disaster game, participants are required to choose the "hard" game level. For the Pandemic 2 game, participants are required to choose the "parasites" role play. Finally, participants are required to complete the game's mission and see their achievements through a report provided by the system.

#### IV. DATA ANALYSIS

The combination of quantitative and qualitative analysis is used to analyze the data from the questionnaire. Using the SPSS software to conduct the quantitative analysis, descriptive analysis was used to analyze user emotional states in three interaction stages as well as the demographic data. Frequency values from the descriptive analysis are used to determine the frequency of users' triggered emotions for each stage of interaction for both persuasive applications. One-way repeated measure ANOVA is used to study the changes in emotional state intensity towards persuasive applications at different interaction stages. This analysis is also used to investigate the changes in user emotions towards persuasive technology as a whole over time. A significant value benchmarked less than .05 is used in determining the statistical significance of the analysis.

For qualitative data, inductive content-analysis coding [26] is used to identify factors that stimulate user emotions when using persuasive technology. It composes four phases; data preparation, data coding, categorization, and theme analysis. In the data preparation phase, questionnaire transcripts for the three interaction stages are arranged by the types of persuasive applications to ensure that the coding process can be distinguished before the general set is generated by merging the existing code. In the coding phase, transcripts text is check and labelled appropriately to create open codes. The inductive analysis method is used in the categorization phase to construct categories based on the obtained open codes. The theme analysis phase performs an overall analysis of the constructed categories to form a theme that shared the same pattern of characteristics.

#### V. FINDINGS

We present the findings into two parts: a) emotional states and b) stimulation factors of emotions to answer the two research objectives.

##### A. Emotional States

Fig. 1 indicates the frequency of user emotions stimulated from the pre-interaction stage for both persuasive applications. From the first impression, the persuasive application AP-1 often arouses emotions such as "interest", "pleasure", "love", "joy", "admiration," and also "disappointment" in the user. Those emotions consist of high-control positive emotions (interest, pleasure and joy), low-control positive emotion (admiration) and also low-control negative emotions (disappointment). Meanwhile, emotions such as "interest", "amusement", "pleasure", "contentment," and "disappointment" dominate the user emotions at first encounter towards persuasive applications AP-2. Table II summarizes the frequency of user emotions for the pre-interaction stage. High-control positive emotions dominate the emotions aroused by AP-1 at the user's first encounter compared to AP-2, where user emotions were dominated by a mixture of both positive and negative emotions but dominated slightly more by high-control positive emotions compared to low-control negative emotions.

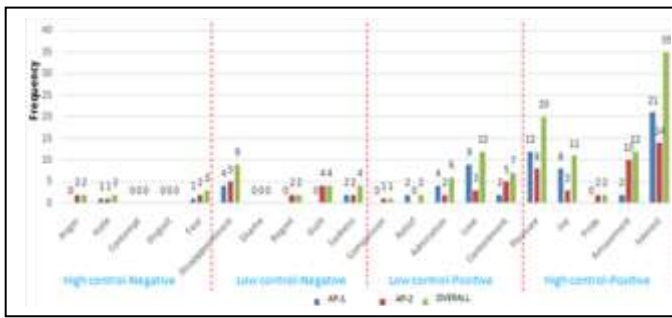


Fig. 1. Frequencies of user Emotions in the Pre-Interaction Stage.

TABLE II. FREQUENCY OF USER EMOTIONS ACCORDING TO GROUPS OF EMOTIONS FOR THE PRE-INTERACTION STAGE

| Emotion Groups                  | Emotion Frequency |      |     | Emotion Intensity Mean |      |      |
|---------------------------------|-------------------|------|-----|------------------------|------|------|
|                                 | AP-1              | AP-2 | ALL | AP-1                   | AP-2 | ALL  |
| High control - Negative Valence | 2                 | 5    | 7   | 2.00                   | 3.33 | 3.17 |
| Low control - Negative Valence  | 6                 | 13   | 19  | 4.13                   | 3.76 | 3.88 |
| Low control - Positive Valence  | 17                | 11   | 28  | 3.66                   | 4.38 | 3.86 |
| High control - Positive Valence | 43                | 37   | 80  | 4.02                   | 4.42 | 4.30 |

Fig. 2 shows the frequency of user emotions in a during-interaction stage. Using AP-1, users often aroused with “interest”, “contentment”, “love”, “relief,” and “disappointment” emotions. These emotions consist of high-control positive emotions (interest), low-control positive emotions (contentment, love and relief) and also low-control negative emotions (disappointment). Meanwhile, the use of AP-2 often stimulated users feeling with emotions of “interest”, “relief”, “sadness”, “disappointment,” and some “amusement”, “joy,” and “pleasure”. All the emotions that the users felt consisted of 3 emotional groups: high-control positive emotions (interest, amusement, joy and pleasure), low-control positive emotions (relief) and low-control negative emotions (sadness, disappointment). Table III summarizes the frequency of user emotions at the during-interaction stage. Positive emotional groups dominate the emotions aroused by both AP-1 and AP-2 with high control.

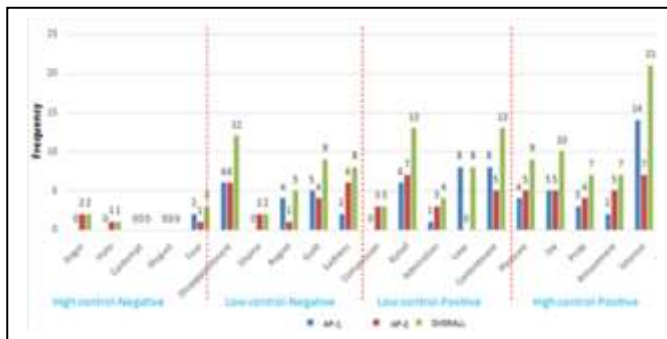


Fig. 2. Frequencies of user Emotions in the During-Interaction Stage.

TABLE III. FREQUENCY OF USER EMOTIONS ACCORDING TO GROUPS OF EMOTIONS FOR DURING-INTERACTION STAGE

| Emotion Groups                  | Emotion Frequency |      |     | Emotion Intensity Mean |      |      |
|---------------------------------|-------------------|------|-----|------------------------|------|------|
|                                 | AP-1              | AP-2 | ALL | AP-1                   | AP-2 | ALL  |
| High control - Negative Valence | 2                 | 4    | 6   | 2.50                   | 4.00 | 3.89 |
| Low control - Negative Valence  | 17                | 8    | 25  | 3.41                   | 3.75 | 3.66 |
| Low control - Positive Valence  | 23                | 9    | 32  | 3.86                   | 3.75 | 3.85 |
| High control - Positive Valence | 28                | 10   | 38  | 3.90                   | 4.05 | 3.95 |

Fig. 3 indicates the frequency of user emotions aroused in the post-interaction stage after six weeks’ usage of both persuasive applications. For AP-1, the aroused emotions consist of the feeling of “interest”, “love”, “pleasure”, “contentment”, and “disappointment”. All of these emotions consist of two groups of positive emotions with high-control (interest, pleasure) and low-control (love, contentment) and also low-control negative emotions (disappointment). For AP-2, user emotions often aroused with the feeling of “interest”, “amusement”, “hate”, “disappointment”, “pleasure”, and “admiration”. Table IV summarizes the frequency of user emotions for the post-interaction stage. The analysis shows that user emotions in using AP-1 are more dominated by positive emotions, but are more often dominated by high-control positive emotions than low-control emotions. Compared to AP-1, users’ emotions towards AP-2 are dominated by low-control emotions but are dominated less by positive emotions than negative emotions.

Additionally, one-way repeated measure ANOVA is conducted to identify the changes in the intensity of user emotional states towards the persuasive applications at Time 1 (pre-interaction), Time 2 (during-interaction) and Time 3 (post-interaction). The mean values and standard deviation are presented in Table V. Result of the analysis shows no statistically significant effect towards time (interaction stages), Wilk’s Lambda = 0.93,  $F(2, 43) = 1.53$ ,  $p > 0.1$ . The result indicates that different interaction stages bring no impact on the intensity of user emotional states.

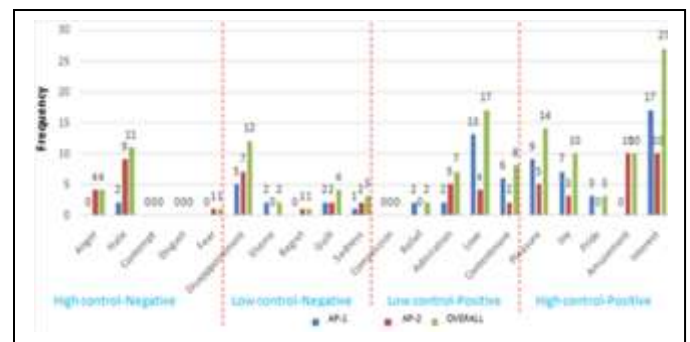


Fig. 3. Frequencies of user Emotions in the Post-Interaction Stage.



TABLE IV. FREQUENCY OF USER EMOTIONS ACCORDING TO GROUPS OF EMOTIONS FOR POST-INTERACTION STAGE

| Emotion Groups                  | Emotion Frequency |      |     | Emotion Intensity Mean |      |      |
|---------------------------------|-------------------|------|-----|------------------------|------|------|
|                                 | AP-1              | AP-2 | ALL | AP-1                   | AP-2 | ALL  |
| High control - Negative Valence | 2                 | 14   | 16  | 5.00                   | 4.14 | 4.17 |
| Low control - Negative Valence  | 10                | 17   | 27  | 2.40                   | 3.96 | 3.48 |
| Low control - Positive Valence  | 23                | 18   | 31  | 4.15                   | 4.32 | 4.14 |
| High control - Positive Valence | 36                | 9    | 45  | 4.04                   | 3.59 | 3.90 |

TABLE V. EMOTIONS DESCRIPTIVE STATISTIC ON THE INTENSITY OF USER EMOTIONAL STATES TOWARDS THE USED OF PERSUASIVE APPLICATIONS AT EACH INTERACTION STAGES

| Time frame                  | N  | Mean | Std. Dev. |
|-----------------------------|----|------|-----------|
| Time 1 (pre-interaction)    | 45 | 4.10 | 0.80      |
| Time 2 (during-interaction) | 45 | 3.86 | 0.94      |
| Time 3 (post-interaction)   | 45 | 3.92 | 1.05      |

However, the changes in user emotions as a whole towards the persuasive applications at different stages of interaction; Time 1 (pre-interaction), Time 2 (during-interaction) and Time 3 (post-interaction) show a different result. Table VI indicates the mean and standard deviation values for each interaction stage. There is a statistically significant effect towards interaction stages, where Wilk's Lambda = 0.48,  $F(2,47) = 25.77$ ,  $p < 0.001$ . The result proved that user emotions changed over time. Nevertheless, the mean value of Table VI explained that even users experienced different emotions at various stages of interaction; the variations in the strength of these emotions were not apparent. Three paired samples t-test were conducted to make post hoc comparisons between the stages. A first paired samples t-test indicated that there was a significant difference between user emotions in pre ( $M=3.00$ ,  $SD=0.70$ ) and during ( $M=3.82$ ,  $SD=0.93$ ) interaction stages;  $t(48) = 5.83$ ,  $p=0.00$ . A second paired samples t-test indicated that there was a significant difference between user emotions during ( $M=3.82$ ,  $SD=0.93$ ) and post ( $M=3.00$ ,  $SD=0.79$ ) interaction stages;  $t(48) = 6.86$ ,  $p=0.00$ . However, a third paired samples t-test indicated that there was no significant difference between user emotions in pre ( $M=2.97$ ,  $SD=0.71$ ) and post ( $M=3.02$ ,  $SD=0.79$ ) interaction stages;  $t(49) = 0.34$ ,  $p=0.74$ .

TABLE VI. DESCRIPTIVE STATISTIC OF EMOTIONAL CHANGES TOWARDS THE USED OF PERSUASIVE APPLICATIONS

| Time frame                  | N  | Mean | Std. Dev. |
|-----------------------------|----|------|-----------|
| Time 1 (pre-interaction)    | 49 | 3.00 | 0.70      |
| Time 2 (during-interaction) | 49 | 3.82 | 0.93      |
| Time 3 (post-interaction)   | 49 | 3.00 | 0.79      |

### B. Stimulation Factors of Emotions

Table VII shows the coding and categorization for each interaction stage. In total, eight categories have been identified that stimulates emotional states: Individual Awareness, Personality, Interface Design, Persuasive Function, Content Presentation, System Quality, Usability and Task. All eight categories were grouped according to the same characteristics based on the factors of assessment of emotional control by [27] and [28] where three factors were identified: User control, System control and Interaction control. User Control is a users' assessment of their emotions on the basis related to themselves. System Control referred to user emotions that were evoked by the persuasive applications itself- particularly the features used to deliver the persuasion process, while Interaction Control referred to user emotions that were evoked by the interaction between the user and the persuasive applications.

User Control factors consist of two categories: Individual Awareness and Personality. Individual Awareness refers to self-assessment that causes the individual to be aware of what is happening to them, and the individual is aware that he or she is experiencing an event, exhibiting behaviour or having unique characteristics [29]. Individual Awareness constitutes aspects of Consciousness Value and Self-Satisfaction. The definition of the two aspects are as follows:

TABLE VII. CODING AND CATEGORIZATION OF EMOTIONS STIMULATION FACTORS FOR EACH INTERACTION STAGES

| Theme               | Category                  | Open-Code                                     |   |   |                     |
|---------------------|---------------------------|---|---|---|---------------------|
|                     |                           | Pre   | During  | Post  |                     |
| User control        | Individual Awareness      | Self-satisfaction                             | Self-satisfaction<br>Consciousness value      | Self-satisfaction<br>Consciousness value      |                     |
|                     | Personality               | Interest<br>Knowledge<br>Skill<br>Relatedness | Interest<br>Knowledge<br>Skill<br>Relatedness | Interest<br>Knowledge<br>Skill<br>Relatedness |                     |
| System control      | Interface design          | Interface attractiveness                      | Interface attractiveness                      | Interface attractiveness                      |                     |
|                     |                           | Layout  | Layout  | Layout  |                     |
|                     | Persuasive function       | Persuasive function                           | Persuasive function                           | Persuasive function                           |                     |
|                     | Content presentation      | Information quality                           | Information quality                           | Information quality                           | Information quality |
|                     |                           |   | Data representation                           | Data representation                           | Data representation |
|                     |                           |   | Multimedia                                    | Multimedia                                    | Multimedia          |
| System Quality      | Reliability<br>Usefulness | Reliability<br>Usefulness                     | Reliability<br>Usefulness                     |   |                     |
| Usability           | Learnability              | Learnability                                  | Learnability                                  | Learnability                                  |                     |
|                     |                           | Ease of use                                   | Ease of use                                   | Ease of access                                |                     |
| Interaction control | Task                      | Feedback<br>Action<br>Expectation             | Feedback<br>Action<br>Expectation             | Feedback<br>Action<br>Expectation             |                     |

- As pointed by [29], individuals need to awake to determine the processing of information either internally or externally. Thus, Consciousness Value refers to the use of a system of alerting users to something they know or do not know. Some of the responses from the participants are as followed:

“It makes me feel relief as I finally know how to deal with natural disaster (well at least basic knowledge)”.

(P3, during-interaction, AP-2)

“I like this application a lot because it helps me manage my weight and exercises.”

(P21, post-interaction, AP-1)

- Self- Satisfaction refers to the achievement achieved by using the system. Achievements achieved by the user through the use of the system either through the activities performed or as a result of the activities performed will stimulate the user’s emotions. The emotions associated with the achievement of an activity or outcome achievement is defined as the emotion of achievement [30]. Some of the responses from the participants are as followed:

“Able to do something that would improve my health and give a good example to my children.”

(P2, pre-interaction, AP-1)

“It makes me feel happy as it teaches me how to react during a natural disaster”

(P3, during-interaction, AP-2)

Personality refers to attributes that distinguish users from using a system. The unique attributes of these individuals influence the user’s emotions towards the system being used. Four aspects referred to Personality in the use of a system which is: Interest, Skill, Knowledge and Relatedness. The following defines the meaning of the four aspects:

- Interest refers to the user’s tendency to use the system. The user tends to have a positive attitude towards the system being used and to encourage the user to use the system. Some of the positive responses given by users are as follows:

“This application makes me interested in doing physical activity.”

(P13, pre-interaction, AP-1)

“Still interest doing the same every day. Very useful information and update version.”

(P8, post-interaction, AP-1)

- Knowledge refers to the knowledge that the user has concerning the system used. For example, if a user does not know the number of calories they eat, it is tough for the user to use applications such as MyFitnessPal that requires users to record the total of the calorie intake of their meals. Here are some responses expressed by participants:

“The calorie input part needed me to estimate the caloric intake for each food that I ate.”

(P26, pre-interaction, AP-1)

“Still unable to put in money for house development after spending for defense.”

(P18, post-interaction, AP-2)

- Skill refers to the ability of the user to control and use the system. When playing Stop Disaster games, users should think of strategies to prevent flooding in lowlands by applying appropriate development. Some of the responses provided by users are as follows:

“I still unable to finish the game successfully within 10-15 minutes’ time frame.”

(P23, during-interaction, AP-2)

“All menus easy to be completed.”

(P5, during-interaction, AP-1)

- Relatedness refers to the shared goals shared between users and the system. Users to lose or control weight have a trusting relationship with the system used and allow for internalization to occur [31]. Some of the responses that illustrate Relatedness are as followed:

“I always look forward to how to maintain my stamina, and the app is giving a lot of interesting training and good and relevant advice from the trainer.”

(P22, pre-interaction, AP-1)

“I feel the system is interesting to me because I want to lose weight...”

(P27, post-interaction, AP-1)

System Control factors consist of four categories: Interface Design, Persuasive Function, Content Presentation, System Quality, and Usability. Interface Design refers to the visual layout of the elements that the user will use to interact with the system. It encompasses two aspects, namely Interface Attractiveness, and Layout. The details of each of these aspects are as follows:

- Interface Attractiveness refers to the use of aesthetic values in system design. For example, the use of dark color and skull images made users feel uncomfortable playing Pandemic 2 game. The example of responses given by the participants was as follows:

“The website has a pleasant look.”

(P26, during-interaction, AP-1)

“I regret playing this game plus the color is not interesting.”

(P11, during-interaction, AP-2)

- Layout refers to interface layouts that are easy to understand. Unstructured interface sets make it difficult for users to find the information or

functionality they want, and this will trigger negative emotions in the system being used.

“...selecting the game features has become easier.”

(P29, during-interaction, AP-2)

“I hate the background of the game, the layout of the game also quite boring.”

(P8, post-interaction, AP-2)

Persuasive Function refers to the functions provided by the system for users to fulfil their tasks in order to achieve their goals. These functions act as a tool that assists the user to monitor their activities progress using the persuasive applications.

“This application has all the functions that I needed, e.g., calorie counter, BMI calculator.”

(P14, pre-interaction, AP-1)

Content Presentation refers to the approach in the delivery of information used by the system in delivering information to users. This category covers three aspects, namely Multimedia, Information Quality and Data Representation. The descriptions of the three aspects are as follows:

- Multimedia refers to the use of multimedia elements to capture the attention of users. The multimedia aspect is more focused on AP-2 targeting applications, games that use simulation methods as a medium to persuade users of the issues being highlighted.

“It amazed me when the flood occurred; I can see good animation that simulates the flood.”

(P16, during-interaction, AP-2)

“It is truly amusing with all graphic impact...”

(P1, post-interaction, AP-2)

- The Information Quality referring to the information provided by the system is complete and accurate.

“Provide knowledge-information in handling flood.”

(P19, pre-interaction, AP-2)

“As this app is telling me more about my diets and all details what I have done with me, I am avoiding the things/food not good for me.”

(P9, during-interaction, AP-1)

- Data Representation refers to the method of data visualization or data for display to the user.

“The app shows result in a good way (graph)”

(P5, pre-interaction, AP-1)

“The game shows the gradual increase in disease spread starting from green color to red.”

(P5, post-interaction, AP-2)

System Quality is a factor that refers to the overall system’s behavior to control and to function correctly in terms

of Reliability and Usefulness aspects. The details of each of these aspects are as follows:

- Reliability refers to a user’s trust towards the system in ensuring that the user can achieve the targeted goal. For examples, the users of Pandemic 2 believed the game in delivering information regarding the diseases spreading. At the same time, through the Fitocracy app, it enables the user to monitor physical activities as one of the ways to lose weight.

“The app enables me to feel good with myself.”

(P26, during-interaction, AP-1)

“It provided many exciting experiences playing in trying to achieve a better result.”

(P1, post-interaction, AP-2)

- Usefulness refers to the system able to help the user achieved its target goal.

“The game allows me to gain knowledge about preventing and handling flood disaster.”

(P19, during-interaction, AP-2)

“I can now practice my favorite exercise with help from this app.”

(P17, post-interaction, AP-1)

Usability is a factor in a system that refers to the quality of attributes in assessing how easy for the user to use the system. It encompasses three aspects which were Learnability, Ease of Use and Ease of Access. The details of each of these aspects are as follows:

- Learnability refers to the ability of the system to allow the user to learn how to use it. According to [32], learnability is when the system is easy to learn.

“I thought I had grasped the idea of how to tackle this game and make the virus more lethal, but it was not proven in my previous game. But in this game, there is an improvement, so I think I have some confidence in my strategy now.”

(P23, during-interaction, AP-2)

“Having difficulties winning the game even at the easiest level.”

(P17, post-interaction, AP-2)

- Ease-of-Use means the system is secure for users to use. This definition is in line with the definition used by [32] and [33] where the system is generally easy to use.

“This app is user-friendly.”

(P4, pre-interaction, AP-1)

“It is easy to use this. I love this app.”

(P3, during-interaction, AP-1)

- The Ease-of-Access aspect refers to the situation in which the system used is easy for the user to access, for

example, through websites that can be accessed using computers and mobile phones or through internet access. Overall, the users have made the accessibility aspect as a factor for the post-interaction phase. Here are the responses provided by the participants:

“It is hard to play the app because it requires Internet connection all the time.”

(P7, post-interaction, AP-2)

“I can use the app anywhere as long that I have Internet.”

(P11, post-interaction, AP-1)

Interaction Control factor consists of Tasks category. Tasks refer to the activities that the user is allowed to do while interacting with the system. The three aspects of Tasks are Action, Feedback, and Expectation. The descriptions of these aspects are as follows:

- Actions are selections of action that users can perform to get feedback from the system. Users are given the freedom to make choices in determining the feedback they get from the system. Among the responses expressed by participants were:

“The app can track record from previous days and allows user to edit and re-edit.”

(P5, pre-interaction, AP-1)

“I found many helpful options that maintained my interest.”

(P9, post-interaction, AP-1)

- Feedback is the response and advice given by the system that stimulates emotions and motivates the user. Among the responses expressed by participants about the action:

“Got zero death and zero injuries.”

(P1, during-interaction, AP-2)

“I like how the app helps us to monitor our health by showing the graph of food intake for the whole month.”

(P11, post-interaction, AP-1)

- Expectation refers to the expectancy that a user can make from using the system, which involves things that users expect the system to do. The following are responses expressed by participants regarding expectations:

“I wish that this game allows for more user control on spreading the diseases.”

(P29, pre-interaction, AP-2)

“The app cannot give a prediction of weight after listing all the foods that the user had eaten.”

(P5, post-interaction, AP-1)

## VI. DISCUSSIONS

### A. Emotional States

This study reveals that, in the pre-interaction stage, most users experience the feeling of “interest”, “pleasure”, “amusement”, “love,” and “joy”. All of these emotions are high-control positive emotions except for “love,” which are low-control positive emotions. Responses from the users towards the AP-1 and AP-2 has always related to the aspect of Interest, the tendency of using the system, making many of the users choosing “interest” and “pleasure” emotion as the emotion that they experienced since interest is the emotion that is activated when one experiences engagement, interest, and curiosity [34]. The findings suggest that people decide to adopt a system when they feel pleasant or having positive thought about the system.

The study found that, in during interaction stage, users experience a variety of emotions mostly the feeling of “interest”, “contentment”, “relief”, “disappointment”, “joy”, “pleasure”, and “guilt”. These emotions are contained in high-control positive emotion (i.e., interest, pleasure, joy), low-control positive emotion (i.e., relief, contentment) and low-control negative emotion (i.e. disappointment, guilt). Findings from the study suggest that System Quality, Content Presentation, Interface Design, Personality, Individual Awareness and Task are the elements that triggered those emotions.

For the post-interaction stage, the study reveals that users experience more emotions of “interest”, “love”, “pleasure”, “disappointment” and “hate”. Those triggered emotions are a combination of four groups, which are high-control positive emotion (i.e., interest, pleasure), low-control positive emotion (i.e., love), low-control negative emotions group (i.e., disappointment) and high-control negative emotion group (i.e., hate). The findings suggest that all eight elements affect users’ emotions at this particular stage. Overall, the triggered emotions for the three interaction stages are dominated by emotions from a high-control positive emotions group indicating that the persuasive applications used in the experiment triggered positive emotions among the users in the persuasion process. These findings agreed with previous studies by [35] and [36] that persuasion is perceived when the user felt positive emotions.

### B. Stimulation Factors of Emotions

In using both types of persuasive applications, it is found that the emotions experienced by the user are under the control of all three factors; User, System and Interaction which encompasses eight elements namely Individual Awareness, Personality, Interface Design, Persuasive Function, Content Presentation, System Quality, Usability, and Tasks. However, what distinguishes these two types of persuasive applications are the aspects that stimulate user emotion found through the open code.

For AP-1 persuasive applications, the pre-interaction stage shows 18 aspects that dominate the user’s emotions except for Multimedia, Ease of Use and Ease of Access elements. However, in the during-interaction stage, aspects such as Skills, Layout, Data Representation, Multimedia, Learnability,

and Ease of Access are not affected by the emotions of the user, and this results in only 14 aspects of the user's emotions. However, in the post-interaction stage, 17 aspects succeeded in affecting the user's emotions except for aspects such as Skills, Layout, and Multimedia. Overall, the multimedia aspect was not seen as a contributing factor to the emotion experienced by users in the use of persuasive applications of tools.

For AP-2 persuasive applications, 12 aspects were found to affect the user's emotions in the pre-interaction stage except for aspects such as Layout, Data Representation, Multimedia, Reliability, Ease of Use, Ease of Access, as well as two aspects of Individual Awareness that include Consciousness Value and Self-Satisfaction. However, in the during-interaction stage found 17 elements that impacted the user's emotions except for aspects such as Persuasive Function, Ease of Use, and Ease of Access. The post-interaction stage also found 17 aspects that influenced the user's emotions except for aspects such as Skills, Learnability, and Ease of Use. Overall, it was found that the Ease of Use aspect was not a contributing factor to the emotions that users experienced in using persuasive applications that fall under the medium category. This finding also found that Ease-of-Access aspect was only evaluated and considered relevant by users in the use of both types of persuasive applications in the post-interaction stage only.

### C. Relation between the Emotional States and Stimulation Factors

The findings show that "interest", "pleasure", and "joy" are the positive emotions that dominate users at each interaction stage. At the same time, "disappointment" is the negative emotions that always triggered in users across the interaction stages. The negative emotions that users felt were caused by the aspects of Task, Individual Awareness, Personality and Interface Design that cover all three factors; User Control, System Control and Interaction Control. The comparison with the outcomes of previously user emotion study [25] reveals a similar pattern in the emotions that users experience even by assessing the different types of stimuli, which in this case a higher learning institution website designed using standard Kansei-based guidelines. However, the study did not identify what specifically trigger those positive emotions. Although the emotion of "interest" in the present study receives the highest values in every interaction stages, this does not mean that the same thing will happen whenever emotions are directed to a tool or device; nevertheless, this offers a test of the significance of these measures to design-related research. The emotion of "interest" is usually correlated with innovative practices, the growth of skills and knowledge, the learning of new abilities, and consistency in an effort [34]. While most previous studies [22] [23] [25] successfully evaluates user emotions to address the user response towards the design of technology, this study, at the same time, identifies the factors and elements that are addressing those user emotions. The findings show that the highlighted elements are the potential antecedents that affect the trust of the user, hence becoming a strategy that needs to be addressed for constructing an emotion-based trust design framework for persuasive technology.

## VII. CONCLUSION

This study has contributed to the findings on what exactly users feel when using a persuasive application and the factors that triggered the experienced emotions to the body of knowledge of Persuasive Technology field interest. User emotions in each interaction stage can vary because of different factors such as User, System and Interaction and the elements of the factor, for examples Individual Awareness, Personality, Interface Design, Persuasive Function, Content Presentation, System Quality, Usability, and Tasks. The findings can be used to construct a design framework of persuasive technology that can bring an emotional impact to the user by focusing on the factors with appropriate design principles or strategies. The designing framework, however, will be reported in the other paper, after triangulation with diary studies' results.

## ACKNOWLEDGMENT

Thank you to the Malaysian government for granting funds on this work, which is part of the research under the Fundamental Research Grant Scheme (FRGS) (FRGS/2/2014/ICT01/UKM/02/3) by the Ministry of Education.

## REFERENCES

- [1] B. J. Fogg, *Persuasive Technology Using Computers to Change What We Think and Do*. Morgan Kaufmann: San Francisco, 2003.
- [2] H. Oinas-Kukkonen and M. Harjuuma, "Persuasive systems design: key issues, process model, and system features," *Communications of the Association for Information Systems*, vol. 24, no. 28, pp. 485-500, Mar. 2009.
- [3] J. Hamari, J. Koivisto, T. Pakkanen, "Do persuasive technologies persuade? A review of empirical studies," in *Persuasive Technology. PERSUASIVE 2014. Lecture Notes in Computer Science*, vol. 8462, A. Spagnolli, L. Chittaro, and L. Gamberini, Eds. Cham, Springer, 2014, pp. 118-136.
- [4] N. Mohamad Ali., S. Z. Abdullah, J. Salim, R. Sulaiman, H. B. Zaman and H. Lee, "Exploring user experience in game using heart rate device," *Asia-Pacific Journal of Information technology and Multimedia*, vol.1, no. 2, pp. 28-36, 2012.
- [5] H. Petrie and J. Precious, "Measuring user experience of websites: think aloud protocols and an emotion word prompt list," *CHI 2010*, pp. 3673-3678, April 2010.
- [6] F. Zhou, X. Qu, M. G. Helander and J. Jiao, "Affect prediction from physiological measures via visual stimuli," *International Journal of Human-Computer Studies*, vol. 69, no. 12, pp. 801-819, 2011.
- [7] I. Lopatovska, and I. Arapakis, "Theories, methods and current research on emotions in library and information science, information retrieval and human-computer interaction," *Information Processing and Management*, vol. 47, no. 4, pp. 575-592, 2011.
- [8] R. E. Petty and P. Brinol, "Emotion and persuasion: Cognitive and meta-cognitive processes impact attitudes," *Cognition and Emotion*, vol. 29, no. 1, pp. 1-26, 2014.
- [9] J. Forlizzi and K. Battarbee, "Understanding experience in interactive systems understanding experience in interactive systems," *5th Conference on Designing interactive systems 2004*, pp. 261-268, August 2004.
- [10] D. Lottridge, M. Chignell and A. Jovicic, "Affective interaction: Understanding, evaluating, and designing for human emotion," *Reviews of Human Factors and Ergonomics*, vol. 7, no. 1, pp. 197-217, 2011.
- [11] K. R. Scherer, "What are emotions? And how can they be measured?" *Social Science Information*, vol. 44, no. 4, pp. 695-729, 2005.
- [12] T. Partala, "Affective information in Human-Computer Interaction," Ph.D. dissertation, Department of Computer Sciences, University of Tampere, Finland, 2005.

- [13] K. Boehner, R. Depaula, P. Dourish, P. Sengers, A. C. Zaidan and U.C. Irvine, "How emotion is made and measured," *International Journal of Human-Computer Studies*, vol.65, pp. 275–291, 2007.
- [14] P. Ekman and W. V. Friesen, *Unmasking the face. A guide to recognizing emotions from facial clues.*, Prentice-Hall: Englewood Cliffs, New Jersey, 1975.
- [15] C. Darwin, *The expression of the emotions in man and animals*, Cambridge University Press, 2009.
- [16] I. Hupont, S. Baldassarri and E. Cerezo, "Facial emotional classification: from a discrete perspective to a continuous emotional space," *Pattern Analysis & Applications*, vol. 16, no. 1, pp. 41–54, 2013.
- [17] E. Krahmer, J. V. Dorst and N. Ummelen, "Mood, persuasion and information presentation. The influence of mood on the effectiveness of persuasive digital documents," *Information Design Journal*, vol. 12, no. 3, pp. 40–52, 2004.
- [18] E. Schaffer, "PET Research: Looking Deeper to Understand Motivations," 2010. <https://humanfactors.com/downloads/whitepapers/PET-research.pdf>.
- [19] B. H. Marcus and L. H. Forsyth, *Motivating People to be Physically Active: Physical activity intervention series*, Human Kinetic, Champaign: IL, 2003.
- [20] W. N. Wan Ahmad, N. Mohamad Ali, "The impact of persuasive technology on user emotional experience and user experience over time," *Journal of Information Communication Technology*, vol. 17, no. 4, pp. 601–628, 2018.
- [21] K. R. Scherer, V. Shuman, J. R. J. Fontaine and C. Soriano, "The GRID meets the wheel: Assessing emotional feeling via self-report," in *Components of emotional meaning: A sourcebook*, J. R. J. Fontaine, K. R. Scherer, and C. Soriano, Eds. Oxford: Oxford University Press, 2013, pp. 281–298.
- [22] J. M. dos. Santos, "Gaia: intelligent control of virtual environments," PhD dissertation, Universidade Tecnica de Lisboa, 2008. <https://dspace.ist.utl.pt/bitstream/2295/236570/1/Dissertacao da Tese de Mest%0Arado - FINAL - Revised - Jorge Santos.pdf>.
- [23] M. T. Longhi, D. F. Pereira, M. Bercht, & P. A. Behar, "An experiment to understand how the affective aspects can be detected in virtual learning environments," *CINTED-UFRGS*, vol.7, 2009.
- [24] G. F. G. Laurans, "On the moment-to-moment measurement of emotion during person-product interaction," Technische University of Delft, 2011.
- [25] P. Turumugon, A. Baharum, N. H. Nazlan, N. A. M. Noh, N. A. M. Noor and E. A. Rahim, "Users' emotional evaluation towards kansei-based higher learning institution website using Geneva Emotion Wheel," *Indonesian Journal of Electrical Engineering and Computer Science*, vol. 16, no. 3, pp. 1547–1554, 2019.
- [26] D. Silverman, *Interpreting Qualitative Data.*, London: Sage Publications Ltd., 2014.
- [27] J. R. Dunn and M. E. Schweitzer, "Feeling and believing: the influence of emotion on trust," *Journal of Personality and Social Psychology*, vol. 88, no. 5, pp. 736–748, 2005.
- [28] J. So, C. Achar, D. Han, N. Agrawal, A. Duhachek and D. Maheswaran, "The psychology of appraisal: Specific emotions and decision-making," *Journal of Consumer Psychology*, vol. 25, no. 3, pp. 359–371, 2015.
- [29] A. Morin, "Levels of consciousness and self-awareness: A comparison and integration of various neurocognitive views," *Consciousness and Cognition*, vol. 15, pp. 358–371, 2006.
- [30] R. Pekrun, A.J., Elliot. And M.A., Maier, "Achievement goals and discrete achievement emotions: A theoretical model and prospective test," *Journal of Education Psychology*, vol. 98, pp. 583–597, 2006.
- [31] R. M. Ryan, H. Patrick, E. L. Deci and G. C. Williams, "Facilitating health behaviour change and its maintenance: interventions based on Self-Determination Theory," *The European Health Psychologist*, vol. 10, pp. 2–5, 2008.
- [32] J. Xu, K. Le, A. Deitermann and E. Montague, "How different types of users develop trust in technology: A qualitative analysis of the antecedents of active and passive user trust in a shared technology," *Applied Ergonomics*, vol. 45, pp. 1495–1503, 2014.
- [33] Z. Yan, R. Kantola and P. Zhang, "A research model for human-computer trust interaction," in *2011 IEEE 10th International Conference on Trust, Security and Privacy in Computing and Communications*, pp. 274–281, November 2011.
- [34] C. E. Izard, *The Psychology of Emotions*. New York: Plenum Press, 1991.
- [35] V. Griskevicius, M. N. Shiota and S. L. Neufeld, "Influence of different positive emotions on persuasion processing: a functional evolutionary approach," *Emotion*, vol. 10, no. 2, pp. 190–206, 2010.
- [36] J. V. Riet, R. A. C. Ruiter, M. Q. Werrij, M. J. J. M. Candel and H. de Vries, "Distinct pathways to persuasion: The role of affect in message - framing effects," *European Journal of Social Psychology*, vol. 40, no. 7, pp. 1261–1276, 2010.

# Decision-Making Analysis using Arduino-Based Electroencephalography (EEG): An Exploratory Study for Marketing Strategy

Ahmad Faiz Yazid<sup>1</sup>, Siti Munirah Mohd<sup>2\*</sup>, Abdul Razzak Khan Rustum Ali Khan<sup>3</sup>  
Shafinah Kamarudin<sup>4</sup>, Nurhidaya Mohamad Jan<sup>5</sup>

Advanced Technology and Sustainable Research Group, Kolej GENIUS Insan<sup>1,3</sup>  
Universiti Sains Islam Malaysia (USIM), Negeri Sembilan, Malaysia<sup>1,3</sup>

Senior Lecturer, Kolej GENIUS Insan, Universiti Sains Islam Malaysia (USIM), Negeri Sembilan, Malaysia<sup>2,5</sup>

Senior Lecturer, Department of Science and Technology, Faculty of Humanity, Management and Science<sup>4</sup>  
Universiti Putra Malaysia (UPM) Bintulu Campus, Bintulu, Sarawak<sup>4</sup>

**Abstract**—Business technology has brought conventional marketing methods to the next level. These emerging integrated technologies has contributed to the growth and understanding of the consumer decision making process. Several studies have attempted to evaluate media content, especially on video advertising or TV commercials using various neuroimaging techniques such as the electroencephalography (EEG) device. Currently, the use of neuroscience in Malaysia's marketing research is very limited due to its limited adoption as an emerging technology in this field. This research uncovers the neuroscientific approach, particularly through the use of an EEG device; examining consumers' responses in terms of brain wave signals and cognition. A proposed theoretical framework on the factors affecting visual stimulus (movement, color, shape, and light) during the decision-making process by watching video advertising had been customized using two conceptual models of sensory stimulation. Ten respondents participated in the experiment to investigate the spectral changes of alpha brain wave signals detected in the occipital lobe. A 2-channel Arduino-based EEG device from Backyard Brains and Spike Recorder software was used to analyze the EEG signal through Fast Fourier Transform (FFT) method. Results obtained from the investigated population showed that there was statistically significant brain wave activity during the observation of the video advertising which demonstrated the interconnection with short-term memory through visual stimulus. Application of the neuroscience tool helped to explore consumer brain responses towards marketing stimuli with regards to the consumers' decision-making processes. This study manifests a useful tool for neuromarketing and concludes with a discussion, together with recommendations on the way forward.

**Keywords**—*Arduino-based electroencephalography (EEG); neuromarketing; short-term memory; TV commercials; visual cognition*

## I. INTRODUCTION

Promoting products is a significant part of business, which is a relentless and ever-evolving condition. The increasing customer density and changing methods of doing business has resulted in huge amounts of data which contrasts with the business research condition from a couple of decades back. For example, the idea of 4Ps (product, place, price, and

promotion) in advertising exercises includes a progression of instruments which organizations use to accomplish their objectives [1]. The promoting blend alludes to an assortment of different devices that organizations regularly use to obtain the ideal reaction from their customers [2]. This practice has been broadly embraced in the field of marketing promotions, particularly for showcasing procedures [3]. Likewise, video-based advertisements or commercials can be viewed as one of the quickest developing forms of online advertising which happen within customers' inner thoughts, practically 80% of new items being promoted will fade into obscurity within a year [4]. Obtaining a better insight into the purchaser's dynamic through procedures via a neuroscientific approach assist organizations in making educated business choices through the assessment of the business viability. It is a basic tenement for advertisers to investigate and grasp human behavior in order to be able to lessen customer dissatisfaction while increasing customer loyalty in order to increase revenue and income stream stability [5].

Up until recent times, the vast majority of organizations still depend on conventional promotional devices, such as studies, examinations, and center gatherings, all of which take a considerable amount of time. These methods are all aimed at trying to provide comprehension in regard to consumer reactions in order to be able to build a better product to satisfy the customer [6]. Thus, in order to reduce the amount of time required to gauge customer reactions while increasing accuracy of data, the application of neuroscience is seen as a possible advancement in marketing methodology. New neuroscience advances have empowered organizations connect more intimately with customers to better cater to their needs. Neuroscience aids in the analysis of a brand or product's potential prior to promotions, which is of particular value to organizations in terms of validity of data and reduction in costs [7]. The advent of "neuromarketing" is fast overtaking statistical surveys as the preferred means to gain customer insight [8].

On the report of Google Trends in 2020, the trend of interest on neuromarketing searched worldwide started from 2004 until present (October 2020) is increasing significantly.

\*Corresponding Author

The findings from several research related to neuromarketing are becoming increasingly important in providing meaningful data to the marketing industry as the area of interest continues to grow. Taken as a whole, it can be inferred that the basic concept of neuromarketing is the use of neuroscience tools to better grasp consumer behavior and to analyze the efficacy of advertising actions from the marketing stimulation, covering the aspects of unconscious and emotional response.

Based on the inquiry targets, the use of neuromarketing can be used by both business activity and research field. Specific to the field of consumer behavior, neuromarketing provides important perspectives and methods into consumer research [9]. This is confirmed by researchers who stated that neuromarketing offers a broad range of knowledge, aggregating in the subconscious mind to behavioral data which guides consumers in their decision-making process [10]. Thus, it was also reported that the research of neuromarketing will measure emotional engagement, memory retention, purchasing intention, novelty, consumer awareness and interest, positioning, product design and creativity, advertisement efficacy, decision-making, online experience, and entertainment performance [11].

In particular, neuroscience helps to understand the role of consumers' inner emotional responses, which play an important part in decision making process [12]. This evidence is supported by the idea of neuroscience seeks to understand the underlying complex of thoughts such as reasoning, decision making, object representation, emotion, memory and consumer responses to marketing [13]. Neuroscience is also associated with the basis for understanding how consumers create, store, recall and relate to information such as brands in everyday life [14]. The findings obtained from studies in neuromarketing give insight into the attitude of the customer that traditional marketing research approaches cannot deliver.

The utilization of neuroscience in Malaysia's is currently limited as it is a new field with few proponents. The fast adoption and advances in the field of neuromarketing in Europe and the United States will inevitably lead to a scramble to adopt this method in Malaysia [15]. Therefore, fast tracking localized studies in this particular field becomes of importance to ensure sufficient preparedness when the time comes. Advertisers need to familiarize themselves with the concepts and practices required of neuroscience in order to stay competitive in providing for the consumers wants and needs.

This research endeavors to inspect local consumers' reactions towards video advertising in terms of their decision-making process; through the use of the neuroscientific method utilizing electroencephalogram (EEG). The present study hopes to encourage advertisers to be more alert towards the customer decision making process in order for them to produce higher impact video-based advertisements which will be more effective in terms of marketing. The use of neuromarketing will reduce disappointments by aiding analysts and organizations to create more effective content which aligns with customer requirements.

This study addresses three main questions as follows; (i) How can neuromarketing be applied effectively in the

processing of visual information through video-based advertising?; (ii) What happens in the human brain during the purchase making decision process through advertising?; (iii) Is it possible to identify the roots of observed and noticed consumer behavior patterns in processes of human brain activity?

The objectives of this study are; (i) To investigate the spectral changes of alpha brain wave signal detected from 2-channel Arduino-based EEG device using the Fast Fourier Transform (FFT) method; and (ii) To identify the contributing factors from video advertising in the activation of brain wave signal towards decision making processes by testing the cognitive response of short-term memory.

The paper is structured as follows. Section II presents the proposed theoretical framework used in this study with relevant literary highlights along with critical and creative reviews towards understanding the research concept. Section III shows the process of acquisition of the EEG signal from the human brain to the computer system by using Spike Recorder software. The project flow in terms of the application of working principles and procedures in conducting the experiment is demonstrated in detail. Section IV discusses the result of this study. The discussion involves the results of the EEG signal which was acquired from test subjects and the result of signal processing through the Fast Fourier Transform (FFT) method. The final discussion details how the EEG signal is affected by several factors applied in video-based advertisements. Finally, Section V and VI explain the conclusion and recommendations for this project. Further research is also addressed pertaining to marketing strategy and other factors involved in advertising required to activate the brain wave signal of consumers towards their decision-making process.

## II. THEORETICAL FRAMEWORK

The proposed framework (Fig. 1) is constructed to show the variables corresponding to the research. This framework is designed by taking into account the work of other researchers and furthering their concepts. The model can be divided into parameter (catalytic factors) selection from the two conceptual models of sensory stimulation (Sensory Stimuli Model [16] and Map of Dynamic Stimuli [17]). The selection of parameters focuses mainly on visualization as it is a major stimulus associated with consumer cognition when watching video advertising or TV commercials. Information related to visualization is analyzed in greater detail to determine factors affecting the stimulus.

It is readily apparent throughout various evidences that assessments of consumer behavior in purchasing decisions were affected by visual stimuli, in which relevant to the usage and quantity in product selection [18]. This results to the role of visual stimuli seem to be more significant, even in the absence of verbal information about an item. Overall, it is easier to concentrate on an item with the use of graphic imagery. This is particularly true in a competitive clutter situation, where visuals provide a sense of value and nurture strong engagements with a brand, thus, contribute to the possible decision to purchase due to positive influence on consumer judgment and purchase decision.



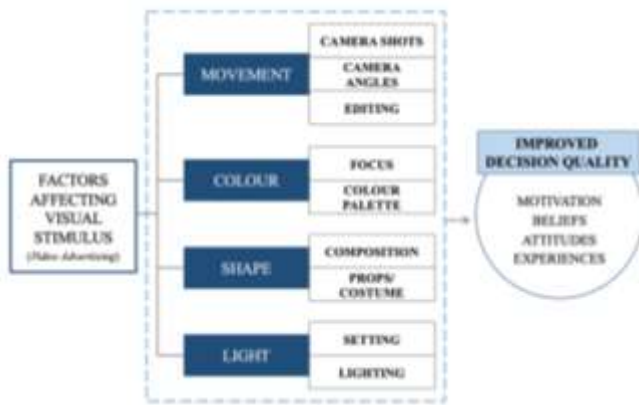


Fig. 1. Theoretical Framework of Relationship between Independent Variable and Dependent Variable.

This particular framework contains both dependent and independent variables. The dependent variable is the decision-making process of consumers (observed through the activation of brain wave signal analysis – neuroscience technique). The independent variable is the factors affecting the visual stimulus, consisting of movement, color, shape, and light.

### III. MATERIALS AND METHODS

#### A. Experimental Design

A total of ten (10) respondents, comprising five male subjects and five female subjects, with an age range of between 16 and 17 years, participated as part in this investigation. All respondents were secondary school students from Kolej PERMATA Insan, Universiti Sains Islam Malaysia (USIM). The participants were well acquainted with promotional videos and notable brands. All participants did not have any history of clinical disease, for example, cognitive/mental/psychological disorders. Prior to the start of testing, an introductory presentation regarding the specific exploration, time length, and motivation behind the study was presented to respondents.

#### B. Experimental Procedures

The experiment was carried out according to stringent protocols with non-disclosure clause and permissions for gathering EEG signals obtained from the participants. A quiet room with dim lighting was used to avoid external interruption. Participants were seated following an ordered seating scheme and the LCD screen was placed 5 meters away from them. The EEG device was introduced to participants at this time and a short clarification about the test was relayed to them orally. The headband of the 2-channel Arduino-based EEG gadget was directly connected to a laptop (Macbook Pro) pre-installed with an open-access version of EEG recording analysis (Spike Recorder) software. All connections were examined prior to the test to ensure no disruptions and full functionality during the analysis.

The trial comprised of four segments; one session for preparation, two sessions for the EEG recording, and an additional session for the survey. Prior to the start of each session, a full white background picture was screened to the respondents to ensure an unbiased state condition. When the

brain wave signal of each respondent reached a constant level, the video advertisement of a notable well-known fast food brand (McDonald's Malaysia) was introduced and participants were made to view a repeat of the commercial during the third session, whereby the commercial was played exactly 30 seconds after which a period of 2.5 minutes was allocated to stabilize brain waves back into a neutral condition. After the repeated viewing of the advertisement, a short-self assessment was conducted. During the final session, respondents completed a poll to assess the commercial. During this final stage respondents were asked to review and recall the specific scenes from the TV commercial for the purpose of evaluating their short-term index. Finally, respondents were required to state or choose the reasons why a particular scene was memorable enough to trigger their short-term memory.

#### C. Data Analysis

The acquired EEG signals were firstly pre-processed using the Arduino Uno microcontroller and Fast Fourier Transform (FFT) method. This was done to derive the statistical features of the alpha wave recorded from each subject. The extracted features were subsequently matched with the film segments to investigate the spectral changes of the subjects' brain wave signal regarding the product advertisements.

The Fourier Transform are the data analysis programs used for the purpose of generating spectrograms, which showed the changes in frequency content of a signal over time. An open-source Arduino Code was used to send data from the computer board running the macOS Catalina (Version 10.15) system. The resulting brain wave signals from this process were detected and analyzed for this study. For the Spike Recorder software, a spectrogram within a certain frequency band was set up to appear below the moving EEG trace of the brain wave signal. The present study also applied a qualitative research methodology through non-probability sampling using a semi-structured survey among the ten individuals who were subjected to the recording of brain wave signals using the EEG device. This survey was conducted for approximately 15 minutes after the EEG experiment to determine the subjects' perceptions and cognitions with regards to the video advertisement. This survey evaluated the subjects' cognitive responses through short-term memory index in order to discover the factors in video advertising that could influence the decision making of consumers via visual stimulus.

### IV. RESULTS AND DISCUSSIONS

This section depicts the investigation of the spectral changes in the alpha brain wave signal detected from the 2-channel Arduino-based EEG device using Fast Fourier Transform (FFT) method and the identification of the contributing factors from video advertising in the activation of brain wave signal towards the decision making process by testing the cognitive response of short-term memory for all participants.

#### A. Spectrum Analysis using Fast Fourier Transform (FFT) Method

The breakdown of results collected from the participants is shown through a progression of figures demonstrating the measurably huge contrasts of brain wave signal activated for

the dataset at the recurrence band (10–40 Hz) of spectrographic analysis. The figures comprise a progression of result boards, each containing two pictures; the top figure represents a frame of the video advertising whereas the lower one shows the corresponding mean of brain wave activity, and the temporal axis beat the time of commercial. In Fig. 2, the primary arrangement of seven film sections taken at every five second from the earliest starting point of the commercial, navigating the entire length of a specific commercial spot. Thorough auditing of each strip highlight shows the tracking of the mean brain wave activity changes according to scenes in the video advertisement viewed by the respondents.

The difference in brain wave activity becomes increasingly obvious in Fig. 3, made out of three boards representing the first, the middle, and the last frame of the video advertisement, respectively. The corresponding mean of the brain wave activity completes each board of the figure. By observing these three panels, it is clear how the middle part of the commercial shows activation of occipital zone with measurable contrasts, while there are two peaks of activity (greater reoccurring green spots) towards the beginning and the end of the clip; revealed the increases of alpha wave power in the brain wave signal.

The investigation of changes in brain wave activity was performed even on shorter spans in order to track their variations over shorter time periods. Subsequently Fig. 4 shows the brain wave activity within the constant frequency band of the spectrogram analysis. The time intervals across the following figures correspond to the starting five seconds, middle five seconds, and the final five seconds of the commercial. These models demonstrate how it is conceivable to obtain measurably critical contrasts in recording alpha brain wave activities of the occipital zone in any event, which in turn decreases the time frame.

It seems clear that there is a great deal of interest in applying new neuroscience research method that have been developed among advertisers to know on how TV commercials actually work. The technique of picture-sorting has been suggested to identify potential branding moments in TV commercials [19], where it proved the picture-sorting technique synchronized to the brain wave measurement and linked to the theory of memory. This discussion has also been supported through the latest study that found positive correlation between user's evaluations and two neuroscience technologies of EEG and eye tracker [20].

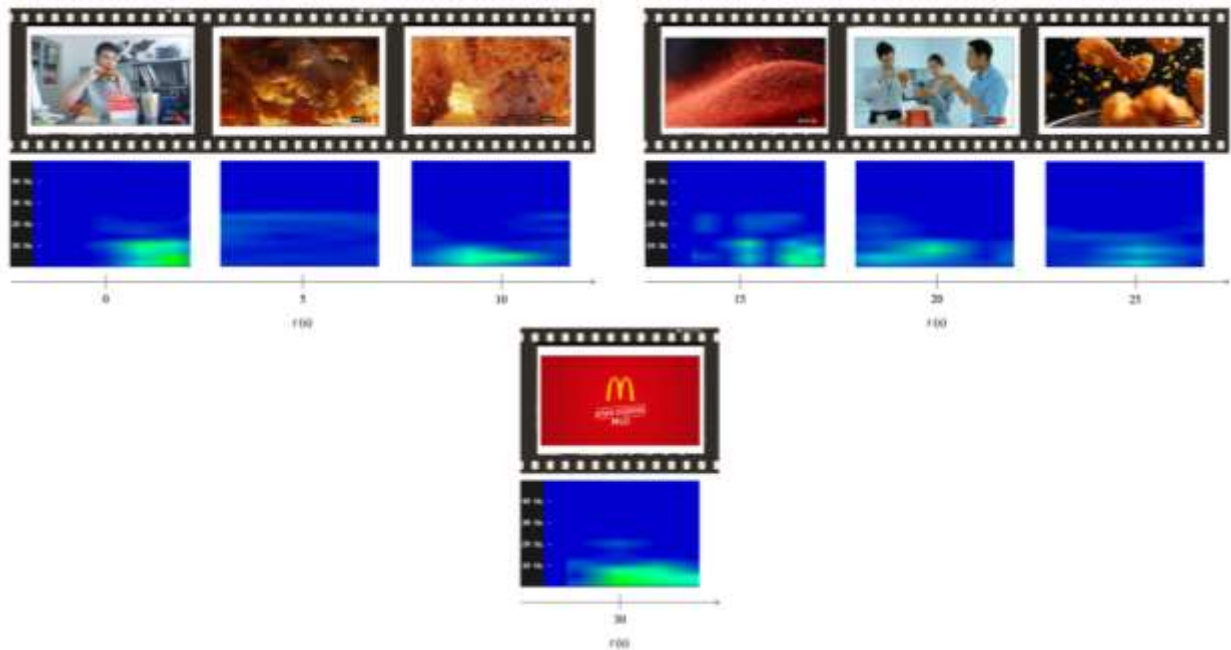


Fig. 2. Tracking of the Brain wave Activity in the Frequency Band within the Time Spots of Every Five Seconds, Across Seven Film Segments.

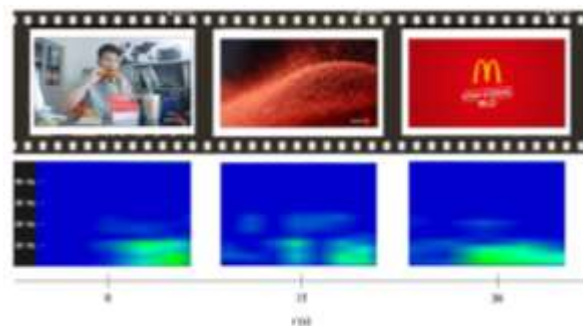


Fig. 3. Tracking of the Brain wave Activity Across Three Panels of the Total Film Sequence.



Fig. 4. Tracking of the Brain wave Activity in the Frequency Band for the First Five Seconds, the Middle Five Seconds, and the Last Five Seconds of the Commercial Spot.

The implementation of neuroscience approach in advertising is growing. It is believed that frequency of branding product exposure in a particular TV commercial might contribute towards the effectiveness of the advertising. This statement is parallel to the experiment that indicated containing multiple exposures of the branding products had a positive impact on the preference for the commercial, in which resulted to higher brain wave activity within several major regions including the frontal, bilateral occipital, and limbic system [21]. Ultimately, there is a need for research to help marketers cultivate effective and memorable marketing campaigns, hence, this study design will contribute to the consumer's decision-making process.

#### B. Cognitive Testing for Short-Term Memory

At the start of the survey, respondents needed to furnish their demographic details. The profile of the respondents is shown in Table I below. Based on Table I, the proportion of male to female respondents were equal 50%:50%. The mean age distribution of the respondents was the age group of above 16 years old as the dominant age group (f=9, 90.0%) compared to the age group of below 17 years old (f=1,

10.0%). All respondents (100.0%) were familiar with video advertising or TV commercials.

Apart from the demographic background questions, all respondents were asked two other questions. Firstly, respondents were asked to rate the advertisement based on several qualities. Secondly, they were asked how memorable the advertisement was. Responses were rated on a five-point Likert-type scale ranging from '1' to '5', which represents in ascending order strongly disagree, disagree, partially agree, agree, and strongly agree. Participants responses towards the commercial are as shown in Table II.

TABLE I. DEMOGRAPHIC PROFILE OF RESPONDENTS

|                                  |          | Frequency | Percentage (%) |
|----------------------------------|----------|-----------|----------------|
| Gender                           | Male     | 5         | 50.0           |
|                                  | Female   | 5         | 50.0           |
| Age                              | > 16 y/o | 9         | 90.0           |
|                                  | < 17 y/o | 1         | 10.0           |
| Familiar with video advertising? | Yes      | 10        | 100.0          |
|                                  | No       | 0         | 0.0            |

TABLE II. REGRESSION ANALYSIS OF THE SURVEY

|  | Mean | SD   |
|--|------|------|
| Q1) Rate the advertisement based on the following: |      |      |
| Comprehension                                      | 4.50 | 0.25 |
| Message  | 4.00 | 0.45 |
| Impact   | 4.80 | 0.16 |
| Clarity  | 4.20 | 0.40 |
| Q2) How memorable was the advertisement?           | 4.90 | 0.30 |

From Table II, it is shown that the majority of respondents chose strongly agrees or agrees options when filling out the questionnaire. All the means computed are above 4. All of the questions have a mean equal to a minimum of 4.00 and maximum of 4.90, which denotes positive responses from participant. Next respondents were asked to list the most memorable scene from the video advertisement they were shown. Based on the responses, the theme and issues of the description are linked and tailored to the four factors affecting visual stimulus as outlined in the theoretical framework of this study.

Table III shows the detail of the responses from each memorable scene according to the factors outlined in the proposed framework. There were five respondents (Respondents 1, 3, 5, 7, and 8) who stated that the element of movement an obvious factor being displayed in the video advertisement. Individual elements include camera shots, camera angles, and editing combined to help consumers to mark and remember specific scenes.

A total of nine respondents expressed their agreement that the color factor is what makes the video advertisement memorable. The focus and color palette significantly affect the color gradient of the whole video frame. The factor of shape consists of the composition of props and costumes worn by the characters in specific scenes. This factor controls how the scene is set or staged. Six respondents noted that the shape factor affected their memory for certain scenes in the video. Another factor which was highlighted by seven respondents is the light factor; this is in line with the research framework. Lighting is a very important aspect of communication which influences the subconscious perception of what the commercial is showing.

As a whole, the EEG technique of spectrogram analysis was deemed to have enabled the researcher to track subjects' brain wave activity while they were observing the commercial. In such a manner, it is possible to acquire a broad measure of the reconstructed brain wave signals by means of a simple FFT method. This allows the researcher to distinguish the different activities within alternating advertising scenes. To sum it up, the present study managed to analyze the brain wave detected brain wave activities, as a result of visual stimulus in response to the video advertisement. This provides proof that the neuroscience technique of EEG can be implemented for marketing research, involving the measurement of consumer's brain wave activity.

Subsequently, it was found that respondents provided differing views on the factors affecting their cognitive responses of short-term memory regarding the video advertisement. All of the four factors are described in depth by

the majority of respondents in terms of the scenes they remember throughout the 30 videos. It was found that the majority of respondents (f=9, 90.0%) agreed that color, including the elements of focus and color palette, played a major role in garnering consumer responses, through the triggering of their short-term memory. In contrast, movement factors such as camera shots, camera angles, and editing garnered the least agreement amongst respondents (f=5, 50.0%) as an influence on consumer's preference.

TABLE III. MEMORABLE ELEMENTS IN THE VIDEO ADVERTISEMENT BASED ON RESPONDENTS' RESPONSES

| Factors (Elements) | Responses  |
|--------------------|--|
| Movement           |  |
| Camera Shots       | "The scene of a female with hijab shows a very interesting close-up. The head and shoulders of the character are clearly focused." (Respondent 3, 7, and 8)  |
| Camera Angles      | "The placement of a chicken bucket at the end of the video has been captured at high angle. The camera is placed above eye level, looking down." (Respondent 1 and 7)  |
| Editing            | "The continuity of each character's scene with a close-up, everyone biting the chicken. When one shot ends and another one begins." (Respondent 5)   |
| Color              |  |
| Focus              | "Footage of people biting crunchy chickens is shown in shallow focus. When one part of the image is in focus, and another part is not." (Respondent 1, 4, 5, 6, 8, and 9)<br><br>"Special effect for chicken being fried in oil. When the overall frame is blurry, out of focus." (Respondent 2, 3, 5, and 10) |
| Color palette      | "Trademark of McDonald's logo at the end of commercial. The range of colors chosen for the scene is well blend; red and yellow." (Respondent 1, 4, and 10)   |
| Shape              |  |
| Composition        | "The first scene describes how things are positioned in the frame." (Respondent 3 and 4)<br><br>"The role of husband and wife characters expresses strong family bond, sharing chicken together." (Respondent 2 and 5)   |
| Props/ Costumes    | "The items used for background and anything worn by actor is appropriate, shot in the office." (Respondent 1 and 6)  |
| Light              |  |
| Setting            | "Footage of chicken is cooked in a frying pan. Encloses where a scene takes place." (Respondent 2 and 7)   |
| Lighting           | "Special effect customized in the scene of chicken covered with spices. Like the chicken being placed in a spotlight." (Respondent 1, 2, 3, 4, 5, 7, and 9)  |

## V. FUTURE SCOPE

Thanks to neurotechnological advancements, analysis of consumer brain wave movement is presently achievable, which means both researchers and advertisers have a suitable

method to gather greater comprehension of consumers behavior. The utilization of an EEG device is paving the way for future consumer research and will improve on the traditional methods using statistical surveying methods.

From the research and post-experiment interviews conducted, researcher can conclude that an advertisement which is more likely to attract and influence consumer's purchasing making decision is the one that fulfills several factors affecting the visual stimulus. "Look More, Like More". This demonstrated by the brain wave activation detected from alpha wave signal in the occipital lobe. The effects of purchaser's behavior, promotion, marketing, pricing, product circulation, and decision making can be analyzed from a much more scientific base Further research on consumer cognition and their visual attention will allow researchers to have greater understandings of human behavior for use in a wide array of fields including marketing, health care, personal traits and wellness.

Further research in this field is expected to motivate researchers, academicians, and professionals to further develop the field of marketing and advertising research of consumers through the application of neuroscience.

## VI. CONCLUSION

This research presents the application of an EEG device as a neuroscience technique in the field of neuromarketing, in conjunction with the development of a theoretical framework based on the two previous conceptual models of sensory stimulation, particularly regarding the visualization aspect. The findings of this research are the results of the fulfillment of the two objectives set at the beginning of the study.

The first objective of this research is to investigate the spectral changes of alpha brain wave signal detected from a 2-channel Arduino-based EEG device using Fast Fourier Transform (FFT) method. The fulfillment of the first objective answers the first and second research questions. The main finding of this research in the context of the first objective is the method of tracking the mean of alpha brain wave activity through spectrographic analysis, as a result of visual stimulus in response to the video advertisement.

The second objective is to identify the contributing factors from the video advertisement in the activation of brain wave signals towards the decision-making process which was triggered by the cognitive response of short-term memory. The implementation of this second objective addresses the third research question. The research finding in the context of the second objective is the evidence of the effectiveness of the theoretical framework development derived from the two preceding conceptual models. Factors affecting visual stimulus underlined in the framework (movement, color, shape, and light) were proven to be major guidelines for the production of video advertisement designed to influence the consumer decision making process.

The theoretical framework of this research shows proof that the visual attention of human beings is reflective of their purchasing behavior. It is hoped that this study will improve the quality of the decision-making process, influenced by the cognitive response of short-term memory. Consumers'

experiences throughout their lives are registered in their memory which means that their prior experiences which are collected in the first stage of memory can have a big impact on their future experiences in terms of how they react and respond.

The current study has usefulness beyond academia as it will be of great use to the industry, in practical terms the study provides a novel application of the EEG technology in the field neuromarketing. The results and methods demonstrated in this study will be helpful to those involved in the field of marketing and advertising. This will provide a more efficient method to obtain usable consumer data to increase marketing effectiveness.

## ACKNOWLEDGMENT

The authors would like to acknowledge the contributions of the anonymous reviewers whose valuable comments and insightful suggestions led to the improvement of the preliminary version of this paper.

## REFERENCES

- [1] Alizade, R., Mehrani, H., & Didekhani, H. (2014). A Study on the Effect of Selected Marketing Mix Elements on Brand Equity with Mediating Role of Brand Equity in Etka Chain Stores-Golestan Province. Kuwait Chapter of Arabian Journal of Business and Management Review, 3(11), 184-193.
- [2] Khan, M.T. (2014). The Concept of 'Marketing Mix' and its Elements (A Conceptual Review Paper). International Journal of Information, Business and Management, 6(2), 95- 107.
- [3] Hedging, T., Knudtzen, C. F., & Bjerre, M. (2009). Brand Management - Research, Theory and Practice. New York: Routledge.
- [4] Cruz, C. L., de Medeiros, J. F., Hermes, L. R., Marcon, A., & Marcon, É. (2016). Neuromarketing and the advances in the consumer behaviour studies: A systematic review of the literature. International Journal of Business and Globalisation, 17(3), 330-351.
- [5] Mehta, B., & Panda, R. (2015). Neuromarketing-Contour between the Proximate and the Ultimate level of Consumer Decision Making. IFRSA Business Review, 5(1), 1-5.
- [6] McDowell, W. S., & Dick, S. J. (2013). The Marketing of Neuromarketing: Brand Differentiation Strategies Employed by Prominent Neuromarketing Firms to Attract Media Clients. Journal of Media Business Studies, 10(1), 25-40.
- [7] Lassere, A. (2014, October 26). The Marketing Corner: 'Brave New World'. Retrieved from The Epoch Times: [https://www.theepochtimes.com/the-marketing-corner-brave-new-world\\_1043889.html](https://www.theepochtimes.com/the-marketing-corner-brave-new-world_1043889.html)
- [8] Ciprian-Marcel, P., Lăcrămioara, R., Ioana, M. A., & Maria, Z. M. (2009). Neuromarketing – Getting Inside the Customer's Mind. Annals of Faculty of Economics, 18(1), 804-807.
- [9] Genco, S. J., Pohlmann, A. P., & Steidl, P. (2013). Neuromarketing for Dummies. Mississauga: John Wiley & Sons Canada, Ltd.
- [10] Colaferro, C. A., & Crescitelli, E. (2014). The Contribution of Neuromarketing to the Study of Consumer Behavior. Brazilian Business Review, 11(3), 123-143.
- [11] Sebastian, V. (2014). Neuromarketing and Evaluation of Cognitive and Emotional Responses of Consumers to Marketing Stimuli. Procedia - Social and Behavioral Sciences, 127, 753-757.
- [12] Solnais, C., Andreu-Perez, J., Sánchez-Fernández, J., & Andréu-Abela, J. (2013). The contribution of neuroscience to consumer research: A conceptual framework and empirical review. Journal of Economic Psychology, 36, 68-81.
- [13] Perrachione, T. K., & Perrachione, J. R. (2008). Brains and brands: Developing mutually informative research in neuroscience and marketing. Journal of Consumer Behaviour, 7(4-5), 303 - 318.

- [14] Ahmad, Z. A. (2010). Brain in Business: The Economics of Neuroscience. *Malaysian Journal of Medical Sciences*, 17(2), 1-3.
- [15] Abd Hamid, A. I., Abdullah, J. M., & Fauzan, N. (2018). The Future of Cognitive Neuroscience. *International Journal of Engineering & Technology*, 7(3.22), 1-4.
- [16] Esmailpour, H., & Zakipour, M. (2016). The Sensory Stimuli Model; Engage with the Consumer Senses for Brand Distinguishes. *Journal of Management Sciences*, 2(4), 212-218.
- [17] Colombo, S., Gomo, R., & Bergamaschi, S. (2013). Enhancing Product Sensory Experience: Cultural Tools for Design Education. *International Conference on Engineering and Product Design Education* (pp. 698-703). Dublin, Ireland: Dublin Institute of Technology.
- [18] White, K., Habib, R., & Hardisty, D.J. (2019). How to SHIFT Consumer Behaviors to be More Sustainable: A Literature Review and Guiding Framework. *Journal of Marketing*, 83(3), 22-49.
- [19] Young, C. (2002). Brain Waves, Picture Sorts®, and Branding Moments. *Journal of Advertising Research*, 42(4), 42-53.
- [20] Wang, L. (2019). Test and Evaluation of Advertising Effect Based on EEG And Eye Tracker. *Translational Neuroscience*, 10, 14-18.
- [21] Wang, R. W. Y., Chang, Y., Chuang, S. (2016). EEG Spectral Dynamics of Video Commercials: Impact of the Narrative on the Branding Product Preference. *Scientific Reports*, 6, 36487.

# Finger Movement Discrimination of EMG Signals Towards Improved Prosthetic Control using TFD

E.F. Shair<sup>1\*</sup>, N.A. Jamaluddin<sup>2</sup>, A.R. Abdullah<sup>3</sup>

Centre of Excellence for Robotic, and Industrial Automation (CERIA)  
Department of Electrical Engineering, Faculty of Electrical Engineering  
Universiti Teknikal Malaysia Melaka, Hang Tuah Jaya, Durian Tunggal, 76100 Melaka, Malaysia

**Abstract**—Prosthetic is an artificially made as a substitute or replacement for missing part of a body. The function of the missing body part can be replaced by using the prosthesis and it can help disabled people do their activities easily. A myoelectric control system is a fundamental part of modern prostheses. The electromyogram (EMG) signals are used in this system to control the prosthesis movements by taking it from a person's muscle. The problem for the myoelectric control system is when it did not receive the same attention to control fingers due to more dexterous of individual and combined finger control in a signal. Thus, a method to solve the problem of the myoelectric control system by using time-frequency distribution (TFD) is proposed in this paper. The EMG features of the individual and combine finger movements for ten subjects and ten different movements is extracted using TFD, ie. spectrogram. Three machine learning algorithms which are Support Vector Machine (SVM), k-Nearest Neighbor (KNN) and Ensemble Classifier are then used to classify the individuals and combine finger movement based on the extracted EMG feature from the spectrogram. The performance of the proposed method is then verified using classification accuracy. Based on the results, the overall accuracy for the classification is 90% (SVM), 100% (KNN) and 100% (Ensemble Classifier), respectively. The finding of the study could serve as an insight to improve the conventional prosthetic control strategies.

**Keywords**—*Electromyography; feature extraction; time-frequency distribution; spectrogram; classification; machine learning*

## I. INTRODUCTION

Nowadays, the world just not depends on current science and medicine but instead it also creates a variety of new technologies. Among the attention of the world is the creation of electronic tools that assist in physiotherapy facilities. In the field of physiotherapy, many tools can help to guide paralyzed or disabled patients during rehabilitation training such as prosthetic hand and leg.

Furthermore, the first prosthetic hand and leg which is used for treatments, especially in physiotherapy, was introduced in the early 16th century. Prosthetic is an artificial made as a substitute or replacement for missing part of the body due to accident or permanent disablement. This prosthetic can help disabled people to do their work or activities easily. Our body uses the muscles to control the limb movement. However, in prosthetic, electromyogram (EMG) signals from an individual muscle are used instead [1].

The EMG signals are to record the electrical activity of muscles. This signal knows the condition of muscles and nerve of the body when movements exist. However, the EMG signals can be affected by several factors, especially during data collection [2]. Thus, several methods can be used to get better accuracy of surface EMG signals for a prosthetic hand.

Currently, there is a high technology that can create and manufacture prosthetics which is used to replace the loss part of the body and being normal again. Nowadays, artificial limbs have advancements in the materials used and the design of artificial. There is for enhancements and comfortable use when using the prosthetic. Also, the electronics have been used as new materials and become common in artificial limbs. The myoelectric limbs have become more common than cable operated limbs to control the limbs. The myoelectric has been used by converting the muscle movements to electrical signals. The myoelectric used the electrodes to convert the signals of muscle movements to electrical signals. However, there still have some technical problems in the process of capturing or analyzing the data [3]. The myoelectric signals can be triggered by internal and external disturbances.

Recent attempts have been made to obtain more dexterous human finger power, given the success of using EMG signals in interpreting the expected forearm gestures. For example, using surface EMG signals to determine when the finger is active and which finger is enabled using only two electrodes positioned on the forearm. There was an experiment that used two electrodes to detect four finger movements by using time distribution and neural networks of good accuracy. However, the performance of the time distribution (TD) features are not satisfactory even though time consumption and dimensions of TD is faster and smaller [4]. On the other hand, frequency distribution (FD) features can be difficult to detect EMG signals for stroke subjects due to the lower power frequency at muscle contraction [5].

Besides that, the classification of individual single finger movement is common but there are only several types of research that have been made for the classification of multiple individuals and combine finger movement in the same finger. To recognize the EMG signals from different classes of the finger movements, a suitable classifier must be employed in the system.

\*Corresponding Author

## II. RELATED WORKS

### A. The Problem of Prosthetic Hand

In the field of medicine research, the necessary function can be continuing by the prosthesis. This new technology and behaviors in the field of a prosthesis can be the modern treatment of diseases such as diabetes, stroke, and peripheral artery disorder. The patient continues to be the final common denominator. The prosthetic is made to replace the function of missing limbs either to walk or moving depending on their desire.

There have a limited performance for prosthetic tools where it is performed with one specific activity or perform by bimanually while the prosthetic hand can perform with multiple activities and tasks. In contrast to the mechanical appearance of prosthetic tools, prosthetic hands appear human-like.

Hybrid devices can be used for high-level amputation for example at or above the elbow with a combination of body-powered and myoelectric elements. This system can control two joints at once where one is body-powered and the other one is myoelectric. This device is cheaper and smaller than a prosthetic composed entirely of EMG controlled components. However, the myoelectrical prosthetic has a disadvantage where the prosthetic is heavy, expensive than the other type of prosthetic, and depending on usage and power consumption to operate. Other than that, prosthetic, not 100% reliable because the EMG sensor sometimes gets "misreads" the user intent when it attached to the skin. Currently, the prosthetic on the market is not a full feedback loop and the input proprioceptive sensor is not fed back to the natural neural pathways of the user [6].

### B. EMG Signals

Electromyogram (EMG) is the electrical activity to define nerve and muscle problems in response to a stimulation of the muscle of nerve. During the test, a small needle (electrodes) are used to pick up the electrical activity through the skin into muscle and displayed on a monitor in a waveform. The EMG measured the muscle during rest, slight contraction, and forceful contraction for the electrical activity. There are two methods to measure EMG signals: invasive and noninvasive. For invasive methods, it uses needle electrodes while a noninvasive method uses electrodes above the skin surface of the patient's body [7].

EMG is an analytical technique involved in the development, recording, and study of myoelectric signals. Myoelectric signals are formed by physiological changes in the state of the membranes of muscle fiber [8]. EMG signals have a wide range of applications in biomedical engineering and it is one of the vital biological parameters, prosthetic devices, and rehabilitation devices [9]. It is a bio-potential signal acquired through the muscle fiber body by electrodes to analyze muscle activity [9] and these signals measure the electrical activity during contraction and relaxation phase of the muscle fiber [7].

The EMG has also been used to find the effect of symptoms such as muscle weakness, deformity, stiffness, and

shrinkage. Other than that, EMG is also used to test the problem of the motor like involuntary muscle twitching and nerve compression, injuries such as carpal tunnel syndrome, injured nerve root, and muscle degeneration.

### C. Feature Extraction

Features extraction is a significant way of collecting useful information contained in the surface EMG signals and eliminating unnecessary sections and interferences. The features of EMG signals are divided into three groups which are time-domain, frequency domain, and time-frequency domain [10]. The advantages and disadvantages of the features are shown in in Table I.

Time-frequency analysis is evaluated in time and frequency domain as shown in Fig. 1. Features taken from time-frequency distribution (TFD) should be reduced before being sent to the classifier. To improve the accuracy of the classification, time-frequency distribution feature is proposed to overcome the limitation of TD features [12].

The function of TFD feature is to identify time-varying system properties from the non-stationary system. The TFD has a major problem which is high dimensionality and high resolution of features vectors and to overcome the problem is to reduce the dimensionality of the data [12]. Furthermore, mathematical functions described in the time domain and the frequency domain are commonly used as dimensionality reduction methods for TD features [13]. There are two techniques for dimensionality reduction which are feature projection and feature selection. Features projection techniques attempt to determine the best combination of the original features and create a new feature set that is generally smaller than the original one. For feature selection, it needs to consider a features vector for numerous specific EMG signal classification [12].

Time-frequency also used for EMG signal processing by past research. During the test, the surface myoelectric signal is compressed towards the lower frequency and the frequency of the signal is continuous changes over time. This can classify the surface myoelectric as slow and fast. For slow nonstationary is because of the electrical manifestations and affects the accumulation of metabolites. Next, fast nonstationary is related to the biomechanics of the task. The modification of the frequency content of the signal is affected by the variations in muscle force [14].

Besides that the mathematical techniques have been advanced to solve the problem in signal processing where there are combination methods of state space and statistical decision theory. These techniques happen to a broad class of nonlinear problems and focus on the presence of additive noise due to the problem of signal processing is nonlinear [15].

### D. Classification

Classification is from the extracted information of the EMG signals to map different patterns and match them appropriately. The classifier is to divide different categories of the features extracted and going to practice being control commands for the controller in the next stage [12]. There are



many techniques to classify EMG data and have their advantages and disadvantages as shown in Table II.

The problem for biomedical is when there are a few applications that can analyse the demand of patients. So, the application of machine learning can solve this problem with the detection and classification of the neuromuscular disorder based on EMG signal processing. From this application, the patients can skip techniques of ultrasound or MRI to diagnose the neuromuscular disorder. Many biomedical used the support vector machine (SVM) in signal classification applications as a machine learning method. The SVM can improve the accuracy of EMG signal classification and classify it into normal, neurogenic, or myopathic. The classification for SVM is applied based on the trained model after generated the training data in the training process [16].

TABLE I. FEATURE EXTRACTION DOMAIN

| Features Domain       | Advantages  | Disadvantages  |
|-----------------------|---|--|
| Time domain           | Low noise environments<br>Lower computational complexity [10]                       | Non-stationary property of EMG signal<br>Changing in statistical properties over time [10] |
| Frequency domain      | Reducing interference<br>Good localization of the signal.<br>Very clean signal [11] | High noise environment [11]  |
| Time-frequency domain | Can overcome the limitation of time-domain features [12]                            | High dimensionality<br>High resolution of feature vectors [12]                             |

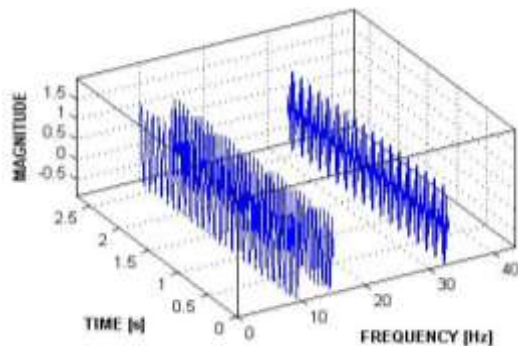


Fig. 1. Graph of Time-Frequency Domain.

TABLE II. CLASSIFICATION TECHNIQUES

| Techniques of classification     | Advantages   | Categorize                              |
|----------------------------------|--|---|
| Support vector machine (SVM)     | Works relatively well when there is a clear margin of separation between classes | Multiple motions                        |
| K-nearest neighbor (KNN)         | Simple implementation  | Hand motion                             |
| Multilayer perceptron (MLP)      | Capable of prescribing nonlinear class boundaries                                | Hand motion and forearm motion          |
| Artificial neural networks (ANN) | Suitable for modeling nonlinear data due can cover the distinctions              | Hand motion (left, right, up, and down) |
| Fuzzy logic (FL)                 | Control techniques in biosignal processing                                       | Biosignal characteristics               |

Recently, k-nearest neighbor (KNN) is a common machine learning tool due to its speed of processing and simplicity in the process of recognition. The concept of KNN is quite simple. The KNN algorithm creates a set of k data points in training data and forecasts test data dependent on the nearest neighbor. However, the significance of k must be carefully chosen because it has a direct effect on the efficiency of the classification. Specifically, the k-value depends mostly on the specification of the data set and model. However, the KNN algorithm is fast, easy, and effective [17].

### III. MATERIALS AND METHODS

#### A. EMG Data

The data for this project is obtained from open source by past research. The data includes EMG signal of finger movement from ten subjects which are six males and four females aged between 20 and 35 years old. The characteristic of the subject is normally limbed with no neurological or muscular disorders. EMG data are taken from this subject by using EMG channels. To firmly stick the sensor to the skin, two of the slot adhesive skin interface (DELSYS DE 2X SERIES EMG SENSOR) was applied to each of the sensors. There are ten classes of finger movements were including the movements for individuals and combined. The duration of every movement is in 5s with a resting period in 3 to 5s between each movement. The positions of the first electrode are adhesive skin interface stick to the skin and second electrode on the wrist. Positions of the electrode are shown in Fig. 2. Fig. 3 shows the data acquisition set up for the EMG data recording.

#### B. Signal Pre-Processing

In the EMG analysis, signal pre-processing which includes eliminating the offset signal, signal segmentation, and detection of onset are required. The signal offset was estimated using the baseline signal means. The offset in each channel was subtracted from the signal to remove the unwanted signal. Segments are calculated to produce signals indicating muscle activation before the features were removed and movement patterns were observed. The auto-segmentation as proposed by [18] was used to segment the EMG signal, thus, helping to reduce the computational complexity of the feature extraction. The magnitude and frequency of muscle activation segments varied from those of muscle activation segments.

The data EMG signals have been filtered in this stage. It is to improve the accuracy of data EMG signals invalidation. Besides, the filtering process is needed to overcome the noise in raw data signals and reduce the artifacts by using the various method. In this stage, the bandpass filter methods have been used to filter the data EMG signals with a range between 20 and 450 Hz. The filtering has been done for separate signals for each movement of EMG signals and to pass the only certain range of frequency and sampling rate is 4000 Hz. The unnecessary noise will be discarded following the range of bandpass filters. However, the noise in signals is very difficult to remove. This stage will give a big impact on classification because the accuracy of signals is depending on the data EMG signals.

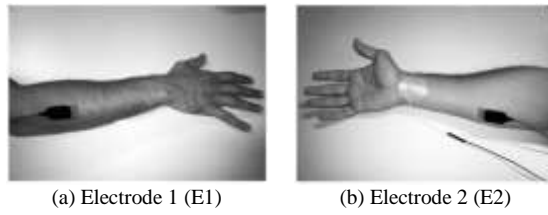


Fig. 2. Electrode Position.



Fig. 3. Data Acquisition Set Up.

### C. Time-Frequency Distribution

In the field of signal processing, feature extraction plays a critical role throughout achieving better identification quality for the detection of movement patterns. Various signal processing methods are applied in raw EMG to achieve an accurate and actual EMG signal. This process involves the conversions of raw EMG signals into a feature vector. Besides that, the characteristics of the EMG signal analysis can be classified into three groups, which includes the time-domain feature, frequency-domain feature, and time-frequency domain feature. The amplitude of the signals depends on the type and state of the muscle during the analysis phase. Most of the research focuses on time-domain to keep the computational complexity low and this feature does not require additional signal transformation. Various signal processing techniques are used on raw EMG to produce a reliable EMG signal.

In this project, the time-frequency distribution (TFD) which is spectrogram is used in feature extraction. The spectrogram is a fundamental component of TFD in the analysis of signals, particularly for noise and artifact reduction. The spectrogram is used to overcome the limitation of time and frequency representation for the non-stationary EMG signal. It is defined as the squared magnitude of STFT as expressed in (1).

$$S(t, f) = \left| \int_{-\infty}^{\infty} x(\tau) \omega(\tau - t) e^{-j2\pi f \tau} d\tau \right| \quad (1)$$

where  $S(t, f)$  is the time-frequency representation,  $x(\tau)$  is the EMG signal, and  $w(t)$  is the observation window.

TFD is preferred to obtain time and frequency information simultaneously. The spectrogram reveals the non-stationary existence of EMG signals in the time-frequency analysis. In TFD, the time and frequency resolution can be adjusted to obtain valuable signal details.

The parameter of the EMG signal was then estimated from the resulted time-frequency representation of the spectrogram.

The root mean square voltage ( $V_{rms}$ ) was measured instantaneously over time and the average values were taken for hand movement prediction. The average RMS voltage can be expressed as.

$$V_{rms(avg)} = \frac{1}{T} \int_0^T V_{rms}(t) dt \quad (2)$$

where

$$V_{rms}(t) = \sqrt{\int_0^{f_{max}} S_x(t, f) df} \quad (3)$$

where  $V_{rms}(t)$  is the instantaneous RMS voltage,  $S_x(t, f)$  is the time-frequency representation, and  $f_{max}$  is the maximum frequency of interest.

### D. Machine Learning

The information derived from the EMG signals will then be fed into the classifier to identify the different patterns and match them properly. Classifiers should be used to distinguish between different classes of features extracted. The obtained classifications will then be used as control commands for the controller in the next stage. Multiple methods are used to identify EMG information such as artificial neural networks (ANN), Bayesian classifier (BC), fuzzy logic (FL), multilayer perceptron (MLP), support vector machines (SVM), linear discriminant analysis (LDA), hidden Markov models (HMM) and K-nearest neighbor (KNN). Recently, several researchers have shown interest in effective ways to identify the origins of EMG signals.

The machine learning algorithm selected to determine the characteristics of the separation of the 10-finger movement in the EMG data signal after the signal processing phase. In this stage, the data EMG signals of  $V_{rms}$  have been separated according to each movement and electrode by each subject. The total  $V_{rms}$  data are 200 of 10 subjects with 10 movements for 2 electrodes. The total  $V_{rms}$  data of EMG signals will be separated into training and testing sets to evaluate the performance of data EMG signals. There are 80% of data for training and the other 20% of data for testing. The training test is to train the machine before getting an accurate value for testing. Next, the data have been imported into apps classification learners to analyse the accuracy of classifiers with train the data. There have various types of classifier in classification learner and the wide classifier has been used is KNN and SVM. Therefore, the best selection of classifiers is depending on the percentage of accuracy classifier.

## IV. RESULTS AND DISCUSSIONS

The EMG signal data consist of ten classes of individual and combined fingers movement. Every subject completed six times of test for 10-finger movements and resting time between the tests is around 3 to 5 seconds. Four of the six-time test is training and two of that is testing. The raw data is obtained from 10 individual movements which are 5 tests for individual movement and 5 tests for combined movement with two electrodes. This part shows the results and discussions for all methods that have been used.

Results of the EMG signal was obtained from individual finger movement and combine finger movement. The EMG

signal data was run using the MATLAB software and the graph of the EMG signal is a voltage (V) versus time (s) and it is shown as the amplitude of EMG signal during the test.

There are different finger movements for individuals and combine finger movements. The individual finger movements are consisting of thumb (T), index (I), middle (M), ring (R), little (L). Then, for combined finger movements are consists of hand close (HC), thumb index (T-I), thumb little (T-L), thumb middle (T-M), and thumb ring (T-R). The different movements of the finger are shown as Fig. 4.

These results indicate that using two channels of electrodes during collected data for individual finger movement. It also shows the signal for both electrodes for each movement. From the signals, the electrode 1 is more informative than the electrode 2 for thumb finger movement. This is due to the location of the second electrode where the electrode mounted as shown in Fig. 2. The located of the second electrode is on the low contraction muscle during finger movement because of that the signal of the second electrode not more informative. The informative signal depends on the contraction or muscle movement during the test of the finger movements. The signal of EMG data during the test of thumb movement for the individual finger movement as shown in Fig. 5.

These results show the signal of combine finger movements and indicate two channels of electrodes during the test of finger movement. These signals are made up of two electrodes used during the test of combined finger movements. For 0.7 s of the signal shows the second electrode gives more information than the first electrode in hand close movement. This is due to the location of electrodes mounted during the test. However, after 0.7 s the signal of both electrodes shows the same or constants informative. In detail, contraction or muscle movement is higher on the location of the second electrode for 0.7 s, and the muscle movement almost the same after that. The EMG signal for hand close in combined finger movements as shown in Fig. 6.



Fig. 4. The different Finger Movements.

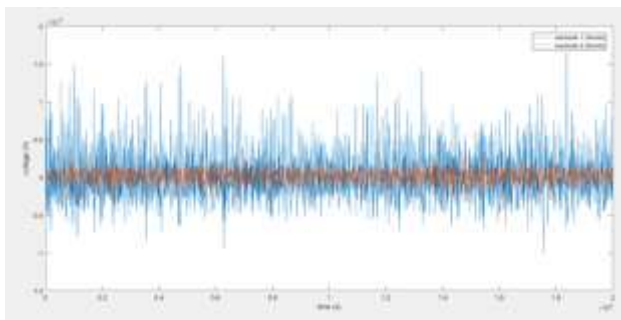


Fig. 5. Raw EMG Signal for Thumb Finger Movement.

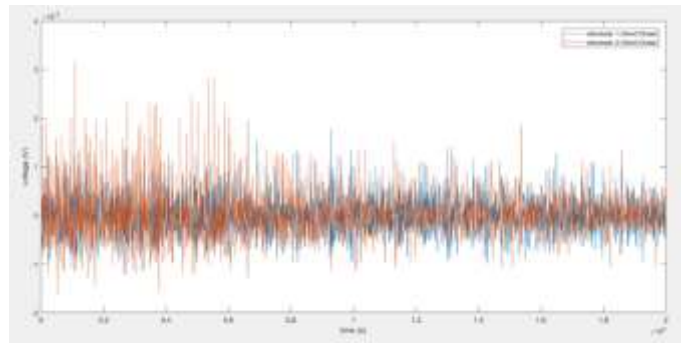


Fig. 6. Raw EMG Signal for Hand Close Finger Movement.

In this part, the signals of EMG are based on the information given on the contraction or muscle movement during the test of finger movement for individual and combined movements. The location of the electrodes is important to collect the data during the test. This EMG signal needs to filter to reduce unnecessary noise before the signal processing stage. The filtering process has been done for each movement and separated by electrodes. The bandpass filter has been used with range 20 Hz to 450 Hz and the sampling rate is 4000 Hz. The signal only passes by the range of frequency and the other will discard. The graph shows the EMG signal before and after the filtering process for the middle finger movement for electrode 1. The filtered middle finger movement as shown in Fig. 7.

In signal processing, the TFD is selected as a fundamental component to analyse the EMG signal, especially for noise. From this, the time and frequency can be measured for RMS. In TFD, to obtain valuable signal details in the EMG signal it can adjust the time and frequency resolution. The extraction of instantaneous RMS voltage EMG signal for electrode 1 of hand close finger movement is shown in Fig. 8.

The average RMS voltage of the EMG signals for individual finger movement is shown in Fig. 9. In the figure, subject 4 shows the highest level for electrode 1 and electrode 2 compare to the other subject in the index finger movement. Thus, during the test of individual finger movement, the subject 4 get more information in signal due to contraction or muscle movement. Subject 4 gives the best signal to control or classify the EMG signal for individual finger movement.

The average RMS voltage of the EMG signals for combine finger movement is shown in Fig. 10. Based on the figure, subject 4 have the highest level of electrode 1 and electrode 2 for average combined finger movement compared to the other subject. From this, the contraction or muscle movement at the location of electrode mounted for subject 4 is higher during testing. The EMG signal for subject 4 has more information about the muscle movement for electrode 1 and electrode 2 in hand close finger movement.

Fig. 11 shows the average  $V_{rms}$  signals for each electrode from all subjects. The average of data EMG signals is to investigate the various levels of acceptance of the suggested. The figure shows the level of information in EMG signals based on the finger movement. The higher level of the graph means the more informative the EMG signals.

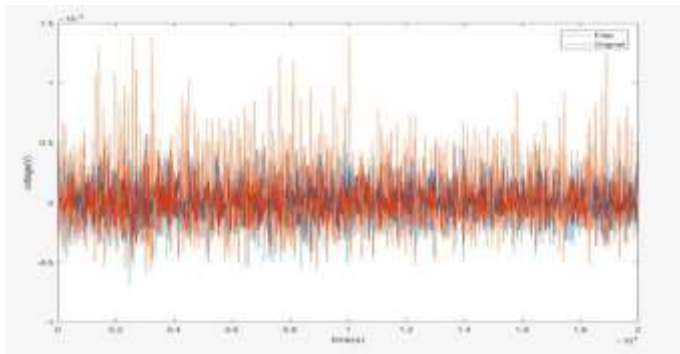


Fig. 7. Raw and Filtered EMG Signal.

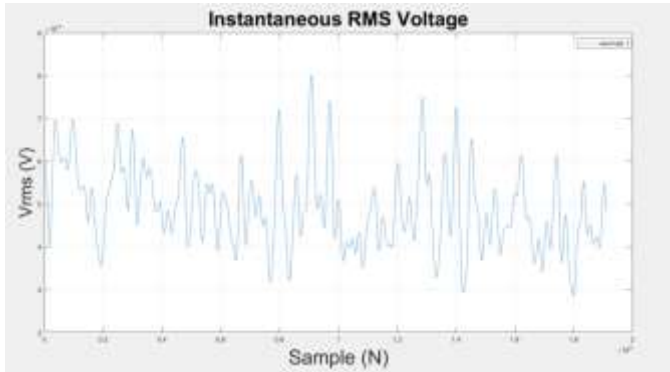


Fig. 8. Instantaneous RMS Voltage for Hand Close Finger Movement.

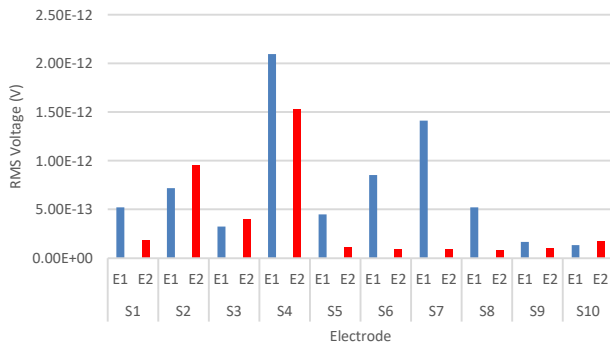


Fig. 9. Average RMS Voltage for Individual Finger Movement (Index).

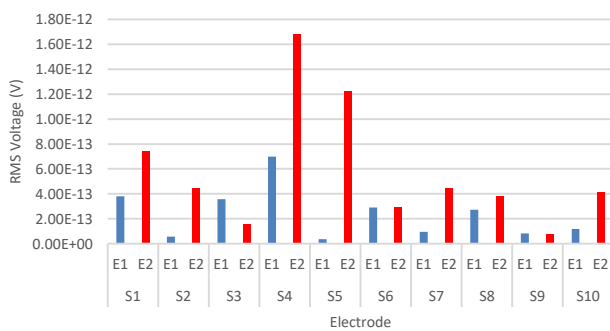


Fig. 10. Average RMS Voltage for Combined Finger Movement (Hand Close).

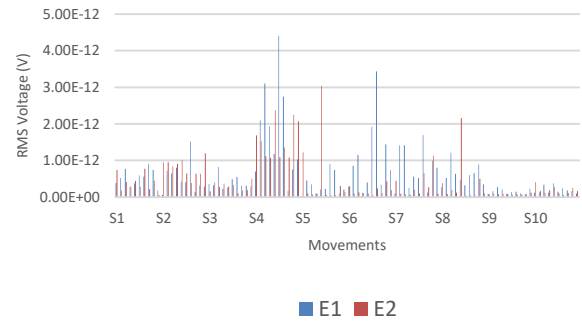


Fig. 11. The Average  $V_{rms}$  EMG Signals for Each Electrode from All Subjects.

The final step in this project is to determine a suitable classifier for data EMG signal from different classes of finger movements. After getting an average of  $V_{rms}$  data for EMG signal, techniques classifiers are used in signal classification to analyse the best machine learning for data EMG signal. The techniques are typically used to avoid confusing the prosthetic controller with different classification decisions and to increase the efficiency of the classifier by avoiding unnecessary classification errors. There is the comparison of techniques classifier for 100% data EMG signal as shown in Table III.

The total  $V_{rms}$  data EMG signal will divide into two-part which is 80% for training and the other 20% for testing. The training data is to train the machine learning before taking the results for testing finger movements. There is the comparison data for 80% and 20% of  $V_{rms}$  EMG signal as shown in Table IV.

From the comparison results, the best classifier for  $V_{rms}$  data EMG signal is the k-nearest neighbor (KNN). This is because the percentage of the accuracy of the KNN classifier is 100% for training and testing which is more accurate from the other classifier. The accuracy of EMG classification is determined based on the percentage in classification learner and can be a plot by a scatter plot and confusion matrix. A scatter plot or scatter graph is displaying the values of two variables from a set of data and identify the type of relationship between variables. The scatter plot for 80% of  $V_{rms}$  data EMG signal as shown in Fig. 12 and the scatter plot for 20% of  $V_{rms}$  data EMG signal as shown in Fig. 13.

Next, the confusion matrix or table of confusion is showing the error matrix for data with predicted class and actual class. The confusion matrix can plot by true positive rates and false-negative rates. The confusion matrix for 80% of  $V_{rms}$  data EMG signal as shown in Fig. 14 and the confusion matrix for 20% of  $V_{rms}$  data EMG signal as shown in Fig. 15.

TABLE III. THE COMPARISON FOR 100% DATA EMG SIGNAL

| Type of classifier                  | Percentage (%) |
|-------------------------------------|----------------|
| SVM (fine gaussian SVM)             | 64 %           |
| KNN (fine KNN)                      | 100 %          |
| ENSEMBLE CLASSIFIER (boosted trees) | 73 %           |

TABLE IV. THE COMPARISON OF 80% (TRAINING) AND 20% (TESTING) OF THE DATA EMG SIGNAL

| Classifier          | 80% of data | 20% of data |
|---------------------|-------------|-------------|
| SVM                 | 75%         | 90%         |
| KNN                 | 100%        | 100%        |
| Ensemble Classifier | 73%         | 100%        |

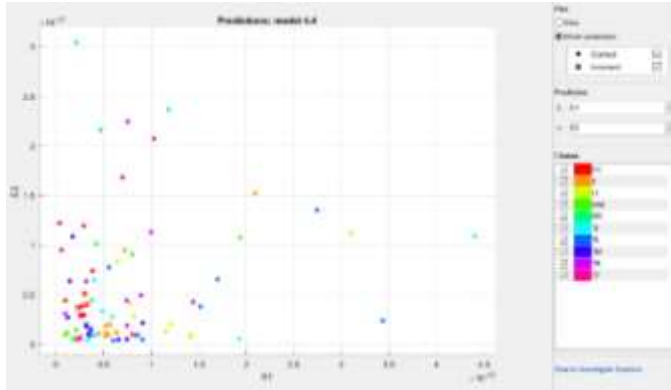


Fig. 12. The Scatter Plot for 80% of  $V_{rms}$  Data EMG Signal (Training).

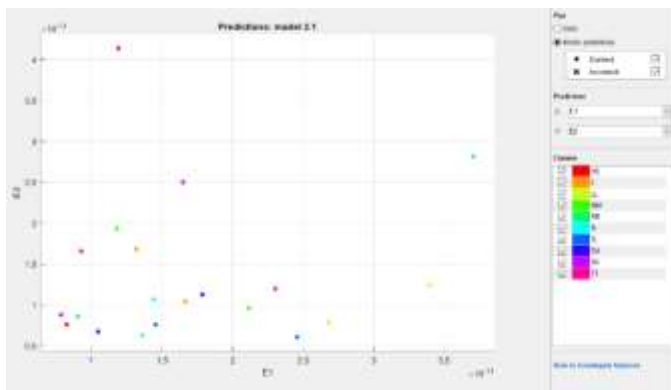


Fig. 13. The Scatter Plot for 20% of  $V_{rms}$  Data EMG Signal (Testing).

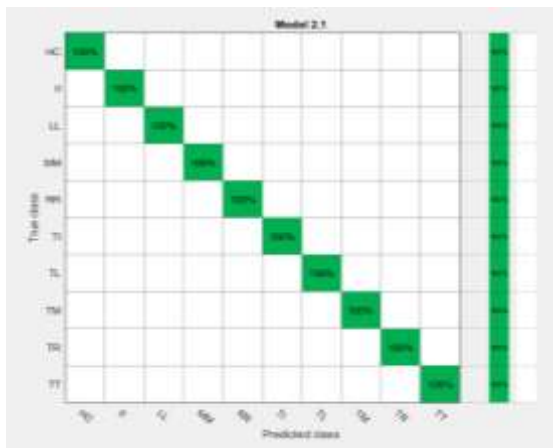


Fig. 14. The Confusion Matrix for 80% of  $V_{rms}$  Data EMG Signal (Testing).

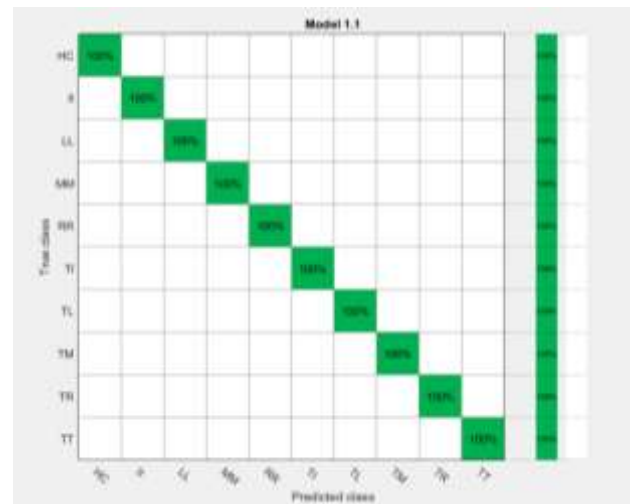


Fig. 15. The Confusion Matrix for 20% of  $V_{rms}$  Data EMG Signal (Testing).

As a result, significant increases in classification precision have been made. In the signal classification, three classifiers have been used to train the training and testing data of the  $V_{rms}$  EMG signal. The most accurate for the  $V_{rms}$  EMG signal is the KNN classifier with 100% accuracy for training and testing. This is because the KNN classifier is easy to implement.

## V. CONCLUSIONS

As a conclusion, the research to analyze the surface EMG signals ( $V_{rms}$ ) in finger movement by using the TFD have been presented. This analysis covered the analysis data of an individual and combined finger movements of EMG signals for prosthetic hand control. The EMG signals have been filtered by using the bandpass filter to overcome the unnecessary noise in signals with a range from 20 to 450 Hz. TFD is then used for the feature extraction to get the average of signal to classification.

In addition, this research classify the accuracy of individual and combine finger movement based on surface EMG signals towards improved prosthetic control. The three classifiers have been used to train all data EMG signal and the most accurate classifier have been chosen as machine learning to conduct the EMG signals. The data of EMG signals have been trained followed by each movement and each subject.

Finally, the performance of the KNN classifier has been compared with other classifiers. The data  $V_{rms}$  EMG signals have divided into two parts which are training and testing. For instance, the training data is to train the machine learning to get accurate data for testing. The EMG datasets are belong to 10 different classes for individual and combined movements collected from 10 subjects by using two channels of electrodes and the accuracy of classifier in the range 64% to 100% with various types of classifiers of the data  $V_{rms}$  EMG signals.

## VI. FUTURE WORKS

For future works, further studies about other finger movements for individual and combined finger movements is essential. It is to get data for other movements towards the prosthetic hand. Next, the EMG signal must be tested with prosthetic hand to make sure the classifier that has been chosen is suitable and can be integrated with the prosthetic hand. This is to ensure the accuracy of the classifier is accurate even after intergrating it with the prosthetic hand.

## ACKNOWLEDGMENT

This project is fully funded by Universiti Teknikal Malaysia Melaka under Short Term Grant Scheme (High Impact) no. PJP/2020/FKE/HI19/S01717.

## REFERENCES

- [1] S. Pancholi and A. M. Joshi, "Portable EMG Data Acquisition Module for Upper Limb Prosthesis Application," in IEEE Sensors Journal, vol. 18, no. 8, pp. 3436-3443, 15 April 2018.
- [2] S. Rezazadeh, D. Quintero, N. Divekar and R. D. Gregg, "A Phase Variable Approach to Volitional Control of Powered Knee-Ankle Prostheses," 2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), Madrid, 2018, pp. 2292-2299.
- [3] V. K. Mishra, V. Bajaj, A. Kumar and G. K. Singh, "Analysis of ALS and normal EMG signals based on empirical mode decomposition," in IET Science, Measurement & Technology, vol. 10, no. 8, pp. 963-971, 11 2016.
- [4] E. F. Shair, S. A. Ahmad, M. H. Marhaban, A. R. Abdullah and S. B. M. Tamrin, "Implementation of Spectrogram for an Improved EMG-based Functional Capacity Evaluation's Core-Lifting Task," 2018 IEEE-EMBS Conference on Biomedical Engineering and Sciences (IECBES), Sarawak, Malaysia, 2018, pp. 13-17.
- [5] R. D. Wilson, S. J. Page, M. Delahant, J. S. Knutson, D. D. Gunzler, L. R. Sheffler and J. Chae, "Upper-Limb Recovery After Stroke: A Randomized Controlled Trial Comparing EMG-Triggered, Cyclic, and Sensory Electrical Stimulation," Neurorehabilitation and Neural Repair, vol. 30, no. 10, 2016, pp. 978-987.
- [6] P. Beckerle, S. Willwacher, M. Liarokapis, M. P. Bowers, and M. B. Popovic, "9 - Prosthetic Limbs," M. B. B. T.-B. Popovic, Ed. Academic Press, 2019, pp. 235-278.
- [7] J. L. Segil, S. A. Huddle and R. F. f. Weir, "Functional Assessment of a Myoelectric Postural Controller and Multi-Functional Prosthetic Hand by Persons With Trans-Radial Limb Loss," in IEEE Transactions on Neural Systems and Rehabilitation Engineering, vol. 25, no. 6, pp. 618-627, 2017.
- [8] P. Konrad, The ABC of EMG - Noraxon, vol. 1, no. March. 2006.
- [9] S. Pancholi and A. M. Joshi, "Portable EMG Data Acquisition Module for Upper Limb Prosthesis Application," IEEE Sens. J., vol. 18, no. 8, pp. 3436-3443, 2018.
- [10] D. Zhou, Y. Fang, J. Botzheim, N. Kubota and H. Liu, "Bacterial memetic algorithm based feature selection for surface EMG based hand motion recognition in long-term use," 2016 IEEE Symposium Series on Computational Intelligence (SSCI), Athens, 2016, pp. 1-7.
- [11] L. Wu, X. Zhang, X. Chen and X. Chen, "Visualized Evidences for Detecting Novelty in Myoelectric Pattern Recognition using 3D Convolutional Neural Networks\*," 2019 41st Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC), Berlin, Germany, 2019, pp. 2641-2644.
- [12] N. Nazmi, M. A. A. Rahman, S. I. Yamamoto, S. A. Ahmad, H. Zamzuri, and S. A. Mazlan, "A review of classification techniques of EMG signals during isotonic and isometric contractions," Sensors, vol. 16, no. 8, pp. 1-28, 2016.
- [13] S. Inoue, M. Oya and H. Ohta, "Finger Joint Dynamics with Myoelectric Signal inputs," 2018 International Conference on Information and Communication Technology Robotics (ICT-ROBOT), Busan, 2018, pp. 1-4.
- [14] N. Jose, R. Raj, P. K. Adithya and K. S. Sivanadan, "Classification of forearm movements from sEMG time domain features using machine learning algorithms," TENCON 2017 - 2017 IEEE Region 10 Conference, Penang, 2017, pp. 1624-1628.
- [15] A. Furui, H. Hayashi, Y. Kurita and T. Tsuji, "Variance distribution analysis of surface EMG signals based on marginal maximum likelihood estimation," 2017 39th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC), Seogwipo, 2017, pp. 2514-2517.
- [16] M. Kilany, E. H. Houssein, A. E. Hassanien and A. Badr, "Hybrid water wave optimization and support vector machine to improve EMG signal classification for neurogenic disorders," 2017 12th International Conference on Computer Engineering and Systems (ICCES), Cairo, 2017, pp. 686-691.
- [17] J. Too, A. R. Abdullah, N. M. Saad, N. M. Ali, and T. N. S. T. Zawawi, "Application of Gabor Transform in the Classification of Myoelectric Signal," vol. 17, no. 2, pp. 873-881, 2019.
- [18] E. F. Shair, A. R. Abdullah, T. N. S. Tengku Zawawi, S. A. Ahmad and S. Mohamad Saleh, "Auto-Segmentation Analysis of EMG Signal for Lifting Muscle Contraction Activities", J. Telecommun. Electron. Comput. Eng., vol. 8, no. 7, pp. 17-22, 2016.

# Autism Spectrum Disorder Diagnosis using Optimal Machine Learning Methods

Maitha Rashid Alteneiji<sup>1</sup>, Layla Mohammed Alqaydi<sup>2</sup>  
Abu Dhabi School of Management  
Abu Dhabi, UAE

Muhammad Usman Tariq<sup>3</sup>  
Assistant Professor, Abu Dhabi School of Management  
Abu Dhabi, UAE

**Abstract**—Autism spectrum disorder (ASD) is the disorder of communication and behavior that affects children and adults. It can be diagnosed at any stage of life. Most importantly, the first two years of life, regardless of ethnicity, race, or economic groups. There are different variations of ASD according to the severity and type of symptoms experienced by people. It is a lifelong disorder, but treatment and services can improve the symptoms. The literature focuses on one of the main methods used by physicians to diagnose ASD. Many types of research and medical reports have been reviewed; however, a few of them only give good medical results for the strong differentiation of ASD from healthy people. This paper focuses on using machine learning algorithms to predict an individual with specific ASD symptoms. The target is to predict an individual with specific ASD symptoms and finding the best machine learning model for diagnosis. Further, the paper aims to make the autism diagnosis faster to deliver the required treatment at an early stage of child development.

**Keywords**—Autism diagnosis; autism disorder; autism detection; machine learning; ASD

## I. INTRODUCTION

Artificial Intelligence has increasing importance in society. The main aim of the paper is to use artificial intelligence and machine-learning models that can help in medical fields by finding the optimal model that can recognize individuals with specific Autistic Spectrum Disorder symptoms [8]. The attention to the Artificial Intelligence field has been grown in the last few years [9]. This interest is not only motivated by the trend of designing models with human thoughts or behaviors, but also for the way of their use in real life [2,6]. The development of such models represents an ambitious and competitive task among many scientists and programmers [3]. The ability to harness artificial intelligence in medical matters has been an important topic since the beginning of this century [1]. Since then, devices and ideas have developed to facilitate the prediction and detection of specific diseases by refining machine education. Based on the importance of artificial intelligence, this paper aims to identify and implement the machine learning process that produces an algorithm capable of detecting whether a person has Autism Spectrum Disorder. Since early intervention affords the best opportunity to support healthy development and deliver benefits across the lifespan, this paper helps parents to recognize if their child has an ASD at an early age [4,7]. It also has a value in the health sector since there is no valid health reason for this disease. Based on autism speaks organization in USA, ASD diagnoses by applying behavioral exams or questionnaires, which require a

lot of time and effort provided from parents and clinicians [12]. Therefore, this work shows the power of machine learning algorithms in detecting individuals with specific Autistic Spectrum Disorder symptoms [5].

Further, this paper aims toward making the diagnosis of autism a faster process that enables delivery of therapy at earlier and more impactful stages of child development using machine-learning algorithms. The introduction is the first section of this paper. It mainly gives a general overview of the autism spectrum disorder. It also represents the target group and the benefits to the community. The second section is the theoretical background concerning this paper. It focuses on Autism Spectrum Disorder diagnosing methods and comparing these methods to choose the optimal diagnosing method database. It also gives an overview of machine learning approaches and the evaluation of machine learning. The next section is the methodology, which fits the knowledge discovery process to the available ASD datasets. The fourth section is the results, which represent in detail the individual results for the entire applied machine learning models, in terms of the performance, and a comparison between them to find the optimal machine learning model. The last sections, which are a discussion and conclusion, summarize the work and a discussion about the research questions. Also, it gives an outlook for possible future papers and improvements on this topic.

## II. THEORETICAL BACKGROUND

Various autism rating scales have been developed over the past 30 years to diagnose the autism spectrum disorder (ASD) at an earlier stage [1,10]. There are various tools and instruments developed by psychologists and neuroscientists to diagnose it at earlier stages [11]. Most of the tools focus on the diagnosis of users through screening methods [3].

### A. Autism Behavior Checklist (ABC)

One of the methods is the autism behavior checklist (ABC) based on identifying ASD in children at early stages [4]. The benefit of the checklist utility is to evaluate the current autistic symptoms of the user with the help of parents in different situations and conditions. It provides a set of questions to evaluate the in-depth condition of the user [6]. The method combines different scales, such as language, object recognition, body, sensory, social, and daily use skills [10]. The item scores are mostly from 1-5 based on the impairment degree [2,5]. The ABC has been in use for more than 30 years as a rapid tool for diagnosing autism in early stages [7,13].

However, there is no general agreement to use the same values as a standard method [1,3, 14].

### B. Child Behavior Checklist (CBCL)

The Child Behavior Checklist (CBCL) is one of the oldest screening tools and most widely used standardized measures in child psychology for evaluating unusual behavioral and emotional problems [15]. The CBCL questionnaire focuses on internalizing and externalizing behaviors such as anxiety, over-control, aggression, and hyperactivity. There are two versions of The Child Behavior Checklist: a preschool version (aged 2 to 5) and a school-age version (aged 6 to 18). The preschool version questionnaire carries 100 questions with three scale responses ranging from 0-2, where 0 represents 'Not True,' and 2 indicates 'Very True.' The school-age version questionnaire carries 118 questions with the same rating scale responses [4, 9, 16].

### C. Social Communication Questionnaire (SCQ)

The Social Communication Questionnaire (SCQ) is an ASD-screening tool used for age four and above. The SCQ consists of 40-items based on parent-report screening measures after a semistructured parent interview with a trained clinician or researcher that can be used for diagnostic ASD symptomatology. There are two different versions of the SCQ. The SCQ 'Current' asks a respondent to indicate whether behaviors have been present during the past three months. The other version is the SCQ 'Lifetime' references complete developmental history and asks respondents to indicate whether behaviors have ever been present [11, 17].

### D. Autism Spectrum Quotient (AQ)

The Autism-Spectrum Quotient (AQ) is among the most widely used scales assessing autistic traits in the general population. The AQ is a self-administered questionnaire for measuring how adults with normal intelligence show autistic traits. It consists of 50 questions, with ten questions assessing five different domains relevant for autistic traits: social skill, attention switching, attention to detail, communication, and imagination. People with a clinical diagnosis tend to score above 32 out of 50 on the AQ [10].

### E. Short Autism Quantitative (AQ-10)

In 2012, Allison et al. create a new version of Autism-Spectrum Quotient (AQ) for adults consisted of 10-items only to make it simpler and more timesaving. AQ-10 has a predictive power similar to the origin AQ version [5,18]. Later on, Allison created shorter versions for adolescents and children with ten items too. Score calculations for the adolescent and child short versions are different from the AQ adult short version. The diagnosis depends on the final questionnaire score in behaving with some genetic information [10, 20].

### F. Comparison of ASD Diagnosing Methods

The paper evaluated different ASD diagnosing methods. Table I provides a comparison between the ASD screening methods discussed in the first section of this section. As noted, most ASD screening tools focus on infants, toddlers, and children. Almost all diagnosing methods used questionnaires to diagnose ASD behaviors [19]. CBCL has the maximum

number of items in the questionnaire with 118 items, while the AQ-10 has the minimum number of questionnaire items with only ten items. Screening is valid if it detects most cases with the target disorder, which gives a high sensitivity rate and excludes most cases without the disorder, which gives a high specificity rate [21, 30].

TABLE I. COMPARISON OF ASD DIAGNOSING METHODS

| Diagnosing method | # Q | Target                 | Specificity | Sensitivity |
|-------------------|-----|------------------------|-------------|-------------|
| CBCL              | 118 | Children & Adolescents | 82%         | 75%         |
| ASSQ              | 27  | Children & Adolescents | 86%         | 91%         |
| ABC               | 57  | Children               | 91%         | 77%         |
| AQ                | 50  | Children               | 95%         | 95%         |
| AQ                | 50  | Adolescents            | NA          | NA          |
| AQ                | 50  | Adult                  | 52%         | 93%         |
| AQ-10             | 10  | Children               | 74%         | 77%         |
| AQ-10             | 10  | Adolescents            | NA          | NA          |
| AQ-10             | 10  | Adult                  | NA          | NA          |
| SCQ               | 40  | Children & Adolescents | 93-100%     | 58-62%      |

In ASD screening methods, sensitivity refers to the true positive rate, which is the ability of the screening tool to identify a person with autism [29]. Specificity refers to the true negative rate, which is the power of the screening tool to identify a person who is control of autism. In terms of validity, almost all the screening methods have acceptable sensitivity rates, ranging from 70%- 100% and specificity between 80% and 100% [6, 22]. As shown in Table I, AQ screening is the most efficient method with only ten questions, which require less time to complete than other methods. Further, AQ deals with many age segments, and each of them has a specific questionnaire, which will be discussed in detail in the methodology section [15, 23].

### G. Knowledge Discovery in Database

Knowledge Discovery in Databases, KDD, is an exploratory analysis and modeling of big data. KDD is the organized process of identifying useful, valid, and meaningful patterns from big databases [29]. The core of the KDD process is data mining, which explores the unknown patterns of the algorithms to develop the models. The model uses for predicting new unknown instances [30,31]. There are nine iterative processes in KDD listed below:

1) *Understanding the application domain*: This process defines the goals of the end-user and the environment in which the KDD process will take place with a full understanding of what should do.

2) *Creating the data set*: This process checks the available data and then integrates it with additional obtained data into one data set for the knowledge discovery.

3) *Preprocessing*: The available data goes to the preprocessing step, which includes handling missing values



and removing noise or outliers or duplicated data to enhance the reliability of the data.

4) *Data transformation*: This step prepares and develops better data to create the best possible model. It includes customizing the data dimension by feature selection and record sampling.

5) *Choosing the appropriate data mining task*: The main goal of this process is to decide on which type of data Mining to use. Data mining types include clustering, classification, and regression, depending on the DM goal, either prediction or description.

6) *Choosing the data mining algorithm*: This step includes selecting the appropriate searching patterns to use. Each algorithm has parameters and tactics of machine learning.

7) *Applying the Data Mining algorithm*: To get a satisfying result in this process, it might require applying it several times.

8) *Evaluation*: The model with the found patterns is evaluated and interpreted concerning the goals mentioned in the first process. This step focuses on the usefulness and comprehensibility of the induced model.

9) *Using the discovered knowledge*: The final step is to try the knowledge into other systems for further action and make changes to the system and measure the effects.

#### H. Recent Machine Learning Research on ASD Screening and Diagnosis

With the fast growth in the big data field, Autism Spectrum Disorder research must benefit from this area as other searching fields. By looking at the available research about using machine learning in diagnosing ASD, it has been clear that there is a positive effect of using machine learning in diagnosing ASD with a database containing an ASD screening method with some genetic information. Machine learning can be sorted as unsupervised and supervised learning. ASD diagnosing is a supervised machine learning that has a dependent variable to be predicted, which is the result of the diagnosis [5, 24]. A successful supervised machine-learning model is the one that can predict the target correctly and generalize new instances predicts. Usually, model validation can be measured by accuracy, which has two subtypes, sensitivity and specificity. Another measurement for the model's accuracy is the area under the receiver operating character curve, AUC [25]. The AUC shows how well a method makes positive and negative categorical distinctions between sensitivity and specificity. According to Kayleigh's research, most of ASD diagnosing models used ADTree and SVM algorithms. ADTree is a classification machine-learning algorithm that consists of an alternation of decision nodes, which specify a predicate condition, and prediction nodes, which contain a single number. An ADTree classifies an instance by following all paths for which all decision nodes are true and summing any prediction nodes that are traversed. [12, 26] The objective of the support vector machine algorithm, SVM, is to find a hyperplane in N-dimensional space (N: the number of features) that distinctly classify the

data points. [2, 5, 27] data, data mining, uses complex mathematical machine learning algorithms [28].

### III. METHODOLOGY

This section starts by describing the paper's technical framework and then identifies the data-mining goal of this paper. The following parts include describing the data collection processes, the specification of a DM approach and algorithm, the optimal machine learning models to be used with supervised cases, and how it will be evaluated. The predefined goal was to develop a model that can recognize an individual with specific ASD symptoms using specific machine-learning methods. The next practical working steps will take place in this paper are:

- Find the right datasets for paper use and apply data mining steps to find the appropriate machine learning models that can expect ASD symptoms.
- Apply several machine-learning algorithms to critically review, evaluate, and compare one another to choose the optimized prediction model.
- Provide a machine-learning algorithm with optimal prediction quality for identifying an individual with specific ASD symptoms without access to the class label.
- Introduce a proper framework for ranking the quality of the diagnosing models.

A significant limitation is that this work focuses on evaluating some ASD characteristics that may affect the diagnostic result, while the simple reasons and factors for having an Autistic Spectrum Disorder are still not clear all over the world.

#### A. Technical KDD Framework

The proposed Knowledge Discovery in Database framework of ASD detection traits is shown in Fig. 1. The Knowledge Discovery process (KDD) process starts with understanding the autism spectrum disorder symptoms and factors, which are discussed through the second section. The next three processes are related to data collection, data preprocessing, and data transformation.

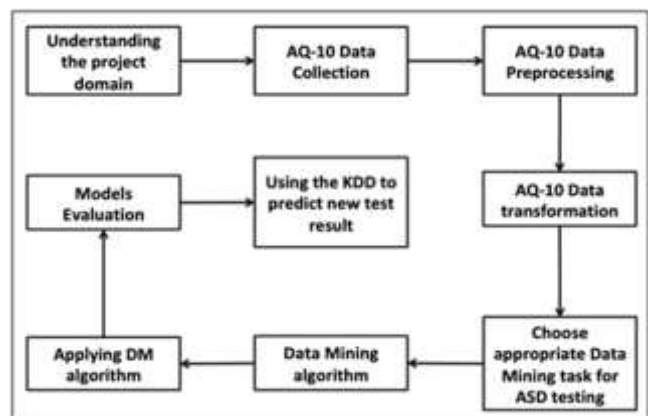


Fig. 1. KDD Process in Database.

### B. Data Collection

The data used in this paper is secondary data; the AQ-10 ASD data. Data was collected through a mobile application using the ASD AQ-10 diagnostic method, and healthcare professionals relied on the final diagnosis. The primary purpose of the app is to collect useful data about ASD cases. This data is stored in a MySQL database to understand the main features that may affect ASD diagnostics using data analysis. Also, the data was divided into three databases, since each age category has a specific diagnosing questionnaire.

### C. Data Preprocessing

The collected data usually cannot be used directly in performing the analysis process. Therefore, the raw data needs to be cleaned and performed in a usable format. Cleaning the data includes replacing or removing missing values and discretization for certain continuous variables such as the age of individuals; this step is called the data pre-processing. Before using the available data, all the redundant or unnecessary variables must be removed from the Database. The Null values, the columns that do not have unique values, must be removed too. This process will improve the model's prediction and raise the accuracy rate. The variables that have been removed in the ASD databases are as listed:

- Case number: the case number was added as a counter of all the cases that used the ASD screening.
- ASD Screening type: this variable was added to split the data into spirit databases based on age categories.
- Reasons for taking the screening: this variable contains texts, and it did not add any value to the data analysis. Also, it will negatively affect the prediction model results.
- Language: since the app is available for people worldwide, the diagnostic test is also available with the most common languages.
- User: this variable represents the person who answers the screening instead of the child.
- Used app before: the users were asked this question to avoid attribute duplication.

### D. Data Transformation

On the other hand, data transformation includes customizing the data dimension by feature selection according to the model needs [31]. These two operations are mandatory to achieve better performance and accuracy in the Machine Learning prediction models. Data mining processes are the next step in the KDD processes that consists of applying data analysis and discovery algorithms that produce a particular enumeration of models over the data. In order to find the optimal machine-learning algorithm that gives the best ASD result prediction and discover stricter rules, the data mining type for the user database must be specified. The rule phase will be discovered by a classification system that will be used to predict the value of the unseen cases. All the possible machine-learning algorithms will be evaluated using the confusion metrics to find the accuracy, sensitivity, and specificity. The last process is sharing the discovered

knowledge with the health professional so it can be used in a professional way to serve society and help parents to discover their kids' status from an earlier age.

### E. Data Description

In this paper, the used databases relating to specific age groups, which are infants, children, and adolescents. The datasets can be divided into ten behavioral questions for each age group, and several variables that influence the final evaluation of the condition are used in the diagnostic database. The influencing variables include age, gender, ethnicity, jaundice, and family history. Table II shows these variables with the data type and the description for each variable.

The next three Tables III, IV, and V show the ten variables details in the toddler, adolescent, and children screening methods. The three databases carry out ten symptoms to be answered with either yes or no. These ten questions are the most noticed symptoms in diagnosing individuals with ASD.

TABLE II. VARIABLES USED FROM ASD DATABASE FOR DIAGNOSIS

| Variable       | Data type   | Description   |
|----------------|-------------|---|
| A1 – A10       | Binary      | YES/NO  |
| Age            | Continuous  | Age of Individual   |
| Gender         | Binary      | Male/Female   |
| Ethnicity      | Categorical | List:<br>White<br>Middle Eastern<br>White European<br>Asian<br>Black<br>Latino<br>Mixed<br>Others |
| Jaundice       | Binary      | YES/NO  |
| Family history | Binary      | YES/NO  |
| Nationality    | Categorical | List (All the worlds' counties)   |
| Target class   | Binary      | YES/NO  |

TABLE III. AQ-10 CHILDREN SCREEN FEATURES

| Variable | Child screening features  |
|----------|---|
| A1       | S/he often notices small sounds when others do not  |
| A2       | S/he often notices small sounds when others do not  |
| A3       | In a social group, s/he can easily keep track of several different people's conversations             |
| A4       | S/he finds it easy to go back and forth between different activities                                  |
| A5       | S/he doesn't know how to keep a conversation going with his/her peers                                 |
| A6       | S/he is good at social chit-chat  |
| A7       | When s/he is read a story, s/he finds it difficult to work out the character's intentions or feelings |
| A8       | When s/he was in preschool, s/he used to enjoy playing pretending games with other children           |
| A9       | S/he finds it easy to work out what someone is thinking or feeling just by looking at their face      |
| A10      | S/he finds it hard to make new friends  |

TABLE IV. AQ-10 ADOLESCENT SCREENING FEATURES

| Variable | Adolescent screening features  |
|----------|--|
| A1       | S/he notices patterns in things all the time   |
| A2       | S/he usually concentrates more on the whole picture rather than the small details                |
| A3       | In a social group, s/he can easily keep track of several different people's conversations        |
| A4       | If there is an interruption, s/he can switch back to what s/he was doing very quickly            |
| A5       | S/he frequently finds that s/he doesn't know how to keep a conversation going                    |
| A6       | S/he is good at social chit-chat   |
| A7       | When s/he was younger, s/he used to enjoy playing games involving pretending with other children |
| A8       | S/he finds it difficult to imagine what it would be like to be someone else                      |
| A9       | S/he finds social situations easy  |
| A10      | S/he finds it hard to make new friends   |

TABLE V. AQ-10 TODDLER SCREENING FEATURES

| Variable | Toddler screening features   |
|----------|--|
| A1       | S/he often looks at you when you call his/her name   |
| A2       | S/he often can easily get eye contact with you   |
| A3       | S/he can easily point to indicate that s/he wants something                                |
| A4       | S/he can easily point to share interest with you   |
| A5       | S/he can easily pretend  |
| A6       | S/he is good at social chit-chat   |
| A7       | S/he can easily follow where you are looking   |
| A8       | When you or someone in the family is upset, s/he can show signs of wanting to comfort them |
| A9       | S/he first words was typical ones  |
| A10      | S/he finds it easy to use simple gestures  |

Based on the available database, the total number of attributes in this paper is equal to 1811 attributes after removing the NULL and redundant attributes. The attributes are divided into three spirit databases since each age group must deal with specific symptoms. The toddler age group carries more than half of the aggregate data, with around 70% diagnosed as having autism spectrum disorder. Child and Toddler age groups distribution divided by half for each gender, which also divided by half for the diagnosing result.

#### F. Specification of DM Approach and Algorithm

This part deals with the KDD-steps 5, 6 and 7, including choosing the appropriate Data Mining task and algorithm then applying it. Evaluating the results of all appropriate algorithms leads to choose the optimal machine-learning model. The data used in this paper mainly focuses on diagnosing individuals with ASD symptoms, based on AQ-10 with several variables that usually affect the diagnosing result. Hence, the prediction model is considered a classification problem that results from either having ASD or not. Therefore, many suitable supervised

models were applied for the given task, the results were analyzed and evaluated. Before applying these machine-learning models, the feature selection process must be applied to support the evaluation results and the models' accuracy by removing the weak variables from the databases. The following feature selection techniques and supervised classification models were considered suitable for the ASD diagnosing task.

## IV. RESULTS AND DISCUSSION

This section gives an overview of the achieved results, the experiment process to solve the research question, and visualization of those results. It focuses on feature selection techniques and results, the machine-learning model's evaluation measures based on feature selection results, and the standard rules related to ASD detection that has been extracted by the best machine-learning model.

### A. Simulation and Implementation

The simulations and operations evaluate the strength of the statistical procedure and identify the machine-learning model's strengths and weaknesses using the confusion matrix that leads to the simulation result. In RStudio, a machine-learning model can give a confusion matrix as a model evaluation tool. Applying the model and extracting the matrix can be used to define a set of mathematical values to determine the efficiency of the model and choose the best model to anticipate the results of the database set. There are sets of values that are taken into account, and they are error rate, accuracy evaluation, sensitivity, specificity, and FN-value. The visualizations associated with these values were also constructed using the Tableau software.

### B. Feature Selection Results and Analysis

As mentioned earlier in Section 3, the ASD databases' prediction model is considered a classification-predicting problem that carries out a categorical label variable, with categorical input variables. The feature selecting algorithms mainly evaluate the relationship between ASD test results independently with each other variable in the database using two filter-based techniques: Chi-Squared and mutual information. The lines of code related to the mutual information techniques, which is the information gain method, and the Chi-Squared were applied to evaluate the autistic trait features in all the available datasets then compared the performance for each technique. All the databases have eight selected features. It shows the highest eight variables in performance for each database based on the result of the mutual information and the Chi-Squared techniques. Both feature selection techniques give the same result with a few differences in the order of the variables based on the efficiency in the toddler database. All databases relied on the AQ questions and showed a high correlation with the ASD diagnosing results. The variables in the feature selection techniques were the only variables used in the machine learning models to improve the model's performance.

Both feature-selecting algorithms show that the A4 variable in the Child database has the highest correlation with the target class, resulting from the ASD test. Also, both show that the A6 variable in the Adolescent database has the highest

correlation, which carries out the “S/he is good at social chit-chat” questionnaire sentence. For the Toddler database, both techniques show that the A9 has the highest correlation with the label class. The A9 in the Toddler database carry out the “S/he first words were typical ones” questionnaire sentence.

C. Machine Learning Model Evaluation Measures based on Feature Selection

The evaluation techniques used in this paper are based on the result of the confusion matrix of each machine-learning model. The models' performance can be evaluated by calculating the error rate, accuracy, sensitivity, and specificity.

D. Error Rate and Accuracy Evaluation

After applying all the machine learning models that fit the ASD classification problem in RStudio, the accuracy rate measurement is described in Table VI.

The above comparison shows that the Neural Network model has the highest accuracy rate measurement in each database compared with the other machine-learning models. The toddler database has the best accuracy results compared to the child and adolescent databases. The number of attributes in the toddler database is much higher than the other two databases, affecting the accuracy rate result. This result indicates that the toddler age group is the best age to diagnose if they have an ASD. Neural networks can learn complex and non-linear relationships. It can infer unseen relationships on unseen data and give the model the ability to generalize and predict unseen data. Fig. 2 depicts each database's trends based on the error rate measurement in percentage and the applied machine-learning models. The color indicates the three databases: the Adolescent database, Child database, and Toddler database. Since the figure deals with the error rate, the lower the value, the better in results performance. It shows that the toddler database controls their error rates better than the other databases, with rates between 0.96% and 5.75%. The adolescent database shows the highest error rates, which may be due to the small size of the database compared to the rest. Also, the child's psychological changes during the adolescent period may significantly affect the validity of expectations. In comparing the machine learning models used in this paper, the Neural Network model gives the lowest error rate measurements comparing with the rest models in all databases. This shows that the Neural Networks model does not perform well only on datasets with significant data attributes, such as the Toddler dataset, but also with datasets with a limited number of attributes, such as the Adolescent dataset.

E. Sensitivity and Specificity

Fig. 3 and 4 display the sensitivity rates and specificity rates derived by the SVM, Naïve Bayes, Neural Network, Random Forest, GBM, XgBoost, AdaBoost, and CV Boosting algorithms on the Child, Adolescent, and Toddler datasets. Both the sensitivity rates and specificity rates results generated by the considered algorithms on all the datasets have shown acceptable levels of performance.

Neural Networks Model has higher sensitivity and specificity rates than most of the remaining algorithms on all the available datasets. For the Child database, Neural Network

Model achieved a 96.3% sensitivity rate and 97.2% specificity rates, while the adolescent database achieved 100% and 90.9% sensitivity rate and specificity rates. For the Toddler database, the sensitivity and specificity rates achieved 98.6% and 100%. As shown in the previous dashboards, some of the machine-learning models cannot perform well in small datasets, while the Neural Network Model shows outstanding sensitivity and specificity rates.

TABLE VI. ACCURACY RESULTS

| ML model           | Accuracy in %  |                     |                  |
|--------------------|----------------|---------------------|------------------|
|                    | Child-database | Adolescent-database | Toddler-database |
| SVM                | 95.41284       | 91.89189            | 96.83544         |
| XgBoost            | 92.07921       | 75.51020            | 97.14286         |
| AdaBoost           | 90.13158       | 93.24324            | 94.60317         |
| CV Boosting        | 92.73084       | 82.37096            | 94.68691         |
| Neural Network     | 96.73203       | 96.10390            | 99.03537         |
| Random Forest      | 87.33333       | 88.60759            | 94.24920         |
| Naïve Bayes        | 92.53438       | 94.75806            | 94.59203         |
| Random Forest- GBM | 93.33333       | 93.67089            | 95.84665         |

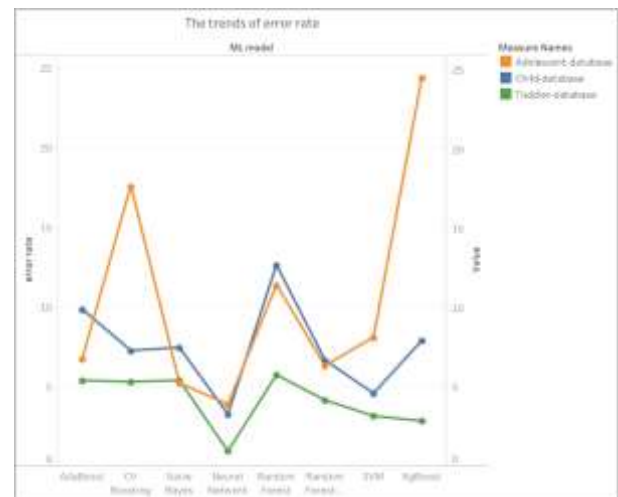


Fig. 2. Error Rate Results.

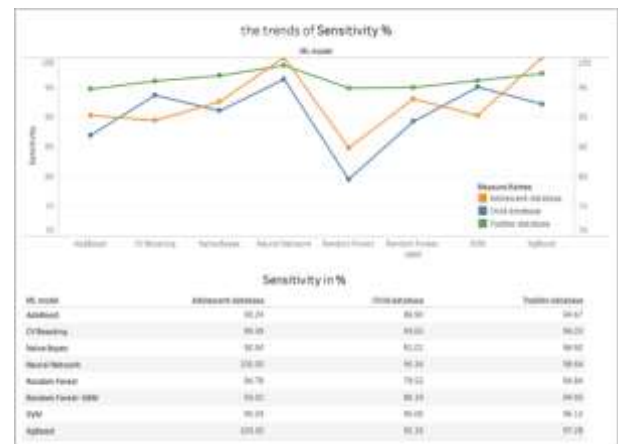


Fig. 3. Sensitivity Rate Results.

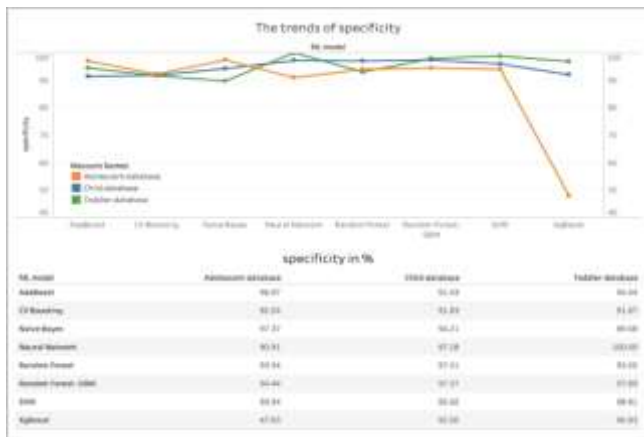


Fig. 4. Specificity Rate Results.

Overall, the results reported higher performance for the Neural Network model on the datasets when compared with the considered machine-learning algorithms, and these results are consistent with the error rate produced earlier and can be attributed to the non-redundant rules sets generated by the Neural Network model that will be discussed in the next section.

#### F. False Negative Rate

In a binary classification medical test, false negative is an error in which the test result incorrectly indicates that there is no case, when present in reality. Overall, many models show good results in the FN rate in the three databases, as shown in Table VII. Neural Networks Model has the best FN rate in two databases and a good result in the third one, the adolescent database.

TABLE VII. FN RATE

| ML model           | Child-database | Adolescent-database | Toddler-database |
|--------------------|----------------|---------------------|------------------|
| SVM                | 1.8%           | 2.7%                | 0.3%             |
| XgBoost            | 3.96%          | 24.5%               | 0.95%            |
| AdaBoost           | 2.6%           | 1.4%                | 2.1%             |
| CV Boosting        | 4.1%           | 3.6%                | 2.6%             |
| Neural Network     | 1.3%           | 3.9%                | 0%               |
| Random Forest      | 1.3%           | 2.5%                | 2.2%             |
| Naïve Bayes        | 2.8%           | 1.2%                | 3.3%             |
| Random Forest- GBM | 1.3%           | 2.5%                | 0.6%             |

#### G. Association Rules Result based on Feature Selection

The Association Rule is a rule-based machine-learning technique that discovers interesting relations between variables in large databases. Since the number of founded rules using the apriori function can be controlled by the values of Support and Confidence parameters, fixing the generated parameters will give balanced rules results. Many variables have frequently appeared within the rules that cover specific ASD characteristics within the databases. Variables A4, A1, and A6 in the Child database strongly influence the class labels. Additionally, variables A3, A4, and A6 appeared in

multiple rules in the Adolescent database, while items A2, A4, and A9 appeared in many rules in the Toddler database. The association rule shows that combining the results of two or three variables gives the best correlation between the variables and the diagnosing result, which is the class variable.

#### V. DISCUSSION

The paper's primary aim was to provide the best machine-learning model that diagnoses individuals with specific Autistic Spectrum Disorder symptoms. Several processes were needed to select the best machine-learning model. It was necessary to choose the most efficient ASD questionnaire diagnosing method and collect a high-quality database for each age group, which was done by the help of previous research studies and surveys conducted among the health professionals. With the available datasets and the applied data mining algorithms, the most accurate model was selected as the best machine-learning model to diagnose ASD symptoms. The first part focused on Knowledge Discovery in Databases processes, which is the best way to predict an individual with specific Autistic Spectrum Disorder characteristics using machine-learning models. The Knowledge Discovery in Databases processes that had been followed in this paper is as the following:

- Understanding the application domain
- Creating the data set
- Preprocessing and Data transformation
- Choosing the appropriate Data Mining task
- Choosing and applying the Data Mining algorithm
- Evaluation
- Using the discovered knowledge

These seven processes were followed in this paper, starting by searching and reading about Autistic Spectrum Disorder symptoms and diagnosing methods to choose the most efficient ASD questionnaire diagnosing method. Collecting and processing a high-quality database for each age group are two essential success factors for any data-mining paper. The data was managing in cooperation with the health professionals to ensure that the diagnosing results in the database are correct. The feature selection process supports the evaluation results and the models' accuracy by removing the weak variables from the databases.

After preparing the dataset, the data-mining task was set to a classification task, since the class label is categorical. Many suitable supervised models were applied for the given job, such as support vector machine, naïve Bayes, neural networks, and ensemble methods. The evolution process showed the result performance for each machine-learning models. All these processes answered the first research question, asking about how to use machine learning to diagnose individuals with specific Autistic Spectrum Disorder characteristics. The second part has been answered in the evaluation process in the Knowledge Discovery in Databases. They used machine-learning models were compared by the performance, which includes the accuracy rate, the sensitivity rate, and the

specificity rate. This comparison concluded that the neural networks model gives the best performance practice for the three databases.

## VI. RESULTS COMPARISON

The results of this paper give a better performance comparing with the papers reviewed in Section 2. The available databases in this paper are considered one of the best available datasets those days, since it deals with each age group as a separate database, and contains a good number of attributes. The toddler database using the Neural Network model gives the best accuracy result, while adolescent and child databases using the Neural Network model also have excellent accuracy results compared with the other models.

## VII. CONCLUSION AND FUTURE WORK

This paper aimed to provide useful and accurate ASD screening models to help parents and interested parties quickly diagnose their children's condition. Unfortunately, some families and adult patients do not have sufficient knowledge of ASD symptoms, so cases of autism spectrum disorder are not dealt with early. Artificial intelligence and machine learning are used at this time in most living areas, and their use in the field of medical diagnosis contributes to a pioneering step in using the available data as a tool for development and progress. All the primary seven processes of the KDD was used and described in this paper. These processes contain data gathering and data preprocessing, choosing an appropriate data mining approach to find patterns among the data and interpret them. Finally, the results were used for further research. The empirical results on the used datasets related to children, adolescents, and toddlers show that the neural networks model yielded the highest performance results compared to the other machine learning models used in this paper concerning predictive power, sensitivity, and specificity.

The development of this paper into an application program will provide families with a quick and straightforward scan tool using the lowest set of elements related to ASD, which contributes to increased accessibility and early detection. In the future, it is possible to develop this paper for use in the health system of the Ministry of Health and Prevention, where data are available for all patients registered in all the hospitals affiliated with the Ministry in this system. It is also possible to provide the departments of schools, kindergartens, and nurseries with an easy-to-use system for this paper to be applied to children for early detection. Another potential area for further use of this study could be the application of machine education and artificial intelligence models in health systems that store patient data for a range of diseases and health symptoms to contribute to the early detection of potential diseases.

## REFERENCES

- [1] American Psychiatric Association . (2013, 10 10). American Psychiatric Association . Retrieved from American Psychiatric Association : <https://www.psychiatry.org/>.
- [2] Ben-David, S. S.-S. (May 2014). Understanding Machine Learning: From Theory to Algorithms. -: Cambridge University Press.
- [3] Becerra-Culqui, T. A., Lynch, F. L., Owen-Smith, A. A., Spitzer, J., & Croen, L. A. (2018). Parental first concerns and timing of Autism Spectrum Disorder diagnosis. *Journal of autism and developmental disorders*, 48(10), 3367-3376.
- [4] Bishop-Fitzpatrick, L., Movaghar, A., Greenberg, J. S., Page, D., DaWalt, L. S., Brilliant, M. H., & Mailick, M. R. (2018). Using machine learning to identify patterns of lifetime health problems in decedents with autism spectrum disorder. *Autism Research*, 11(8), 1120-1128.
- [5] Bravo Oro A., N.-C. M. (2014). *Autistic Behavior Checklist (ABC) and Its Applications*. New York: Springer.
- [6] Carla A. Mazefsky, R. A. (2011, March -). PubMed Central. Retrieved May 30, 2012, from [ncbi.nlm.nih.gov: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3362998/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3362998/).
- [7] Duvetkot, J., van der Ende, J., Verhulst, F. C., Slappendel, G., van Daalen, E., Maras, A., & Greaves-Lord, K. (2017). Factors influencing the probability of a diagnosis of autism spectrum disorder in girls versus boys. *Autism*, 21(6), 646-658.
- [8] Gandhi, R. (2018, June 7). *towardsdatascience*. Retrieved June 7, 2018, from [towardsdatascience.com: https://towardsdatascience.com/support-vector-machine-introduction-to-machine-learning-algorithms-934a444fca47](https://towardsdatascience.com/support-vector-machine-introduction-to-machine-learning-algorithms-934a444fca47).
- [9] Greenhalgh, T. (1997). *How to read a paper. Papers that report diagnostic or screening tests*. London: University College London Medical School.
- [10] Grove, J., Ripke, S., Als, T. D., Mattheisen, M., Walters, R. K., Won, H., ... & Awashiti, S. (2019). Identification of common genetic risk variants for autism spectrum disorder. *Nature genetics*, 51(3), 431-444.
- [11] Hershy, A. (2019). Gini Index vs Information Entropy. Retrieved from [medium.com: https://towardsdatascience.com/gini-index-vs-information-entropy-7a7e4fed3fcb](https://towardsdatascience.com/gini-index-vs-information-entropy-7a7e4fed3fcb).
- [12] Hyde, K. K., Novack, M. N., LaHaye, N., Parlett-Pelleriti, C., Anden, R., Dixon, D. R., & Linstead, E. (2019). Applications of supervised machine learning in autism spectrum disorder research: a review. *Review Journal of Autism and Developmental Disorders*, 6(2), 128-146.
- [13] Ibrahim, S., Djemal, R., & Alsuwailem, A. (2018). Electroencephalography (EEG) signal processing for epilepsy and autism spectrum disorder diagnosis. *Biocybernetics and Biomedical Engineering*, 38(1), 16-26.
- [14] Jung, H. (2018). *medium.com*. Retrieved from <https://towardsdatascience.com/adaboost-for-dummies-breaking-down-the-math-and-its-equations-into-simple-terms-87f439757dcf>.
- [15] LeBarton, E. S., & Landa, R. J. (2019). Infant motor skill predicts later expressive language and autism spectrum disorder diagnosis. *Infant Behavior and Development*, 54, 37-47.
- [16] Lindner, L.-O. L. (2017). Is the Autism-Spectrum Quotient a Valid Measure of Traits Associated with the Autism Spectrum? A Rasch Validation in Adults with and Without Autism Spectrum Disorders. *Journal of Autism and Developmental Disorders*, 47.
- [17] Liu, W., Li, M., & Yi, L. (2016). Identifying children with autism spectrum disorder based on their face processing abnormality: A machine learning framework. *Autism Research*, 9(8), 888-898.
- [18] Lord, C., Elsabbagh, M., Baird, G., & Veenstra-Vanderweele, J. (2018). Autism spectrum disorder. *The Lancet*, 392(10146), 508-520.
- [19] Marvin, A. R. (2017). *Analysis of Social Communication Questionnaire (SCQ) Screening for Children Less Than Age 4*. US: [springer.com](http://springer.com).
- [20] Mason, Y. F. (1999). *The Alternating Decision Tree Learning Algorithm*. ICML '99 Proceedings of the Sixteenth International Conference on Machine Learning (pp. 124-133). San Francisco: Morgan Kaufmann Publishers Inc.
- [21] Mohammad Moshirpour, B. H. *Applications of Data Management and Analysis*.: Calgary, Canada: [springer](http://springer.com).
- [22] NIH. (2018, March 1). Retrieved from The National Institute of Mental Health Information Resource Center : <https://www.nimh.nih.gov/health/topics/autism-spectrum-disorders-asd/index.shtml>.
- [23] Omar, K. S., Mondal, P., Khan, N. S., Rizvi, M. R. K., & Islam, M. N. (2019, February). A machine learning approach to predict autism spectrum disorder. In *2019 International Conference on Electrical, Computer and Communication Engineering (ECCE)* (pp. 1-6). IEEE.
- [24] Pagnozzi, A. M., Conti, E., Calderoni, S., Fripp, J., & Rose, S. E. (2018). A systematic review of structural MRI biomarkers in autism spectrum

- disorder: A machine learning perspective. *International Journal of Developmental Neuroscience*, 71, 68-82.
- [25] Peebles, F. T. (2019). Early Autism Screening: A Comprehensive Review . *International Journal of Environmental Research and Public Health — Open Access Journal* , PMC6765988.
- [26] Sandin, S., Lichtenstein, P., Kuja-Halkola, R., Hultman, C., Larsson, H., & Reichenberg, A. (2017). The heritability of autism spectrum disorder. *Jama*, 318(12), 1182-1184.
- [27] Sharma, S. R., Gonda, X., & Tarazi, F. I. (2018). Autism spectrum disorder: classification, diagnosis and therapy. *Pharmacology & therapeutics*, 190, 91-104.
- [28] Stevens, E., Dixon, D. R., Novack, M. N., Granpeesheh, D., Smith, T., & Linstead, E. (2019). Identification and analysis of behavioral phenotypes in autism spectrum disorder via unsupervised machine learning. *International journal of medical informatics*, 129, 29-36.
- [29] Tan, C. D. (2018). “I’m a normal autistic person, not an abnormal neurotypical”: Autism Spectrum Disorder diagnosis as biographical illumination. *Social Science & Medicine*, 197, 161-167.
- [30] Thabtah, F. F. (2017). uci. (Manukau Institute of Technology) Retrieved from uci.edu: <http://archive.ics.uci.edu/ml/datasets/Autism+Screening+Adult>
- [31] Thabtah, F., & Peebles, D. (2020). A new machine learning model based on induction of rules for autism detection. *Health informatics journal*, 26(1), 264-286.

# Educational Tool for Generation and Analysis of Multidimensional Modeling on Data Warehouse

Elena Fabiola Ruiz Ledesma<sup>1</sup>, Elizabeth Moreno Galván<sup>2</sup>

Enrique Alfonso Carmona García<sup>3</sup>, Laura Ivoone Garay Jiménez<sup>4</sup>

Instituto Politécnico Nacional, Escuela Superior de Cómputo, Ciudad de México, México<sup>1</sup>

Instituto Politécnico Nacional, Unidad Profesional Interdisciplinaria en Ingeniería y Tecnologías Avanzadas  
Ciudad de México, México<sup>2,3,4</sup>

**Abstract**—The curricular inclusion of topics, study plans, and teaching programs related to the study of Data Science has been trending mostly in higher-level education for the last years. However, the previous knowledge requirements for students to adequately assimilate these lessons are more specialised than the ones they obtain during secondary education. On the one hand, the interaction with complex techniques and materials is needed, and on the other, tools to practice on-demand are required in the current learning. So, this is an excellent opportunity for the creation of data analysis tools for educational purpose that could be considered as a starting point of a broad area of application. This paper presents a pedagogical support tool aimed to facilitate the student approach to the basic knowledge of data mining through the practice of the analysis of online analytical processing (OLAP). It is a prototype that allows the visualisation of the multidimensional cubes generated with all possible combinations of the dimensions of the data set, as well as their storage in databases, the recovery operations for views, and the implementation of an algorithm for the selection of the optimal view set for materialising the set of records resulting from a search of the database, and computing the materialisation costs and total records recovered. The prototype also carries out and present recurrent patterns and association rules while considering factors such as support variables and reliability. All of this is done explicitly to aid the students to comprehend the generation process of data cubes in the data mining discipline.

**Keywords**—Educational data mining; data cube; view materialization; educational software

## I. INTRODUCTION

The extensive analysis of Big Data [1] comes from a branch on statistical analysis that companies used to identify spending trends, which is used to predict the consumer's behaviour and to analyse commercial activities. From this initial idea, several data handling techniques have been created, such as Data Mining (DM) or Machine Learning (ML). They have been successfully applied to a range of human effort areas including detection of illness and mobile health [3], [4], [5], [6], environmental and pollution studies [7], [8], [9], [10], being these only examples of the many areas to which these techniques have been recently applied. In the educational field, there are different applications such as educational data mining and learning analytics [11], [12], whose purpose is oriented towards the designing of algorithms, methods, and models, that will allow exploring data from learning environments.

In order to study the data, there is a constant need for a data warehouse. Gathering data from a company or organisation in a single database helps analysts and managers to support their decisions or to find valuable data, but reduce the extraction time and cost are constant requirement [13], [14], [15]. Moreover, the design and construction of a data warehouse require the application of extraction, integration, transformation, and data cleaning processes [16]. So the data warehouse increases their dimension, and a multidimensional model is going to be required. This model is going to be described and defined along with the optimal view set to be materialised. The last process consists of determining the necessary tools for viewing data.

The most traditional tools for data mining and automated learning are becoming insufficient as the tendencies in technology progress, prompting the creation of new and increasingly powerful, complex solutions. In the academic sense, from the students' point of view, the vast availability of these new tools and methods, and its many overlapping uses represents a challenge. An excess in variety makes selecting a study and comprehension strategy even harder than it usually is, especially in their primer approach.

This document is focused on the data preparation from a data warehouse, and it uses a software tool that automatically designs a multidimensional model and helps into the creation and storing of the data warehouse in a database. Once the outline and the hypercube materialisation through previous calculation are ready, the results of the information retrieval are sped up. Also, the optimal view set selection algorithm proposed by Harinarayan is applied in materialising [17]. In the last step, the tool can determine the frequent patterns present in data, and it also calculates the association rules, considering predetermined support and reliability.

Parameters about the performance of this tool, such as the effectiveness and time-saving in calculating, shows an improved performance over the manual method of the same procedure that it is commonly presented to the students in the courses about essential data mining topics.

## II. LITERATURE REVIEW

Data analysis is the process of working on data in order to discover useful information for business decision making. Data Analytical Tools are software developed to perform data analysis tasks such as process and manipulate data, analysing



relationships and correlations, as well as identifying patterns and trends for interpretation.

In recent years, a large variety of Data Analytical Tools have been created to carry out data science tasks. However, in this section, we will present some of the most used and accessible tools that are currently available for teaching this topic, all of them related to Educational Data Mining [18] or Learning Analytics [11][19] instead of the broader array of tools that could be used for the most modern statistics analysis, since these are complex for the beginner student.

Two characteristics that are commonly present in Data Analytical Tools on the educational field; they are the integration of the functionalities of data mining, and the application of techniques dedicated to data mining for didactic purposes [20]. Nevertheless, innovation in education has become relevant, so several projects appear to apply it. For example, the Hadoop Ecosystem has been designed to help researchers and students in all aspects of typical data analysis and automatic learning processes (Machine Learning) [21].

In the revision work by S. Slater et al. [22], there is an analysis of several didactic and research tools for educational data mining, classifying them as follows:

- Data Analytical Tools for the handling, cleaning, and formatting of the data, per example: Microsoft Excel and EDM Workbench.
- Data Analytical Tools for model selection and testing, also identification, mapping, exploration and analysis of relationships such as RapidMiner, Weka, KEEL, KNIME, Orange, and SPSS.
- Data Analytical Tools for visualisation of the structure of the tree methodology as Tableau, d3js, and InfoVis.

#### A. Data Analytical Tools

S. Yadav and Urbina provide a list of analytical tools, and their descriptions, as well as the definition [2], [27]. In Table I, the most relevant characteristics are listed, and the tools that could be used to educational purposes are identified.

TABLE I. TOOLS FOR DATA ANALYSIS

| Tool                          | Description  |
|-------------------------------|--|
| WEKA*                         | Implements algorithms for data preprocessing, classification, regression, grouping, association rules, and viewing. It is free to use under the Public License GNU, and it contains a wide range of modelling and data processing techniques. [2].   |
| Orange*                       | This software allows preprocessing, information filters, data modelling, evaluation, and exploration of modelling techniques.  |
| Rapid Miner*                  | Predictive analysis tool. It is sturdy, easy to use, and has a broad open-source community where the users can integrate their self-made specialised algorithms. It provides the user with learning schemes, models, and algorithms from WEKA and R [2].   |
| Rattle                        | Free, open-source data mining toolset that is written in the statistic language R. It presents visual and statistical data summaries [2].  |
| Knime*                        | Integrative software allows data processing, analysis, and exploration as well as advanced prediction algorithms and machine learning.   |
| CLUTO                         | This tool is for grouping high and low dimension data analysis. It has multiple classes of algorithms for grouping such as partition, agglomeration, and graphic-based partition, similarity/distance functions as Euclidian distance, cosine, correlation coefficient, extended Jaccard, and even self-defined functions [2]. |
| Jaspersoft BI Suite           | An open-source suite produces reports based on database columns, reducing the data from the sources to tables and interactive graphics.  |
| Pentaho Business Analytics    | Software platform that simplifies the information inclusion from the different sources   |
| Talend Open Studio            | Offers a development environment for linking Hadoop data processing works.   |
| Splunk                        | This tool creates a data index such as a book's structure or a text block.   |
| Apache Storm                  | A distributed computing system that allows the user to process unlimited data in a reliable real-time way.   |
| Apache Drill                  | It is an SQL search engine for Big Data exploration. It has been designed from scratch to allow high-performance analysis in semi-structured data.   |
| Cassandra                     | It is used by large active data sets, coming from Netflix, Twitter.  |
| HBase                         | A distributed database management system built upon the Hadoop file system and oriented towards columns format,  |
| Neo4j                         | It is a native graphics database management system which uses the data relations as the first-class entities. It has upgraded performance in comparison to relational databases.   |
| CouchDB                       | A database management system that is wholly dedicated to web applications, storing data in JSON files.   |
| OrientDB                      | It combines the flexibility of file databases with graphic databases.  |
| FlockDB                       | It is an open-source database management system that uses a wide but shallow network graphics. It was designed to store social graphics.   |
| MOA (Massive Online Analysis) | It is a project designed in partnership with WEKA that offers flow analysis online for various WEKA algorithms and with the same user interface [24].  |
| MADlib                        | It is a collection of SQL-based algorithms; it includes grouping, classification, regression and themed models as well as validation tools [23].   |
| Dato, before GraphLab         | It is an independent product that can be connected to Hadoop for graph analysis and machine learning tasks [25].   |

\* For educational purpose, (Modified from Source: [2],[23], [24],[25])

In the literature review, some of the tools are focused on the learning-teaching area of data science. However, their focus is not entirely didactic for non-experienced users. Because they are for specialised tasks and depending on the data set type, a choice of tool is made. So, these tools have more complexity than required for a basic implementation. In many cases, these tools are developed for more experienced users or area professionals, what it offers an opportunity area for pedagogical tools in this field.

### III. ARCHITECTURE AND DESIGN OF THE TOOL

The Data Analytical Tool presented in this paper is developed in Java, and it generates the multidimensional variant for Online Analytical Processing (OLAP) named multidimensional structure (MOLAP). For this purpose, it uses relational database modelling technology for the construction of the data warehouse (DW) in a MySQL system by reading a data source in CSV format separated by commas. In this way, the system creates a multidimensional model generated by a table for each dimension or column in the data source.

The Data Mining based on Lattice is a method to organise data in domains determined by combinations of the dimensions of a dataset [26]. These combinations can be determined by information retrieval with SQL structures named views. The system can calculate the views from the MySQL database and managing a cache memory file by making use of a linked list structure [27], and a modified B-Tree [28], where each node in the tree constitutes a view also named cuboid from the Lattice. This system generates both the logic and the visual representation of a data cube.

Materialised views are structures that improve data access time by precomputing intermediary results; an effective technique for improving query performance is using indexing [29]. In this sense, the algorithm proposed by Harinarayan [17],[30] is used to improve the efficiency of the model for determining the pre-calculated views to materialise, and his algorithm is implemented into the Data Analytical Tool. Given the educational approach of this development, the execution tests were carried out in a synthetically created data cube; that is, with data sources created by an additional computer program, and the measurements were randomly generated.

#### A. Architecture of the Data Analytical Tool

The general tool's architecture is shown in Fig. 1, and the specific elements are described and detailed in the following paragraphs.

The data are obtained from an external data source in CSV format, which is used to fill out a database structured into tables, matching each file's header with a field of the table. The cube is comprised of a vector (D1, D2, ..., DN) with N dimensions corresponding to the attributes of the database. The generated cube or Lattice L contains a cartesian product of the N dimensions. Each cuboid represents a possible useful aggregation from structured language queries (SQL) known as views. The materialisation of these views and its efficiency is based on the Harinarayan algorithm implementation. In this work, queries and aggregations are optimised, first, by application of operations like slicing in OLAP that generates data columns corresponding to single values with at least one

dimension. Then, it helps the visualisation and recompilation of information about a specific dimension. Finally, dice operation in OLAP is provided, which selects a subset of dimensions considering a specific values range in each dimension.

In the analysis of the patterns, the data source must be a binary matrix, this means that the columns constitute dimensions and the rows are corresponding with the records. So if a record has the dimension, the intersection is filled with a value 1, or 0 on the contrary case.

#### B. Lattice Construction

For the software implementation, making use of the modified B-Tree structure, each view is represented by a node, as shown in Fig. 2. In this case, dimensions, the node attribute is previously determined by the number of attributes of this node; ancestors are the previous nodes, and descendants are the subsequent nodes. Both last sets contain at least one shared dimension with the actual one.

For example, for the three-dimensional data source,  $N=3$  and  $2^N$  possible groupings (nodes) are generated. Their relationship can be appreciated in Fig. 3, and it is essential to note the complexity involved in that excepting the apex and base (Node 7), each node possesses various descendants and several ancestors at the same time.

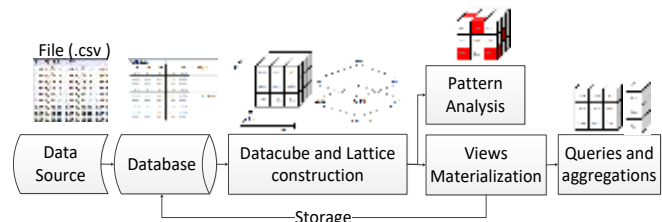


Fig. 1. The Architecture of the Proposed Tool).

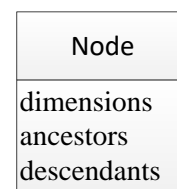


Fig. 2. Node Structure.

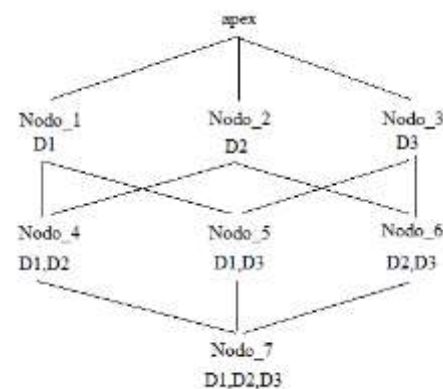


Fig. 3. Three-Dimensional Lattice Structure.

The information of the Lattice's constitution is stored in each node as has been shown in Fig. 2. So apex is 0 dimension, Nodes 1,2, and 3, have one dimension and Node 4,5,6, have two dimensions associated. Node 7 is considered the base cuboid. In this example, Fig. 4 shows the structure of Node 5 as an example; it has two ancestors with dimensions D1 and D2, and its descendants that have its dimensions and D2.

The order of the interconnected node net structure described by the Lattice contains all views that can be used to get any query related to a business question, as well as materialising or to pre-calculating the cuboids. However, it is crucial to know the physical space limitations in the storage unit. It is also recommended to materialise the base cuboid (full detail, apex) as it can be used to respond to any question, and then move on to the less costly views, which results in less time and resources to obtain the desired answers. The generated structure allows the application of other algorithms, so the computational cost that is inverted is justified.

The pseudocode for the creation of the logical structure for the data cube is shown in Fig. 5.

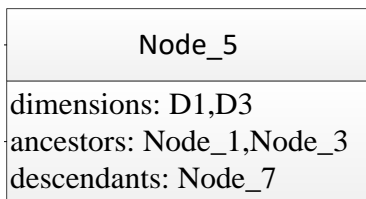


Fig. 4. Representation of the Values of an Example Node.

```

// Creating a dimension list
for of each dimension
ListDimensions <- name dimension

// Creation of apex node
create a new node Nodeactual
Nodeactual.dimensions <- apex
Nodeactual.ancestors <- null
Nodeactual.descendants <-null

// Creation of Nodes for storing one dimension cuboids
for each dimension in ListDimensions
create a new node Nodeactual
Nodeactual.dimensions <- name dimension
Nodeactual.ancestors <-apex
Nodeactual.descendants <-null

//Creation of nodes for storing cuboids of 2 or more dimensions
Until all levels are covered
for each dimension in ListDimension
create new node Nodeactual
Nodeactual.descendants <-null
for all groups of n dimensions in ListDimensions
for all nodes in previous level
if dimension in Nodenivel -1
    Add Nodenivel -1 to the list of ancestors from the actual
    node
    Add Nodeactual to the list of ancestors from Nodenivel -
    1
    
```

Fig. 5. Multidimensional Cube Generation Pseudocode.

### C. View Materialisation

The formulation of business questions can be carried out through structure query language expressions (SQL), based on the database previously stored. Furthermore, it is possible to use operators on numeric-type dimensions of the views resulting in answering expressions, for example:

*SELECT field, Op field FROM table GROUP BY field*

Where, *field*: is a subset of database attributes or dimensions D1, D2, ..., DN and *Op (field)* is an operation in a numeric-type dimension, as COUNT, SUM, MAX, MIN.

The cost of generation of a view, represented by C(v) is associated to the computational cost of using a view considering that it decreases with the number of dependent relationships, thus, the calculation cost is divided between the total number of dependent relationships.

Therefore, the construction of the net is modified, as shown in Fig. 6.

An advantage of this new structure is to get a straightforward application of Harinarayan's Greedy algorithm proposal [17]. An efficient view generation is done, as shown in the pseudocode in Fig. 7. After selecting a view set named S, the benefit of the view v, denoted by B(v, S) is calculated. B is the difference in cost of storing a descendant view and the cost of its ancestor view and then multiply the difference by the number of relationships dependent to view v. The only views that benefit from the materialisation of v are the ones that can be calculated from v, including the v itself. The list of these views is named as w.

Therefore, the total benefit is the sum of all the benefits from the w set. In Fig. 7, the pseudocode can be appreciated.

The for testing the of use of the developed Data Analytical Tool and its effectiveness in the classroom, it was used in production for the analysis of data from several experimental datasets, speeding up the obtaining of results for the respective case studies was reported. Then a follow session was done, and the results are presented in the next section.

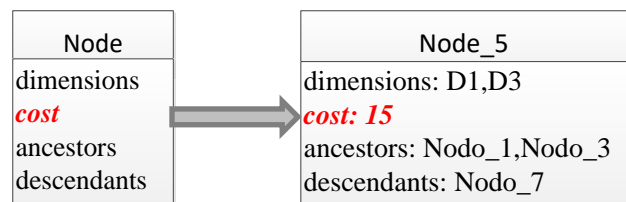


Fig. 6. Addition of the Cost Attribute to the Nodes.

```

S = {top view};
for i=1 to k do begin
    select that view v not in S such that B(v,S) is maximized;
    S = S union {v};
end;
resulting S is the greedy selection;
    
```

Fig. 7. Greedy Algorithm (Source: Harinarayan[17]).

#### IV. RESULTS

In this section are presented the efficiency of the tool to show the results about the calculus of the cuboid, definition of the association rules, and MOLAP visual representation using synthetic data as input.

The test data consists of a synthetic input composed of five dimensions named A, B, C, D and E respectively fulfilled with four records with random 1 and 0 numeric data values. The small size of the synthetic dataset was selected for display purposes as shown in Fig. 8, but the data capacity is limited by MySQL restrictions whose consists on the storage engine such as InnoDB that supports a maximum of 65,535 bytes per row limited by the data type that it hosts, that is approximate of 1.073.741.824 rows.

The corresponding lattice representation must look as shown in Fig. 9, where a labelled node represents each combination of dimension.

In the interface, this Lattice is represented by a button for each node. It replaces the names of the dimensions and respectively views by numbers, preventing dimension names from being longer, so the example lattice is shown in Fig. 10.

**Dataset**

| A | B | C | D | E |
|---|---|---|---|---|
| 1 | 1 | 1 | 0 | 0 |
| 1 | 1 | 0 | 1 | 0 |
| 1 | 0 | 0 | 1 | 1 |
| 1 | 1 | 0 | 1 | 0 |

Fig. 8. Test Synthetic Dataset.

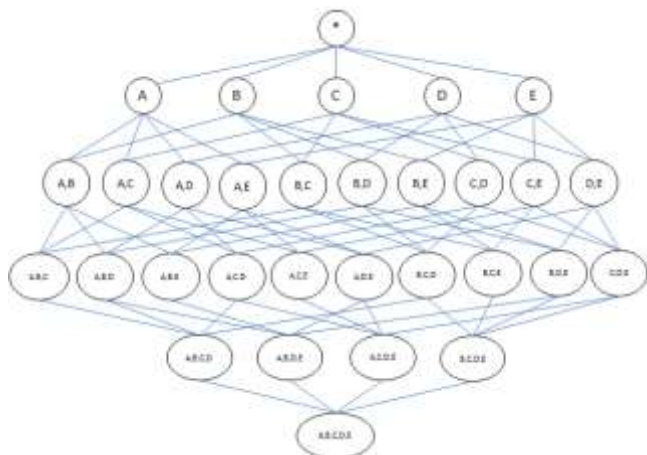


Fig. 9. Lattice Representation.



Fig. 10. Lattice Representation of the Synthetic Dataset Created by the Tool.

After the lattice generation, all the simple views whose consists of SQL sentences such as *select a,b from table*, are calculated automatically. Then, the visualisation of any view is obtained clicking on the button that represents the node of interest. A new window will appear showing the records that comply with the query. In Fig. 11, can be appreciated the result of clicking the button with numbers 0,1,2 that correspond to the node that relates the dimensions A, B and C.

Once, all the possible views to generate are available, the analysis of which ones are adequate for being materialised is carried out using the Harinarayan algorithm, the calculation provided by the tool is shown in Fig. 12 with N=4.

The result is the materialised view set as  $S = \{V2, V4, V7, V5, V10, V6, V12, V3\}$ , according to the employed method. In this output table, the student could analyse the procedure of optimisation whose manual calculations would have been complicated and time-consuming.

In data mining, the technique used to find item sets, subsequences, or substructures that appear in a data set frequently (patterns), requires the following definitions and operations.

Considering  $I = \{I_1, I_2, \dots, I_m\}$  as a set of items or dimensions, D the task-relevant data, T a set of items such that  $T \subseteq I$ . Let A be a set of items, a transaction T is said to contain A if and only if  $A \subseteq T$ . An association rule is an implication of the form  $A \Rightarrow B$ , where  $A \subset I, B \subset I$ , and  $A \cap B = \emptyset$ . The rule  $A \Rightarrow B$  has confidence c in the transaction set D, where c is the percentage of transactions in D containing A that also contains B. It is taken to be the conditional probability,  $P(B|A)$ . Then

$$\text{support}(A \Rightarrow B) = P(A \cup B) \tag{1}$$

$$\text{confidence}(A \Rightarrow B) = P(B|A) \tag{2}$$

Rules are called strong when (1) and (2) are satisfied with a, a minimum support threshold (min sup) and a minimum confidence threshold (min conf). If the relative support of an itemset I satisfy a pre-specified minimum support threshold, then I is a frequent itemset [31].

To show this calculation by the program, the frequent patterns that comply with the minimum support requirements (i.e. min sup: 2) are highlighted with black buttons in the interface, as shown in Fig. 13, so the students could visualise the combination of frequent dimensions in the dataset.

From Equation (2), we have

$$\text{confidence}(A \Rightarrow B) = P(B|A) = \frac{\text{support}(A \cup B)}{\text{support}(A)} = \frac{\text{support count}(A \cup B)}{\text{support count}(A)} \tag{3}$$

Equation (3) shows that the confidence of rule  $A \Rightarrow B$  can be derived from the support counts of A and  $A \cup B$ , and it is straightforward to derive the corresponding association rules  $A \Rightarrow B$  and  $B \Rightarrow A$  [31]. Then, the association rule for a relation e.g.  $\{A,D\}$ , is calculated as follows:  $\text{conf}(\{B\} \rightarrow \{A,D\})$  in the interface:  $\text{conf}(\{1\} \rightarrow \{0,3\}) = \text{supp}(\{0,3\}) / \text{supp}(\{1\}) = 3 / 3 = 1.0$ . At last, the association rules are determined for those who

completed the minimum confidence value (i.e. 60% or 0.60), and the interface displays the result, as shown in Fig. 14.

The students reported that the intuitive interface of the tool focused on concrete operations supported them not to spend extra time in software configurations or learning a complicated interface for the same purpose.

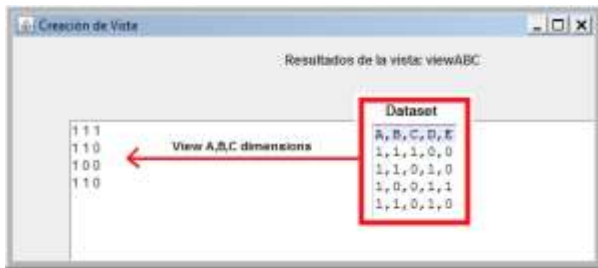


Fig. 11. Test of a Tree-Dimension view Execution.

| Views | Relation-ships | C(v) | Benefit 1 <sup>st</sup> iteration | Benefit 2 <sup>nd</sup> iteration | Benefit 3 <sup>rd</sup> iteration | Benefit 4 <sup>th</sup> iteration |
|-------|----------------|------|-----------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|
| V1    | A              | 4    | DF (4,4) * 7 = 0                  | DF (4,4) * 2 = 0                  | DF (4,4) * 2 = 0                  | DF (4,4) * 1 = 0                  |
| V2    | B              | 3    | DF (4,1) * 7 = 21                 | ---                               | ---                               | ---                               |
| V3    | C              | 1    | DF (4,3) * 8 = 8                  | DF (4,3) * 2 = 2                  | DF (4,3) * 2 = 2                  | DF (4,3) * 1 = 2                  |
| V4    | D              | 3    | DF (4,1) * 7 = 21                 | ---                               | ---                               | ---                               |
| V5    | E              | 1    | DF (3,1) * 7 = 14                 | DF (3,1) * 2 = 4                  | ---                               | ---                               |
| V6    | A,B            | 3    | DF (3,2) * 8 = 8                  | DF (3,2) * 3 = 4                  | DF (3,2) * 3 = 3                  | ---                               |
| V7    | A,C            | 1    | DF (3,0) * 7 = 21                 | ---                               | ---                               | ---                               |
| V8    | A,D            | 3    | DF (1,0) * 8 = 8                  | DF (1,0) * 3 = 3                  | DF (1,0) * 1 = 1                  | DF (1,0) * 1 = 1                  |
| V9    | A,E            | 1    | DF (1,0) * 7 = 7                  | DF (1,0) * 2 = 2                  | DF (1,0) * 1 = 1                  | DF (1,0) * 1 = 1                  |
| V10   | B,C            | 1    | DF (3,1) * 8 = 16                 | DF (3,1) * 2 = 4                  | ---                               | ---                               |
| V11   | B,D            | 2    | DF (4,1) * 3 = 9                  | ---                               | ---                               | ---                               |
| V12   | B,E            | 0    | DF (4,1) * 4 = 12                 | ---                               | ---                               | ---                               |
| V13   | C,D            | 0    | DF (4,0) * 3 = 12                 | DF (4,1) * 1 = 3                  | DF (4,1) * 1 = 3                  | ---                               |
| V14   | C,E            | 0    | DF (4,0) * 4 = 16                 | ---                               | ---                               | ---                               |
| V15   | D,E            | 1    | DF (4,0) * 3 = 12                 | ---                               | ---                               | ---                               |
| V16   | A,B,C          | 1    | DF (4,1) * 4 = 8                  | ---                               | ---                               | ---                               |
| V17   | A,B,D          | 1    | DF (3,0) * 4 = 12                 | DF (3,0) * 1 = 3                  | ---                               | ---                               |
| V18   | A,B,E          | 0    | DF (3,0) * 3 = 9                  | ---                               | ---                               | ---                               |
| V19   | A,C,D          | 0    | DF (3,0) * 4 = 12                 | ---                               | ---                               | ---                               |
| V20   | A,C,E          | 0    | DF (1,0) * 4 = 4                  | DF (1,0) * 1 = 1                  | ---                               | ---                               |
| V21   | A,D,E          | 1    | DF (4,0) * 2 = 8                  | ---                               | ---                               | ---                               |
| V22   | B,C,D          | 0    | DF (4,0) * 2 = 8                  | ---                               | ---                               | ---                               |
| V23   | B,C,E          | 0    | DF (4,0) * 2 = 8                  | ---                               | ---                               | ---                               |
| V24   | B,D,E          | 0    | DF (3,0) * 7 = 8                  | ---                               | ---                               | ---                               |
| V25   | C,D,E          | 0    | ---                               | ---                               | ---                               | ---                               |
| V26   | A,B,C,D        | 0    | ---                               | ---                               | ---                               | ---                               |
| V27   | A,B,D,E        | 0    | ---                               | ---                               | ---                               | ---                               |
| V28   | A,C,D,E        | 0    | ---                               | ---                               | ---                               | ---                               |
| V29   | B,C,D,E        | 0    | ---                               | ---                               | ---                               | ---                               |
| V30   | A,B,C,D,E      | 0    | ---                               | ---                               | ---                               | ---                               |

Fig. 12. Output of Benefit Calculation by the Harinarayan Algorithm Implementation.

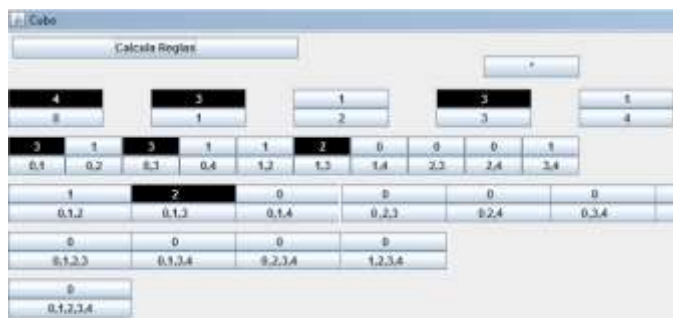


Fig. 13. Frequent Patterns Visualisation by Analytical Tool.

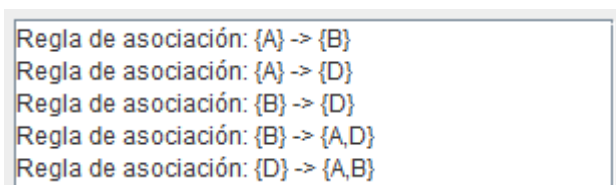


Fig. 14. Output Program of Association Rules Calculation.

## V. CONCLUSIONS

Among the tendencies that have been found in a revision of the techniques reported, they are not oriented toward the pedagogic aspect of teaching data science topics; on the contrary, they use data science as an automatic learning tool for other areas.

In the learning of data mining, the analysis of data cubes is a technique that has prevailed as an efficient form of data analysis. However, its manual design is a very challenge task to be carried out within large volumes of data, so that is why the automation of the processes through computational tools constitutes an excellent aid for the data analysts. In this sense, the presented tool helps automatise the data cube tasks plus the storage model in vectors/matrixes. Usage of the tree structure gives a natural indexation and provides an efficient extraction of the data thanks to the pre-structuring of the added data. All the advantages of automated calculation can be explained in work out session with easy examples to explore the method in various cases as well as testing the solutions for exercises.

The visual presentation and the interaction with the consequence of the changes could improve the understanding of data mining because it constitutes a reinforcement to the constructivism approach in education, that is why the tool is developed with a visual interface focused on a data analysis task. Even though the implemented algorithms are not the only ones that can be used to perform these tasks, they are considered as the basis for a well understanding of more complex proposals. In this preliminary results with students, they could explore the system capabilities to analyse their dataset being able to obtain results in less time and effort than manually, as well as obtain a new data warehouse in MySQL for future tests.

## VI. FUTURE WORK

Even though this tools us completely functional, some improvements could be made. Firstly, testing the tool in control and observation student groups to have feedback of the student and learn about the effect of this digital resources in a virtual class. Besides, adding more functions and embedding the description of the processes in the interface could made this tool a self-learning tool.

In the technical approach, if the memory is limited then, the structure's baseload could be accelerated via data chunks, and it could improve the time consumption for large datasets. Moreover, a further study of alternative algorithms for the data cube creation algorithm using the tree structure could be implemented and could help the students to compare the performance in their practice in this learning tool.

## REFERENCES

- [1] Mohanty S., Jagadeesh M., Srivatsa H, "Big Data" in the Enterprise. In: Big Data Imperatives. Apress, Berkeley, CA, 2013.
- [2] A. B. Urbina, & De la Calleja, J., "Brief review of educational applications using data mining and machine learning", Revista Electrónica de Investigación Educativa, 19(4), 84-96. <https://doi.org/10.24320/redie.2017.19.4.1305>, 2017.
- [3] P. Vijayakumar, S. M. Ganesh, L. J. Deborah, and B. S. Rawal.: A new SmartSMS protocol for secure SMS communication in m-health environment, Comput. Electr. Eng., vol. 65, pp. 265–281, 2018.

- [4] Y. Kazemi and S. A. Mirroshandel.: A novel method for predicting kidney stone type using ensemble learning, *Artif. Intell. Med.*, vol. 84, pp. 117–126, 2018.
- [5] M. Echeverría, A. Jimenez-Molina, and S. A. Ríos.: A semantic framework for continuous u-health services provisioning, *Procedia Comput. Sci.*, vol. 60, no. 1, pp. 603–612, 2015.
- [6] U. R. Acharya et al.: Data mining framework for breast lesion classification in shear wave ultrasound: A hybrid feature paradigm, *Biomed. Signal Process. Control*, vol. 33, pp. 400–410, J. Wang, R. Boesch, and Q. X. Li.: A case study of air quality - Pesticides and odorous phytochemicals on Kauai, Hawaii, USA, *Chemosphere*, vol. 189, pp. 143–152, 2017.
- [7] Q. Wang, J. Wang, M. Z. He, P. L. Kinney, and T. Li.: A county-level estimate of PM<sub>2.5</sub> related chronic mortality risk in China based on multi-model exposure data, *Environ. Int.*, vol. 110, no. February 2017, pp. 105–112, 2018.
- [8] D. Uni and I. Katra.: Airborne dust absorption by semi-arid forests reduces PM pollution in nearby urban environments, *Sci. Total Environ.*, vol. 598, pp. 984–992, 2017.
- [9] M. A. Bari and W. B. Kindzierski.: Ambient fine particulate matter (PM<sub>2.5</sub>) in Canadian oil sands communities: Levels, sources and potential human health risk, *Sci. Total Environ.*, vol. 595, pp. 828–838, 2017.
- [10] K. R. Malik, Y. Sam, M. Hussain, and A. Abuarqoub.: A methodology for real-time data sustainability in smart city: Towards inferencing and analytics for big-data, *Sustain. Cities Soc.*, vol. 39, no. April, pp. 548–556, 2018.
- [11] Sushil S. Chaurasia, Devendra Kodwani, Hitendra Lachhwani, Manisha Avadhut Ketkar, “Big data academic and learning analytics: Connecting the dots for academic excellence in higher education”, *International Journal of Educational Management*, ISSN: 0951-354X, 2018.
- [12] Siemens George, “Learning Analytics: The Emergence of a Discipline”, Volume: 57 issue: 10, page(s): 1380-1400, 2013.
- [13] L. Zhang and J. Wen.: A systematic feature selection procedure for short-term data-driven building energy forecasting model development, *Energy Build.*, vol. 183, pp. 428–442, 2019.
- [14] F. Wang and J. Liang.: An efficient feature selection algorithm for hybrid data, *Neurocomputing*, vol. 193, pp. 33–41, 2016.
- [15] Y. Lin, H. Wang, S. Zhang, J. Li, and H. Gao.: Efficient quality-driven source selection from massive data sources, *J. Syst. Softw.*, vol. 118, pp. 221–233, 2016.
- [16] Z. Manbari, F. AkhlaghianTab, and C. Salavati.: Hybrid fast unsupervised feature selection for high-dimensional data, *Expert Syst. Appl.*, vol. 124, pp. 97–118, 2019.
- [17] V. Harinarayan, A. Rajaraman, and J. D. Ullman.: Implementing data cubes efficiently, *SIGMOD*, 1996.
- [18] A. Peña-Ayala.: Educational data mining: A survey and a data mining-based analysis of recent works, *Expert Syst. Appl.*, vol. 41, no. 4 PART 1, pp. 1432–1462, 2014.
- [19] Stefan Slater, Srećko Joksimović, Vitomir Kovanovic, Tools for Educational Data Mining: A Review, *Journal of Educational and Behavioral Statistics* 42(1):85-106 · January 2017.
- [20] Muyesser Eraslan Yalcin, Birgul Kutlu , “Examination of students’ acceptance of and intention to use learning management systems using extended TAM”, *Volume 50, Issue 5*, pp. 2414-2432, 2019.
- [21] Sara Landset, Taghi M. Khoshgoftaar, Aaron N. Richter and Tawfiq Hasanin, “A survey of open source tools for machine learning with big data in the Hadoop ecosystem”, Landset et al. *Journal of Big Data*, 2015.
- [22] S. Slater, S. Joksimovic, V. Kovanovic, R. Baker, and D. Gasevic.: Tools for educational data mining : a review,” January, 2017.
- [23] Sonam Yadav, “Open Source Big Data Databases Tools”. PCQuest; Gurgaon, 2016.
- [24] Hall, M., Frank, E., Holmes, G., Pfahringer, B., Reutemann, P., & Witten, I. H., “The WEKA data mining software: an update”, *SIGKDD Explorations*, 11(1), 2009.
- [25] Demsar, J., Curk, T., Erjavec, A., Gorup, C., Hocevar, T., Milutinovic, M., et al., “Orange: data mining toolbox in python”, *Journal of Machine Learning Research* 14, 2349-2353, 2013.
- [26] Yichang, “Data Mining method based on Lattice”, 2nd International Conference on Consumer Electronics, Communications and Networks (CECNet), China, 2012.
- [27] Karuna, Gupta G. “Dynamic Implementation Using Linked List”, *International Journal Of Engineering Research & Management Technology*, Volume 1, Issue-5, pp. 44-48, 2014.
- [28] Petra G., Miroslav B., “Analysis of B-tree data structure and its usage in computer forensics”, *Conference: Central European Conference on Information and Intelligent Systems*, 2010.
- [29] Aouiche K., Darmont J., “Data Mining based materialised view and index selection in data warehouses”, *Journal of Intelligent Information Systems* Vol. 33, DOI: 10.1007/s10844-009-0080-0, 2007.
- [30] J. Han, J. Pei, G. Dong, and K. Wang.: Efficient computation of Iceberg cubes with complex measures, *ACM SIGMOD Rec.*, vol. 30, no. 2, pp. 1–12, 2005.
- [31] Jiawei H., Micheline K., “Data Mining Concepts and Techniques”, 2<sup>nd</sup> Edition, Ed. Morgan Kaufmann Publishers, pp. 230, 2006.

# The Effects of Speed and Altitude on Wireless Air Pollution Measurements using Hexacopter Drone

Rami Noori<sup>1</sup>

Faculty of Information Science and Technology (FTSM)  
Universiti Kebangsaan Malaysia (UKM)<sup>1</sup>

Dahlila Putri Dahnil<sup>2</sup>

Center for Software Technology and Management  
(SOFTAM) Faculty of Information Science and Technology  
(FTSM), Universiti Kebangsaan Malaysia (UKM)<sup>2</sup>

**Abstract**—Air pollution has a severe impact on human beings and one of the top risks facing human health. The data collection near pollution sources is difficult to obtain due to obstacles such as industrial and rural areas, where sensing usually fails to give enough information about the air quality. Unmanned Aerial Vehicles (UAVs) equipped with different sensors offer new approaches and opportunity to air pollution and atmospheric studies. Despite that, there are new challenges that emerged with using UAVs, in particular, the effect of wind-generated from UAVs propellers rotation on the efficiency and ability to sense and measure gas concentrations. The results of gas measurement are affected by the propellers rotation and the wind resistance. Thus, the effect of changing UAV speed and altitude on the gas measurement both vertically and horizontally need to be performed. The aims of this paper is to propose a new mobile-wireless air pollution system composed of UAV equipped with low-cost sensors using LoRa transmission. The proposed system is evaluated by studying the effect of changing altitude and speed on the measured gas concentrations CO, LPG, H<sub>2</sub>, and smoke when flying in horizontally and vertically directions. The results showed that our system is capable of measuring CO, LPG, H<sub>2</sub>, and smoke in the vertical mode in both hovering and deploying scenarios. While in horizontal mode the results showed that system can detect and measure gas concentrations at speeds less than or equal to 6 m/s. While at high speed of 8 and 10 m/s there will be an impact on its performance and accuracy to detect the targeted gases. Also, the results showed that the LoRa shield and Radio transmitter AT9S can successfully transmit up to 800 m horizontally and 400 feet vertically.

**Keywords**—Unmanned Aerial Vehicles (UAVs); low-cost sensors; air pollution; LoRa; air quality; radio transmitter; atmosphere

## I. INTRODUCTION

The continuous changes in ambient air that are associated with both natural and anthropogenic emissions (such as aerosols or gaseous pollutants) has a significant effect on air quality and consequently on human health [1]. Recent studies on the “Global Burden of Disease,” identified air pollution as one of the top 10 risks facing human beings. Many cities have consistently violated the recommended concentration ranges of air pollutants and the direct implication of this violation is several air pollutant-related premature deaths [2], [3].

Haze is one of the most common phenomena associated to air contamination which happens almost every year within the past decades in Southeast Asia including Malaysia due to the forest fires particularly in Sumatra and Kalimantan, Indonesia

[4]. According to [5], Malaysian economic loss was around 4,471 USD per haze day and the study predicted in next 20 years, the losses will be from 1 million USD to 1.6 million USD per year. The biomass burning causes serious air pollution and increases the haze condition [6], also considered as a major source of carbon monoxide (CO) [7]. The CO is one of the most hazardous air pollutants that causes severe health problems and more than 400 co-related deaths reported yearly in the U.S. [8].

However, the conventional air contamination monitoring systems (ACMS), despite being efficient and reliable for measuring a wide range of pollutants, still have some drawbacks in terms of their size, weight, and cost. One of the drawbacks is that the monitoring stations may not be able to cover all locations. Hence, there are several un-monitored locations need to be estimated [9]. These drawbacks necessitate the wide deployment of air monitoring stations [10]. Furthermore, the analysis of the data and its deployment in conventional ACMS is too slow. So, high spatial-temporal resolution with a real-time system is fundamental because of the restricted information accessibility and non-versatility of the traditional ACMS [9].

Low-cost sensors change the traditional way of measuring air pollutants [2], they can be used with higher spatial & temporal resolutions [11]. However, the detecting range of the low-cost sensors is much lower than the traditional observing tools due to the direct interaction of the sensor surface with a small volume of the chemical compounds. Hence, a stationary sensor network is not applicable in most cases from both economical and deployment-related perspectives [12]. As a result of that, many researches have been trying to utilize mobile objects in both ground and air [13]. Mobility can fill the gap between the traditional monitoring tools and the air quality measurement models, especially in the areas without monitoring stations where the data about pollution is achieved via air quality modeling or predictions [9].

Although the mobile ground objects offer some advantages, they still have limitations, such as in industrial and rural areas where sensing usually fails to give enough information to acquire sensible measurements with the required granularity level [13]. These concerns necessitate the use of UAVs to monitor such areas [14]. UAV technology has gradually gained popularity over the years [15], and a vast amount of information has been collected for air pollution that spread across scientific and non-scientific databases [16]. It

plays an essential role in many other fields such as food security that help to combat food insecurity due to the COVID-19 pandemic [17]. However, UAVs provide new challenges, in terms of payload capacity, power consumption and stability [10]. Also, the limitation on sensors selection which needs to be suitable and small enough to mount them on board UAVs and that may lead to select sensors with less sensitivity and selectivity [1]. Furthermore, data captured by low-cost sensors must be critically evaluated due to their heavy dependence on numerous factors, especially the impact of wind generated by the UAVs propellers rotation [10]. Also [18] shows that using UAVs for air pollution measurement can be only effective if the location point of the air sensor has been optimized. Moreover, the accurate results of gas measurement rely highly on the contribution of propellers rotation and wind resistance that need to be assessed.

In this study, we aimed to develop a mobile-wireless air pollution system to measure gases concentration such as CO, LPG, H<sub>2</sub> and smoke by developing a detection system based on low-cost sensors (subsystem-1), and an affordable and open-source UAV (subsystem-2). The system is evaluated on the effect of changing altitude and speed on gas concentration when flying horizontally and vertically. The experiment conducted is also to assess the communication range of the system stations using LoRa Shield and Radio transmitter AT9S.

## II. RELATED WORKS

The air pollution monitoring systems and low-cost sensors are presented in this section. The existing works are grouped into three categories based on the carriers of the low-cost sensor nodes. The disadvantages of each category are highlighted in each subsection. This section is organized as follows: i) the static low-cost sensors used for air pollution measurement; ii) low-cost sensors with mobile-ground objects; and iii) low-cost sensors with mobile flying objects.

### A. Static Low-Cost Sensors

The study by [19], used a micro controller-based toxic gas to detect and alert the presence of hazardous gases like LPG and propane emission. When these gases exceed their safety level, an alarm is generated and send an SMS message through GSM modem to an authorized person. Their system consists of PIC 16F877 as a Micro controller and MQ-2 and MQ-6 gas sensors. The analog signal sensed from the sensors represents the concentration of the hazardous gases and converted to digital signals through ADC in their micro-controller. The system provides fast response with accurate results which leads to faster diffusion in emergency cases. The limitations of this system is its only applicable for indoor air quality monitoring and the sensor nodes are constantly on sleep mode as there is no updating of data in the same location continuously.

The author in [20] developed a WSN for outdoor monitoring of air pollution, the designed prototype was tested on a real-time basis. The sensors capture O<sub>3</sub>, NO<sub>2</sub>, CO, H<sub>2</sub>S pollution data while the sensed data is transmitted via GPRS to the server. Solar panel is used to provide power to the stationary sensor nodes. Customized mobile and web apps are

provided for making the air pollution data available to the public. The challenges facing this system is in terms of temporal-resolution when the number of deployed stations increases. The increment override the low spatial-resolution that leads to congestion in a single cellular base station that serves a large number of monitoring stations.

Another research [11] developed a low-cost air quality system known as DiracSense to measure gas pollutants that are indexed by Malaysian ambient air quality standards. They used CO-AF, OX-AF, and NO<sub>2</sub>-AF sensors that are manufactured by Alphasense. The electrochemical sensors measure CO, O<sub>3</sub>, and NO<sub>2</sub> gases. PTU300 sensor measures temperature, pressure, and relative humidity. Also, they used Raspberry Pi as a micro-controller. DiracSense collects, analyzes, and shares air quality data using wireless communication. An android mobile phone application was used to display the data of air quality. They calibrated the sensors by using laboratory and field test experiments. They used an adaptive neuro-fuzzy inference system (ANFIS) as the calibration model and used a multi-layer perceptron (MLP) to assess the capability of the ANFIS as the calibration model. The results showed that the ANFIS model is promising as a calibration tool due to its ability to enhance the accuracy and performance of the low-cost electrochemical sensors.

The study by [21] presented an air pollution and monitoring model that consists of Bluetooth micro controller for transferring the values of the sensors from ADC to a server, MOS (MQ-7, MQ-5) sensors, and server to save all the data collected. They also presented the ID3 algorithm to calculate the sensor values saved on the server based on probability. They proved that the model can predict the air pollution in some areas. Research [22] developed an air pollution measurement and prediction system for measuring CO and H<sub>2</sub>. Their system consists of Beagle bone Black as a micro controller, MOS (MQ-7, MQ-11) sensors, and GPS module for tracking the concentration of pollution. The data collected from the sensors are uploaded on Azure Cloud via Python SQL. They applied a machine learning service on the data saved in cloud for predicting the pollution. This study shows that the cloud data can be used for prediction of air pollution. The challenges facing these two systems is in terms of spatial-resolution that the number of node sensors should also be increased to cover a larger spatial resolution.

### B. Low Cost Sensors in Ground-based Mobile Objects

Many researchers utilize mobile ground objects [13] for air pollution measurement to overcome drawbacks of facing static systems in terms of low spatial-temporal resolution, deployment of sensor nodes, maintenance and calibration obstacles.

In ground mobile objects, [23] presented a real-time WSN-based pollution monitoring. The sensors sense the concentration of CO, CO<sub>2</sub>, and O<sub>2</sub> gases deployed on sensor nodes that have been calibrated. The project implementation was done in the industrial area of Hyderabad city. The study deployed a multi-hop data collection algorithm while the collected air pollutant data from the designed test beds are made available onto the internet through dedicated web interface. The developed system is capable of obtaining the



fine-grain pollution data on a real-time basis. The challenges facing this type of systems are in terms of uncontrolled or semi-controlled mobile and redundant sampling issues.

The author in [24] developed a system called Air-sense for air quality monitoring in both outdoor and indoor. Their system was designed with 4 layers. The first layer for collecting the data through the people carrying a portable Air Quality Monitoring Device (AQMD); the second for treating and formatting the data collected and transmitted through the first layer; third layer is responsible for communication between the cloud server and the smart phone and the fourth layer is responsible for analyzing and storing the data. Their AQMD consists of Arduino Pro Mini board as a micro controller, MOS (MQ-7, MQ-135) for detecting CO, and monitoring air quality and Bluetooth module HC-05 for transmitting the data from AQMD to smart phones. The study concluded that this system will encourage the citizens to be part in the crowd sensing action, which could be a backbone of any smart city. The challenges facing this type of systems in terms of spatial-to-temporal resolution trade-off (higher spatial coverage at the expense of lower temporal resolution), also low data accuracy and reliability.

The author in [25] presented a low-cost system for air quality monitoring using a vehicular sensor network. This system processes the data collected by sensors located on public vehicles. The system consists of Arduino as a micro-controller connected to MOS sensors for detecting (NO<sub>2</sub>, CO<sub>2</sub>, CO, and Ozone), as well as to measure air quality. The acquired data is transferred to a server on Raspberry Pi board through Xbee-based Access Points installed on the road.

The author in [26], suggested a low-cost portable system for air pollution monitoring by using IoT to create awareness to the public about the air quality, enabling them to make better choices regarding traveling routes or purchasing of houses in a better area.

### C. Low Cost Sensors in Flying-based Mobile Objects

Researchers have been trying to utilize mobile flying objects such as UAVs to overcome the drawbacks facing mobile-ground systems in terms of low spatial resolution. Also, it's proven that air pollution changes abruptly even at a small relative distance both horizontally and vertically [13].

In flying mobile objects, [10] presented an Air Pollutants Monitoring Using UAVs (ARIA) project, with the aim of finding a toll to measure air quality vertically at different heights. They presented an overview about their project and the low-cost air pollution measurement without the experimental results. For air pollution monitoring, they presented a system that consists of Raspberry Pi 3 as micro controller, Alpha-sense Gas sensors, and Particulate and Volatile Organic Compounds (VOCs) sensors. They suggested placing the sensors inside the drone to avoid the airstream generated from the propeller rotation, with mention of the drawbacks of this configuration which can lead to losses of signal along the wires. They stored data on an On-Board storage control unit connected with the measurement system and when the drone land off, they can download the data from the board.

The author in [8] suggested an airborne WSN system called AIRWISE for automated measurement and monitoring of ambient air pollution. The system comprised of an unmanned aerial vehicles (UAVs) and a pollution-aware wireless sensor (PAWS) network. The system was designed for monitoring of ambient air pollutants in 3D spaces without any form of human intervention. They suggested two schemes for autonomous monitoring of a 3D area of interest; their PAWS consist of Wasp mote node as a micro-controller, GPS module and integrated sensors for measuring gases concentration and air quality. The targeted gases were O<sub>2</sub>, NH<sub>3</sub>, CH<sub>4</sub>, and CO<sub>2</sub>.

The author in [27] presents a data acquisition system and its coupling into a UAV to facilitate air pollutants monitoring. The collected data is transmitted via RF to the ground station for processing. The results are displayed on a web page that can be accessed using any mobile device or computer. Their monitoring unit consists of UAV S500 quadcopter with the pix-hawk flight computer, Arduino, Shield for XBee antenna connection, XBee PRO S2B Antenna, and DAQ with sensors to measure air quality. The sensors used in this system are MQ7 to measure Carbon monoxide, MQ8 for Hydrogen, MQ131 for Ozone, and MQ135 for Carbon dioxide. The ground station consists of a PC, Arduino, and Receiver antenna 3DR used for monitoring the behavior of variables. The flying was stable when the load was located in the center of its frame, but the flight time decreased from 13 to 10 minutes. The results show that there is no lost in communication or interference within the range of 203 meters.

### III. PROPOSED SYSTEM

The system developed and evaluated in this study consists of two subsystems: the air pollution detection system (subsystem-1) and the UAV (subsystem-2). Fig. 1 illustrates the methodology in developing the system.

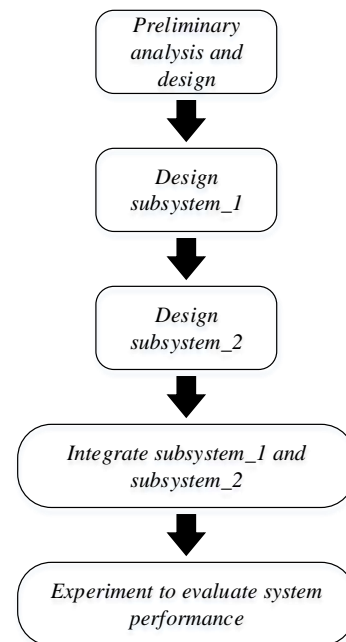


Fig. 1. Proposed Methodology.

### A. Air Pollution Detection System (Subsystem\_1)

The proposed system consists of two subsystems. For subsystem\_1, we used Arduino UNO over the other devices due to its power consumption and simple connectivity [1], LoRa shield is used for data transmission since it offers a long-range transmission compared to other wireless transmission techniques as shown by [28-29], and MQ sensors due to their low cost, compatible with Arduino, and their ability to sense various chemical gases.

The subsystem\_1 composed of two stations, the first station is the detection station (DS\_subsystem\_1) to detect and measure targeted gas concentrations. The second station is the monitoring station (MS\_subsystem\_1) to receive and display the targeted gas concentration measurements.

The DS\_subsystem\_1 composed of MQ2 sensor to sense smoke and CO, MQ6 to sense LPG, and MQ8 to sense hydrogen. The sensors were coupled to the Arduino UNO R3 micro-controller and RFM LoRa Shield to facilitate wireless communication.

The MS\_subsystem\_1 composed of Arduino UNO R3 board with RFM LoRa Shield, connected to a PC with Arduino IDE (Integrated Development Environment) system via a USB port for real-time visual monitoring.

### B. Configuring the MOS sensor (MQ2, MQ6, and MQ8)

As an analog sensor, the response values of the MOS gas sensor are the outputs of the analog-to-digital conversion (ADC). Based on the sensors datasheet, the resistance of the sensor (RS) will be calculate as follows [30]:

$$R_s = \frac{V_c - V_{out}}{V_{out}} \times RL \quad (1)$$

Where

$$V_{out} = \frac{ADC * V_c}{1023} \quad (2)$$

where  $V_c$  is a micro-controller board voltage,  $V_{out}$  is the output voltage of the sensor in the sample space. RL is the sensor load resistance, ADC is analog-to-digital value conversion. And based on the provided chart in the MQ datasheet,  $R_s$  in clean air is constant under standard temperature & humidity. The  $R_s/R_0$  ratio in clean air is 9.8 as described in the datasheet. After calibrating the sensor by placing it in clean air and getting the  $R_0$  value by dividing it with the  $R_s/R_0$  value in clean air, the targeted gas can be sensed using the  $R_s/R_0$  ratio. To calculate the concentration of the targeted gas in ppm, the datasheets provided the sensitivity characteristics of each sensor with ppm (gas concentration) as x-axis and RS/RO as y-axis. So the curve of the particular gas in the sensitivity characteristics of each sensor is used to calculate the slope [31] as follows:

$$Slope = \frac{(Y2 - Y1)}{(X2 - X1)} \quad (3)$$

Where X1 is the logarithm of the first point for targeted gas curve in x-axis (PPM), X2 is the logarithm of the last point for targeted gas curve in x-axis(PPM), Y1 is the logarithm of the first point for targeted gas curve in y-axis(RS/RO), and Y2 is the logarithm of the last point for aimed gas curve in y-axis(RS/RO). The slope of each MOS sensor (MQ2, MQ6, and MQ8) was calculated separately.

Then, we calculated the concentration of the targeted gas in ppm from:

$$Gas\ Concentration = \frac{\left(\frac{R_s}{R_0} - Y1\right)}{(Slope + X1)} \quad (4)$$

### C. Unmanned Aerial Vehicle (Subsystem\_2)

Drones are built for mobile and data-gathering capability. It can be recognized as a flying robot, which can be flown autonomously without any human interaction via software-controlled flight plans or can be remotely controlled [32]. The flying mode makes it easier to reach nodes in a timely manner and also makes it easy to hover and collect data at specific nodes [10]. The Hexacopter drone is chosen in our design as it has a large payload capacity of more than 2 kilogram payload; it also has better manoeuvrability and in-flight stability in comparison to quadrotors. This made them ideal for UAVs and studies on air quality where there is a need to carry various sensors and sustain a fixed in-flight position.









Our subsystem\_2 consists of two stations, Hexacopter drone (HD\_subsystem\_2) and UAV control station (CS\_subsystem\_2) as shown in Fig. 2. The purpose of the CS\_subsystem\_2 is for communicating wirelessly with the HS\_subsystem\_2 via a Radio Controller (RC) or a PC with a mission planner application connected to a Telemetry Radio. The UAV kits in this project are customized based on their functionalities and capabilities as shown in Table I. The UAV are self-assembled from scratch by studying all parts and communicating from the supplier of drone experts.



(a) Hexacopter Drone. (b) UAV Control Station.

Fig. 2. UAV Stations.

TABLE I. UAV KITS

| Model  | Frame   | Description  |
|--|---|--|
| DJI F550 frame                                 |    | Hexacopter drone is strong enough to carry loads. They are built with an increased number of motors; hence, they are relatively stable in-flight, and safer with 6 motors at 120° apart, whereby flight can be maintained even if one fails  |
| Radiolink Pixhawk                              |    | PIXHAWK is an open source flight controller made by Radio link; it has less interference of inner components, less noise, safety while flying, more accurate compass, and supports various flight modes.   |
| Radiolink M8N GPS SE100                        |    | Radiolink M8N GPS is compatible with all the open source flight controllers its accuracy is 50-centimeter position and positions 20 satellites within 6 seconds at open ground. It is capable of valley station-keeping: Max height is 50000 m and Max speed is 515m/s.                      |
| Radiolink AT9S Transmitter with R9DS Receiver  |   | This is an affordable radio system with a range of 900 m on the ground and 1500 m in air (environment-dependent). It includes a 2.8 inches LCD screen that can show a real-time data telemetry, such as GPS, SPEED, voltage etc. and it is suitable for all multicopters, boats, and cars    |
| FPV Radio Telemetry                            |  | used to receive a Real-time data about drone while its flying into the mission planner software in pc.   |
| (2212 900kv) Brushless Motor and "9" Propeller |  | Fig. 1. The feature of the motor is light weight, extreme torque, low current and low temperature  |
| Hobbywing 20A ESC                              |  | The electronic speed controller for translating a pilot's controls into specific instructions to be transmitted to the motors for movement control. Hobbywing 20A ESC is designed especially for multi-copters; It is compatible with various flight-control systems and lightly weight 14 g |
| Imax B6AC Charger                              |  | This is a versatile and fast battery charger, balancer, and discharger. It is compatible with different types of batteries. It is equipped with an internal independent Li-battery balancer for ensuring balanced charging & discharging of 2-6 cells.                                       |

D. Integration of Subsystem 1 and Subsystem 2

The best location to place the DS\_subsystem\_1 is by considering the stability of the drone at the bottom of Hexacopter (HS\_subsystem\_2) as shown in Fig. 3. In order to achieve drone stability during flying and hovering mode, as well as to decrease the effect of the wind generated by the propellers, we placed it as far as possible from the propellers (28 cm from propellers) with an extension to 5 cm from the bottom of the Hexacopter drone. This is also to ensure safety of the DS\_subsystem\_1 by using landing skid in case of fall-down or land-off. By integrating both subsystems we achieved our first goal and developed a mobile wireless air pollution system which consist of two stations: mobile station (DS\_subsystem\_1 and HD\_subsystem\_2) and ground station (MS\_subsystem\_1 and CS\_subsystem\_2).



Fig. 3. Mobile Station.

IV. SYSTEM EVALUATION

In this study, three types of experiments were conducted to evaluate the effect of changing altitude and speed on the system. It detected and measured identified gases and assessed the communication range of the system stations. The experiments were performed as follows: i) To measure the range of wireless communication between the system stations (mobile and ground stations). ii) To evaluate the performance of the mobile station to detect and measure identified gas in vertical mode. iii) To evaluate the performance of the mobile station to detect and measure identified gas in horizontal mode.

A. Range of Wireless Communication

We tested the range of wireless communication of the system by flying the mobile station in an open field to the allowable maximum range according to the safety guidelines. The experiment was conducted to:

- 1) Examine the horizontal and vertical range of the mobile station to communicate with the ground station manually by Radio transmitter AT9S (RC) and autopilot by Radio telemetry with mission planner software.

2) Examine the horizontal and vertical range of the mobile station to transmit (CO, H2, LPG, and smoke) readings by LoRa Shield.

### B. Vertical Evaluation

The experiments were conducted to evaluate the effect of changing altitude and speed on the mobile station performance to detect and measure targeted gas, with the following scenarios:

1) Hovering mode, the first readings were taken at one identified spot at ground level. After recording a number of packets, we raised the mobile station to 5 m. Then the RC was changed into hovering mode to stabilize it at its position and altitude. Then the level was increased to 10m and the measurements are recorded. The same steps are repeated on subsequent level as shown in Fig. 4.

2) Deploying mode, the initial same spot was used as described on previous experiment. The mobile station takes off from the ground to 20 m altitude at different speeds ranging (2-10) m/s with a 2 m/s step size as shown in Fig. 5.

### C. Horizontal Evaluation

Experiments were carried out to evaluate the effect of changing speed on the mobile station performance to detect and measure identified gas. The experiments were set up by a flight plan installed to mobile station through the ground station (using mission planner and Radio telemetry). The flight plan was conducted by using 4 way-points as shown in Fig. 6. The way-points were as follows:

1) The first point was located above the ground station, for the mobile station to take off. After the mobile station took off at 15 m altitude, it flew at speed that was set through (WPNAV\_SPEED in mission planner software) towards the second point.

2) The second point was located at distance of 150 m from the ground station. When the mobile station reaches this point, the drone hovered for 10 seconds. The gas readings received during these 10 seconds will be ignored in our measurements.

3) The third point had the same location as the second point. At this point, the "Do-change-speed" command was used; to change the mobile station speed to 2 m/s, then the mobile station flew towards the fourth point. The mobile station speed changed for the comparison purpose.

4) The fourth point was located at the ground station, where the mobile station landed.

The purpose of using this method was to ensure that the mobile station will fly at the same path in nearly same time at two different speeds, which were:

a) From first to second point, we tested different speeds ranging (2-10) m/s with a 2 m/s step size, to evaluate their impact on the mobile station performance to detect and measure targeted gas.

b) From third to fourth point, we used the fixed speed in the experiment (2 m/s).

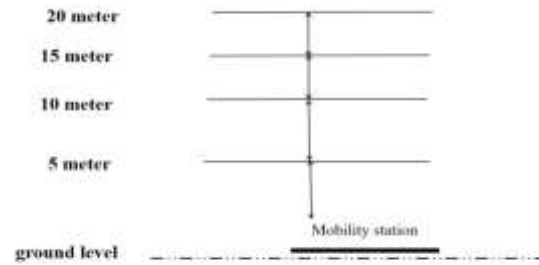


Fig. 4. Hovering Scenario.

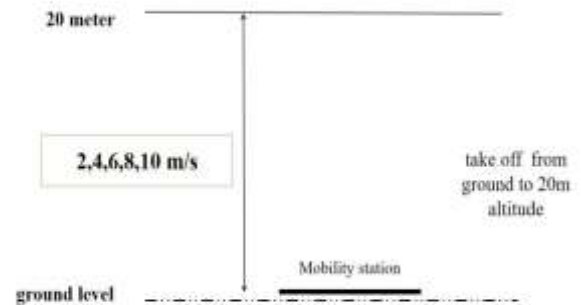


Fig. 5. Deployment Scenario.

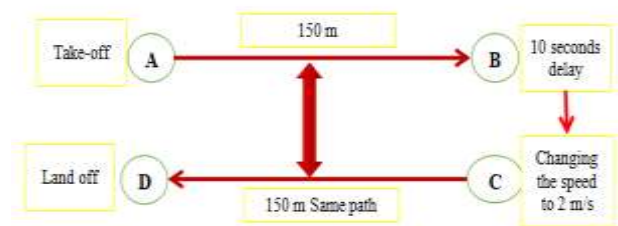


Fig. 6. Flight Plan Scenario.

## V. RESULTS AND DISCUSSION

The results of testing the range of wireless communication for the system stations are shown in Table II. The GPS used in the mobile station was M8N GPS SE100 and each position is accurate within (1 to 2.5) m, and its speed is within (0.1 to 0.5) m/s.

TABLE II. THE WIRELESS COMMUNICATION RANGE BETWEEN SYSTEM STATIONS

| Communication method                 | Horizontal range   | Vertical range      |
|--------------------------------------|--------------------|---------------------|
| Radio transmitter AT9S (RC)          | 800 m horizontally | 400 feet vertically |
| Radio telemetry with mission planner | 100 m horizontally | 100 feet vertically |
| LoRa Shield                          | 800 m horizontally | 400 feet vertically |

The average of 15 readings were taken for each gas at each altitude in the experiments for hovering mode for vertical evaluation as shown in Fig. 7. The average of 10 readings were taken for each gas at each speed for the deploying mode in the experiments as shown in Fig. 8.

Finally, for the results of the horizontal evaluation, we took the average readings for each gas in each speed experiment as shown in Fig. 9.

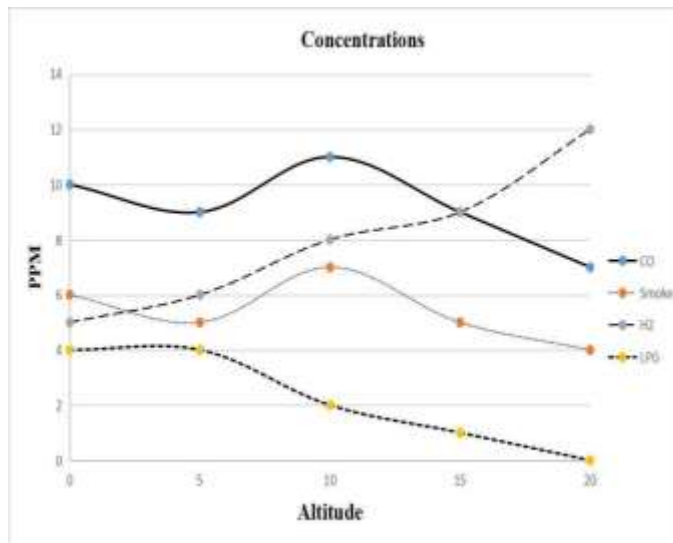


Fig. 7. Gases Concentration in Hovering Mode.

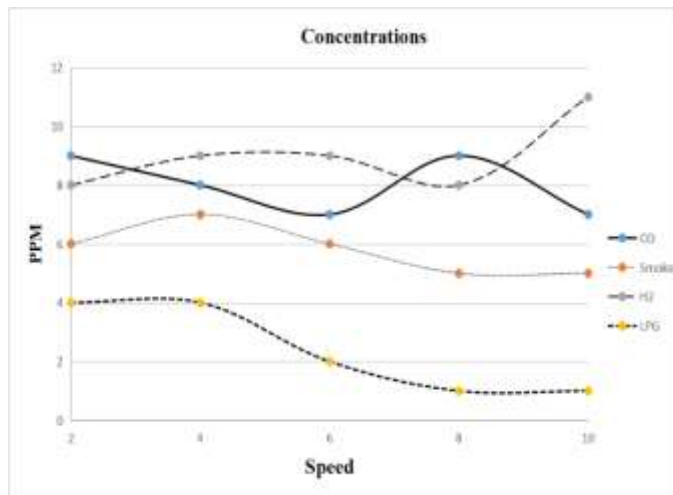


Fig. 8. Gases Concentration in Deploying Mode.

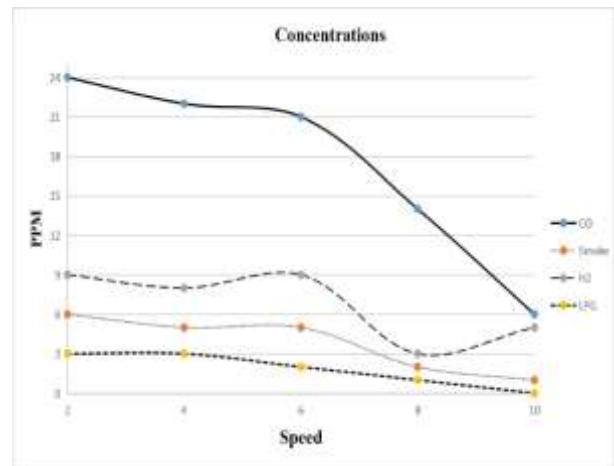


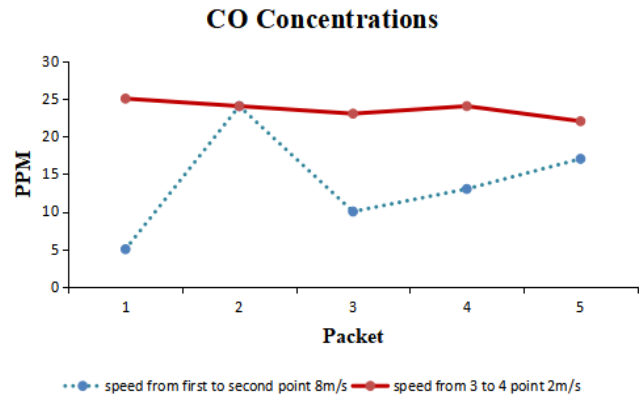
Fig. 9. Gases Concentration in Horizontal Mode.

Fig. 9 shows that there is a big change in average readings for each gas with respect to speed, especially for the CO gas at 8 and 10 m/s. We flew the mobile station at the same path at two different speeds. We compared the packets received at both speeds to ensure that changes in gases readings were due to changes in speed rather than the actual gas concentrations in that path. The results were as follows:

1) In the first case, the mobile station was taken from first to the second point at the following speed of 2, 4, 6 m/s respectively. Then, from third to the fourth point at 2 m/s. We compared the packets received at each point. The results indicated that both gas concentrations were almost in the same range.

2) In the second case, the mobile station was taken from first to the second point at 8 m/s. Then, from third to the fourth point at 2 m/s speed. We compared the packets received at each point as shown in Fig. 10. The results indicated a huge difference between packets received at each speed.

3) In the third case, the mobile station was taken from first to the second point at 10 m/s speed. Then from third to the fourth point at 2 m/s speed. We compared the packets received at each point as shown in Fig. 11. The results have indicated that there was a significant gap in gas concentrations.



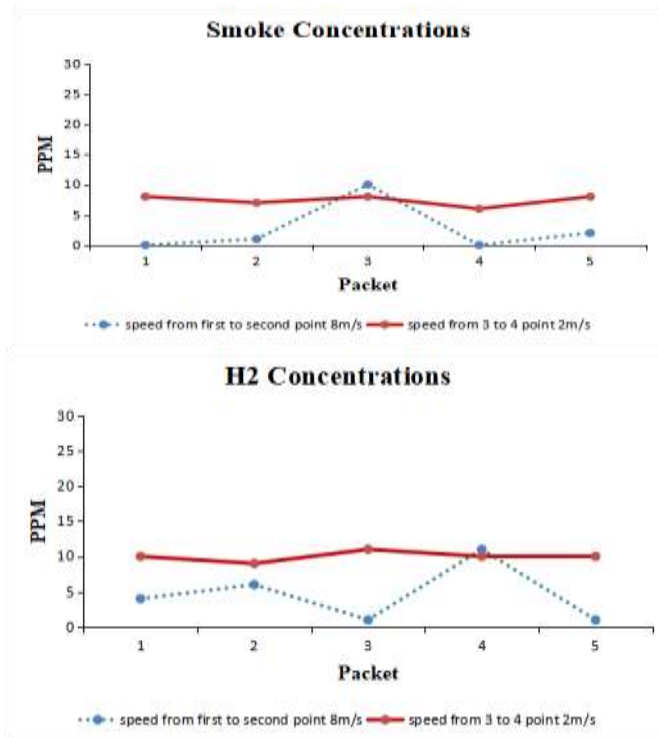


Fig. 10. Gases Concentration at 8 m/s and 2 m/s.

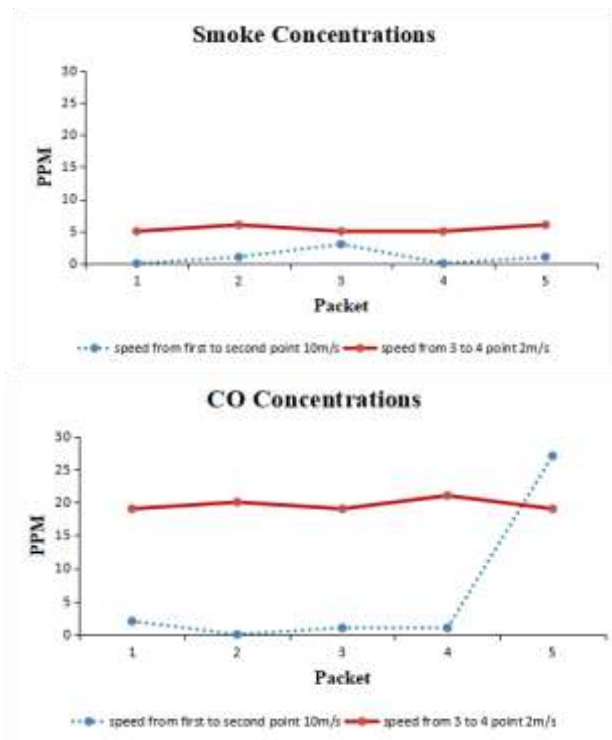


Fig. 11. Gases Concentration at 10 m/s and 2 m/s.

## VI. CONCLUSION AND FUTURE WORK

In conclusion our system is effective to measure (CO, LPG, H<sub>2</sub>, and Smoke) in the vertical mode in both hovering and deploying scenarios. There was no significant impact due to changes in altitude and speed to the system performance on

detecting and measuring the gases concentration. In horizontal mode, the results showed that the system is effective to detect and measure gas concentrations at speeds less than or equal to 6 m/s. While at high speed (8 and 10 m/s) there is an impact on its performance and accuracy to detect the gases. The ineffectiveness of the mobile station performance at high speed is due to the wind generated from the rotation of the propellers in horizontal-mode. The wind generated affects the air samples around the sensors. The low-range detection for the low-cost sensors needs direct exposure to the gases. The experiments indicated that LoRa shield and RC can successfully transmit up to 800 m horizontally and 400 feet vertically.

Our future work will focus on measuring the wind generated from the propeller's rotation at each speed especially at high speed (8 and 10) m/s and its effect to the air quality readings.

## ACKNOWLEDGMENT

The research is financially supported by *Cabaran Perdana* Research Grant Scheme [Grant No.: DCP-2018-001/2], University Kebangsaan Malaysia (UKM). [www.ftsm.ukm.my/softam](http://www.ftsm.ukm.my/softam), Faculty of Information Science and Technology.

## REFERENCES

- [1] T. Villa, F. Gonzalez, B. Miljievic, Z. Ristovski, and L. Morawska, "An overview of small unmanned aerial vehicles for air quality measurements: Present applications and future Prospectives", *Sensors*, vol. 16, no. 7, 1072, 2016.
- [2] P. Kumar, L. Morawska, C. Martani, G. Biskos, M. Neophytou, S. Sabatino, M. Bell, L. Norford, and R. Britter, "The rise of low-cost sensing for managing air pollution in cities", *Environment International*, vol. 75, pp. 199-205, February 2015.
- [3] S. Duangsuwan, and P. Jamjareekulgarn, "Development of Drone Real-time Air Pollution Monitoring for Mobile Smart Sensing in Areas with Poor Accessibility", *Sensors and Materials*, vol. 23, no. 7, pp. 511-520, January 2020.
- [4] M. Sahani, N. Zainon, W. Mahiyuddin, M. Latif, R. Hod, M. Khan, N. Tahir, and C. Chan, "A case-crossover analysis of forest fire haze events and mortality in Malaysia", *Atmospheric Environment*, vol. 96, pp. 257-265, October 2014.
- [5] J. Montgomery, C. Reynolds, S. Rogak, and S. Green, "Financial implications of modifications to building filtration systems", *Building and Environment*, vol. 85, pp. 17-28, February 2015.
- [6] H. Zhao, X. Zhang, S. Zhang, W. Chen, D. Tong, and A. Xiu, "Effects of Agricultural Biomass Burning on Regional Haze in China: A Review", *Atmosphere*, vol. 8, no. 5, May 2017.
- [7] Y. Sawa, H. Matsueda, Y. Tsutsumi, J. Jensen, H. Inoue, and Y. Makino, "Tropospheric carbon monoxide and hydrogen measurements over Kalimantan in Indonesia and northern Australia during October 1997", *Geophysical research letters*, vol. 26, no. 10, pp. 1389-1392, May 1999.
- [8] E. Orestis, and J. Rolim, "An Airborne Wireless Sensor Network for Ambient Air Pollution Monitoring", *SENSORNETS*, pp. 231-239, September 2015.
- [9] W. Yi, K. Lo, T. Mak, K. Leung, Y. Leung, and M. Meng, "A Survey of Wireless Sensor Network Based Air Pollution Monitoring Systems", *Sensors*, vol. 15, no. 12, pp. 31392-3142, December 2015.
- [10] G. Bolla, M. Casagrande, A. Comazzetto, R. Moro, M. Destro, E. Fantin, G. Colombatti, A. Aboudan, and E. Lorenzini, "ARIA: Air Pollutants Monitoring Using UAVs", In 2018 5th IEEE International Workshop on Metrology for AeroSpace (MetroAeroSpace), pp. 225-229, Jun 2018.
- [11] K. Alhasa, M. Nadzir, P. Olalekan, M. Latif, Y. Yusup, M. Faruque, F. Ahamad, H. Hamid, K. Aiyub, S. Ali, M. Khan, A. Samah, I. Yusuff, M.

- Othman, T. Hassim, and N. Ezani, "Calibration Model of a Low-Cost Air Quality Sensor Using an Adaptive Neuro-Fuzzy Inference System", *Sensors*, vol. 18, no. 12, December 2018.
- [12] M. Rossi, and D. Brunelli, "Autonomous Gas Detection and Mapping with Unmanned Aerial Vehicles", *IEEE Transactions on Instrumentation and Measurement*, vol. 65, no. 4, pp. 765-775, December 2015.
- [13] O. Alvear, N. Zema, E. Natalizio, and C. Calafate, "Using UAV-Based Systems to Monitor Air Pollution in Areas with Poor Accessibility", *Journal of Advanced Transportation*, August 2017.
- [14] M. Dunbabin, and L. Marques, "Robots for Environmental Monitoring: Significant Advancements and Applications", *IEEE Robotics & Automation Magazine*, vol. 19, no. 1, pp. 24-39, February 2012.
- [15] G. Rohi, O. Ejofodomi, and G. Ofualagba, "Autonomous monitoring, analysis, and countering of air pollution using environmental drones", *Heliyon*, vol. 6, no. 1, e03252, January 2020.
- [16] J. Burgués, and S. Marco, "Environmental chemical sensing using small drones: A review", *Science of The Total Environment*, vol. 748, 141172, December 2020.
- [17] U. Panday, A. Pratihast, J. Aryal, and R. Kayastha, "A Review on Drone-Based Data Solutions for Cereal Crops", *Drones*, vol. 4, no. 3, September 2020.
- [18] T. Villa, F. Salimi, K. Morton, L. Morawska, and F. Gonzalez, "Development and Validation of a UAV Based System for Air Pollution Measurements", *Sensors*, vol. 16, no. 12, 2202, December 2016.
- [19] V. Ramya, and B. Palaniappan, "Embedded system for Hazardous Gas detection and Alerting", *International Journal of Distributed and Parallel Systems (IJDPS)*, vol. 3, no. 3, pp. 287-300, May 2012.
- [20] A. Kadri, E. Yaacoub, M. Mushtaha, and A. Abu-Dayya, "Wireless sensor network for real-time air pollution monitoring", In 2013 1st International Conference on Communications, Signal Processing, and their Applications (ICCSPA), pp. 1-5, February 2013.
- [21] S. Raipure, and D. Mehetre, "Wireless sensor network-based pollution monitoring system in metropolitan cities", *International Conference on Communications and Signal Processing (ICCSP)*, pp. 1835-1838, April 2015.
- [22] N. Desai, and J. Alex, "IoT based air pollution monitoring and predictor system on Beagle bone black", *International Conference on Nextgen Electronic Technologies: Silicon to Software (ICNETS2)*, pp. 367-370, March 2017.
- [23] M. Pavani, and P. Rao, "Monitoring Real-Time Urban Carbon Monoxide (CO) Emissions Using Wireless Sensor Networks", *International Conference on Information and Communication Technology for Intelligent Systems (ICTIS)*, vol. 2, pp. 290-297. Springer, March 2017.
- [24] J. Dutta, C. Chowdhury, S. Roy, A. Middy, and F. Gazi, "Towards Smart City: Sensing Air Quality in City based on Opportunistic Crowdsensing", In *Proceedings of the 18th international conference on distributed computing and networking*, no. 42, pp. 1-6, January 2017.
- [25] G. Re, D. Peri, and S. Vassallo, "Urban Air Quality Monitoring Using Vehicular Sensor Networks", In *Advances onto the Internet of Things*, pp. 311-323, 2014.
- [26] T. Ahuja, V. Jain, and S. Gupta, "Smart Pollution Monitoring for Instituting Aware Traveling", *International Journal of Computer Applications*, vol. 145, no. 9, pp. 0975-8887, July 2016.
- [27] J. Vega, E. Varela, N. Romero, C. Santos, J. Cuevas, and D. Gorham, "Internet of Things (IoT) for Monitoring Air Pollutants with an Unmanned Aerial Vehicle (UAV) in a Smart City", *Smart Technology*, pp. 108-120. Springer, 2018.
- [28] N. Husein, A. Rahman, and D. Dahnil, "Evaluation of LoRa-based Air Pollution Monitoring System", *(IJACSA) International Journal of Advanced Computer Science and Applications*, vol. 10, no. 7, 2019.
- [29] M. Hossinuzzaman, and D. Dahnil, "Enhancement of Packet Delivery Ratio during Rain Attenuation for LoRa Technology", *(IJACSA) International Journal of Advanced Computer Science and Applications*, vol. 10, no. 10, 2019.
- [30] D. Wijaya, R. Sarno, and E. Zulaika, "Gas concentration analysis of resistive gas sensor array", *International Symposium on Electronics and Smart Devices (ISESD)*, pp. 337-342. IEEE, 2016.
- [31] Kim, "Mq-8 Hydrogen/H2 Sensor Module", <https://sandboxelectronics.com>, p.196, 24 December 2019.
- [32] A. Rahman, W. Jaafar, K. Maulud, N. Noor, M. Mohan, A. Cardil, C. Silva, N. Che'Ya, and N. Naba, "Applications of Drones in Emerging Economies: A case study of Malaysia", *International Conference on Space Science and Communication (IconSpace)*, pp. 35-40. IEEE, 2019.

# An Improved Image Retrieval by Using Texture Color Descriptor with Novel Local Textural Patterns

Punit Kumar Johari<sup>1</sup>, Rajendra Kumar Gupta<sup>2</sup>

Department of CSE and IT  
Madhav Institute of Technology and Science, Gwalior, India

**Abstract**—This paper proposes a new local descriptor of color, texture known as a Median Binary Pattern for color images (MBPC) and Median Binary Pattern of the Hue (MBPH). These suggested methods are extract discriminative features for the color image retrieval. In the surrounding region of a local window, the suggested descriptor classification uses a plane to a threshold that distinguish two classes of color pixels. The Median Binary Patterns of the hue features are derived in the color space from HIS, called MBPH to maximize the discriminatory power of the proposed MBPC operator. In addition to MBPC, MBPH are fused to extract the MBPC+MBPH resulting in an efficient image recovery method combined with color histogram (CH). The structure of the two suggested MBPC and MBPH descriptors are combined with the other fuzzyfied based color histogram descriptor that formed MBPC+MBPH+FCH to improve the performance of the suggested method. The proposed methods are applied on datasets Wang, Corel-5K, and Corel-10K. Experimental results depicted that results of proposed methods are better than existing method in terms of retrieved accuracy. The significant recognition accuracy obtained from the proposed methods which is 60.1 and 63.9 for Wang dataset, 41.88 and 42.47 for Corel-5K and 32.89 and 33.89 for Corel-10K dataset. This hybrid proposed method greatly deals with different textural patterns as well as able to grasp minute color details.

**Keywords**—Image retrieval; Binary pattern; feature extraction; Median Binary Pattern for Color (MBPC) image; Median Binary Pattern for Hue (MBPH)

## I. INTRODUCTION

One of the demanding research domain in the context of intelligent system and computer vision is Content Based Image Retrieval (CBIR). In the past, with the explosion of digital technologies such as multimedia sharing platforms, social networks, and priceless technology available with at most every people that produces millions of images in different scenario and replicate via hosting services [1, 2].

The relevant information is possible only with searching and indexing the massive volume of accessible digital images [3, 4] is only possible with more-and-more likely information retrieval system. Now a days at the advancement of the research domains the CBIR system has paying attention towards various researchers to improve methods which gives high recovery proportion inside less recovery rate. In the last decade, CBIR frameworks were boosted on the grayscale images, at that point by methods for the broad utilization of color image over the various systems, that improvement of a color characteristics for acknowledgement and getting reason. Now it is being joined to upgrade the recovery framework

execution. Therefore, designing precise and fast system has become demanding research domain in the field of recognition of pattern and AI. The CBIR system mainly work on two criteria: extraction of feature and matching feature. The key criteria is feature extraction as it requires with very small variation to be signified by vastly discriminated features from the image. These features discriminate within the class images and major variations between the other existing class images. Essentially building block of the CBIR framework gets request image as input from the expected user and for the purpose of feature extraction from the query image it utilizes a descriptor (may be combination of image content) [5,6]. Different indexing methods has been utilized and the query image highlights is contrasted and the arrangement of highlight vectors of various picture database. The retrieval of images based on most related images from an image database and delivered to the user.

Mainly researchers are focused the features extraction procedures according to the application requirement in different domains. The main goal is to extract distinct features in the viable time. Features extraction is an important task for any multimedia retrieval process. In recent years, characteristic of extraction has been thoroughly investigated [7,8,9,10]. Global characteristic include depiction of the contours in the image, the style definition, texture characteristics and local characteristics [11,12]. Some example of global descriptors are from matrices and invariant moments as well as histogram dependent gradients [13]. Some problem relates to occlusion, viewing and lighting changes and local characteristics of imaging are being dealt with by global methods due to their insufficiency.

The extraction methods through native region that extract features from the local region of the image are well adapted to these issues. Such regions may be as small as easy or chosen by key points as image portions. Texture and color offers knowledge of significance information form the development of efficient features by confirming the strong recovery system output. The classic texture characteristic was derived from the grayscale images. The scale-invariant features transform (SIFT) is more productive, viable and precise descriptor in a cutting edge acknowledgement and characterization framework, among the diverse nearby local descriptor for grayscale images [14]. In order to capture the texture from color images, several variants are available. The SIFT color was evaluated and found to outperform a number of color descriptors. However, the SIFT color is an intensive computation, particularly when scaling the image or the



dimension of the database is increased. Ojala et al. [15], the LBP, often a significant descriptor of texture, has been found in several pattern recognition systems and in computer vision applications to be effective and strong. It preserves characteristics of local texture that are invariant for changes in lighting [16,17]. Most of the efforts have been put in texture recognition [18-21], face analysis [22-24], identification of facial expression [25-28], image recovery [29,30] and so forth, the LBP operator has been successfully used. Most research has been performed on the standard LBP operator and their gray-scale imaging models [31-34]. Color images on the internet are increasingly demanded and used for many realistic applications. Researchers have created descriptor images which denote color texture designs as well as LBP operators for gray scale images [35-39].

In this paper suggested operators MBPC and MBPH for color images which extracts color image structures that imitate the gray-scale texture extracted by the LBP operator. A vector of  $m$  components is called a color pixel and we create a hyperplane for that reason. The hyperplane is used as a threshold boundary with two classes of dividing color pixels. The value is assigned 1 if it is on a plane or above and 0 if it is below a plane, in a  $3 \times 3$  neighborhoods of the current pixels. These operators proposed thus establishes spatial relationships between color pixels, which represent local texture characteristics. In a manner that matches LBP operator histogram for gray images, that may calculate binary pattern histograms extracted from color images. There are 256 histogram bins, while suggested operator's uses eight color pixels in the neighborhood, as features reflecting local image texture patterns. Proposed method is based on the channel color histogram (CH) from the H-specific (HIS) model, which is fusion with MBPC + MBPH. Among the best color image descriptors, most researcher frequently chosen the color histogram. Thus the MBPC+MBPH+CH solution proposes. Similarly, also suggested another solution MBPC + MBPH + FCH, where fuzzyfied color channel histogram [40] based on the H channel color histogram (CH) from the HSV color space is used. The purpose of this study to improve image retrieval rate of related images from each categories of different dataset efficiently.

Although the descriptor output is greatly interrelated in both settings, substantial discrepancies have to be taken into consideration while choosing the descriptors for large-scale jobs. For the best descriptors, a correlation review is given, indicating the best possible combination of its use.

This paper is categorized in the following way: In Section 2, given a color image description of those binary descriptors of the current state of art focused on the local trends. Section 3 derived the suggested MBPC and MBPH operator. In Section 4 one can consider the characteristics and the fusion of the MBPC and MBPH operator with the color histogram (CH) and fuzzyfied Color Histogram (FCH) respectively. Section 5 outlines the different dimension of similarity measurements and results used in applications of image retrieval. In Section 6 presented a comprehensive experimental analysis on different color image databases for their retrieval performance. In Section 7 the conclusion and future work are remarked finally.

## II. A TAXONOMY OF RELATED WORK

In this section, presents an outline of the firmly correlated works which have a place with the varieties of Local-Binary Patterns like strategies produced for the color images to use the benefit of the relationship between the color channels. An explained below, every method has its own advantages and disadvantages.

### A. Local Binary Pattern (LBP)

Firstly proposed a method that describe local information very efficiently within an image [15]. The traditional LBP generates an 8-digit binary number that generates a binary pattern. For the convenience a binary pattern is transformed to be code usually a number in decimal. Within the grayscale image, a  $3 \times 3$  block  $B_m$  has a central pixel  $C_m$ , compared to the surrounding pixels of  $C_m$ , the LBP coding will be done as follows:

$$LBP_{C_m} = \sum_{n=0}^{|B_m|-1} h(C_n - C_m) 2^n \quad h(t) = \begin{cases} 1, & t \geq 0. \\ 0, & t < 0. \end{cases} \quad (1)$$

Where  $|B_m|$  is the number of elements in  $B_m$ .

### B. Uniform Local Binary Pattern (ULBP)

To reduce the feature vector scale, the uniform binary pattern is used in contrast with the LBP [15]. The LBP estimation of pixel  $(i, j)$ , let's represent in the  $LBP_{p, radius}(i, j)$ . In order to  $S$  mean the string of twofold qualities. Only those patterns of  $2^p$  binary patterns (denoted as ULBP) which fulfil this requirements are referred to as Uniform:

$$\sum_{i=1}^p |S_i - S_{i-1}| + |S_0 - S_p| \leq 2 \quad (2)$$

where  $S_i$  denote the  $i^{\text{th}}$  bit of the string. The rest of the patterns are farmed regular and located into a single group. In all the uniform patterns is  $P(P - 1) + 3$  that is less than  $2P$ , especially if  $P$  is big. For  $P = 8$  there are 256 and 59 patterns. Respectively in total LBP and ULBP pattern for example.

### C. Multispectral Local Binary Pattern (MSLBP)

Mäenpää et al. [41] using RGB color channels and six adversary color pairs. Those opposing colors, which provide cross-linkage between values of colors and spatial relations, are used to obtain two color-structure characteristics. The three LBP-functional vectors can be accomplished by treating each channel of an RGB image like gray images [19], analogous to the LBP functionality of the gray scale. The following are derived from the six opposing LBP vectors:

$$MSLBP^{(i,j)}(x_c, y_c) = \sum_{p=0}^{P-1} S(v^i(x_p, y_p) - v^i(x_c, y_c)) \times 2^p \quad (3)$$

where  $(i, j) = \{(1,2), (2,3), (3,1), (2,1), (3,2), (1,3)\}$ , and

$$S(v^i(x_p, y_p) - v^i(x_c, y_c)) = \begin{cases} 1, & \text{if } (v^i(x_p, y_p) - v^i(x_c, y_c)) \geq 0 \\ 0, & \text{otherwise} \end{cases} \quad (4)$$

Here  $v^i(x_c, y_c)$  is the center of pixel intensity of the image component  $j^{\text{th}} = 3$  window,  $v^i(x_p, y_p)$  is a neighborhood pixel intensity value of the image  $i^{\text{th}}$  color variable. There are 2304 MSLBP features, i.e. the function vector after nine LBP

operators have been concatenated. A high recognition tare is given by the method. Nevertheless, the vector size is too larger to decrease the recovery rate.

**D. Local Color Vector Binary Pattern for Face Recognition (LCVBP)**

In LCVBP method color angular patterns and color norm patterns are the two discriminative patterns [22].

$$f_{cn}^p = \varphi_{cn}(h_{cn}^p) \text{ and } f_{ca_{i,j}}^p = \varphi_{ca_{i,j}}(h_{ca_{i,j}}^p), \text{ for } i < j \quad (5)$$

$i = 1, \dots, K - 1$ , and  $j = 2, \dots, K$ .

Similar calculations may be performed for lower-sized features  $f_{cn}^g$  and  $f_{ca_{i,j}}^g$  of  $h_{cn}^g$  and  $h_{ca_{i,j}}^g$ . The lower dimension features of the LBP histogram are joint at feature level by concentrating the lower dimension features in column order respectively in the following order for the proposed LCVBP feature of  $I^p$  and  $I^s$ .

$$f_{LCVBP}^p = [(f_{cn}^p)T (f_{ca_{1,2}}^p)T \dots (f_{ca_{1,K}}^p)T \dots (f_{ca_{K-1,K}}^p)T]T \quad (6)$$

$$f_{LCVBP}^g = [(f_{cn}^g)T (f_{ca_{1,2}}^g)T \dots (f_{ca_{1,K}}^g)T \dots (f_{ca_{K-1,K}}^g)T]T \quad (7)$$

**E. Quaternionic Local Ranking Binary Pattern (QLRBP)**

This is another descriptor for color image. A QLRBP was developed by Lan et al. [42] to integrate multi-spectral channel color knowledge in color pictures, and a local quaternionic rating binary pattern was adopted. The operator QLRBP extracts the quaternionic representation (QR) of color image. Without having to treat every color channel individually, the QLRBP may manage all color channel directly within the quaternionic field. The ranking function can be expressed as:

$$R_{QLRBP}(q_m, q_n) = \delta_{CTQ}(q_n, p_1) - \delta_{CTQ}(q_m, p_1) \quad (8)$$

$$QLRBP_{q_m} = \sum_{n=0}^{|S_m|-1} h(R_{QLRBP}(q_m, q_n)) 2^n \quad (9)$$

**F. Multichannel Decoded Local Binary Pattern (MDLBP)**

MDLBP based on decoder based local binary pattern  $mdLBP_{t_2}^n(i, j)$  for pixel  $(i, j)$  from multichannel decoder map [43]  $mdM^n(i, j)$  and  $t_2$  can be computed as:

$$mdLBP_{t_2}^n(i, j) = \begin{cases} 1, & \text{if } mdM^n(i, j) = (t_2 - 1) \\ 0, & \text{otherwise} \end{cases} \quad (10)$$

For  $\forall t_2 \in [1, 2^c]$  and  $\forall n \in [1, N]$ .

**G. Color Histogram Creation Method based on Fuzzyfication (FCH)**

A motivation behind picking up the HSV color space is that it is recognizably uniform and approximates how individuals perceives. The primary expiation, however, is that HSV color space has been discovered to be stronger than other color spaces in various retrieval experiments. The fuzzy link based approach is used to construct a histogram in the descriptor. More than one histogram output is specified by the term “fuzzy-linking” [40]. The input channels are described in the following fuzzy sets:

- The channel Hue(H) is divided into 10 fuzzy regions,
- The channel Saturation (S) is divided into 3 fuzzy regions,
- The channel Value (V) is divided into 3 fuzzy regions.

$$\mu_A(x) = \begin{cases} 0, & (X < t) \text{ or } (X > v) \\ \frac{x-t}{u-t}, & t \leq X \leq u \\ 1, & u \leq X \leq w \\ \frac{v-x}{v-w}, & w \leq X \leq v \end{cases} \quad (11)$$

Where t is the lower limit, v is upper limit, u is lower support limit and w is upper support limit, also  $t < u < w < v$ . The membership features are displayed in Fig. 1 to 3. For the H channel, the final fuzzy histogram includes only 10 bins out of the 12 bins. The prevailing hues in each image can promptly be taken note. As portrayed in the accompanying area, the histogram in the proposed framework has demonstrated to be an instrument for accurate image recovery.

Histogram bins that are shown in Fig. 1 to 3 is concerning: (1) Red, (2) Orange, (3) Yellow, (4) Light Green, (5) Green, (6) Spring Green, (7) Cyan, (8) Azure, (9) Blue, (10) Violet, (11) Magenta and (12) Rose Red for H Channel. Correspondingly, for S channel receptacles are: (1) Low, (2) Medium and (3) High Saturation. What’s more, for V Channel receptacles are: (1) Dark, (2) Light and (3) Bright.

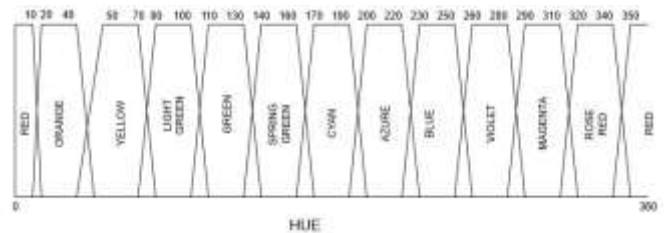


Fig. 1. Membership Function of Hue Color.

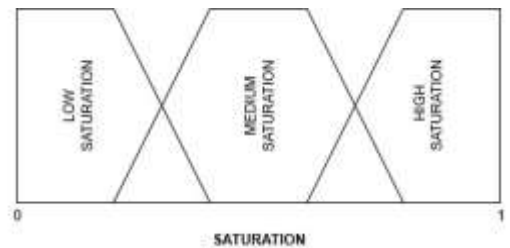


Fig. 2. Membership Function of Saturation Channel.

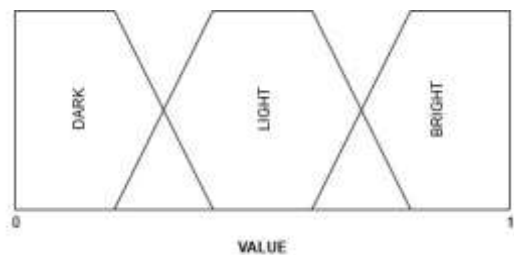


Fig. 3. Membership Function of Value Channel.

### III. PROPOSED WORK

The Median Binary Pattern for color (MBPC) is first briefly described in this section.

#### A. Median Binary Pattern (MPB)

For gray scale images, the classic MBP [50] operator is described. For a circular symmetric neighbor group of  $P$  members, the common form of the operator is represent as:  $R(x_p, y_p)$  and median  $(x_m, y_m)$ , as defined by  $MBP_{P,R}(x_m, y_m)$ :

$$MBP_{P,R}(X_m, Y_m) = \sum_{p=0}^{P-1} S(I(x_p, y_p) - I(x_m, y_m)) \times 2^p \quad (12)$$

Where  $P$  is the quantity of neighbors and has the estimation of solidarity. The operator of the MBP is invariant with repetitive gray scale changes as the edge doesn't rely upon the intensity. The example examined is the result of spatial associations in the specific area. On the off chance that a specific neighborhood has no correlation, it is known as a spot.

Where

$$S(I(x_p, y_p) - I(x_m, y_m)) = f(x) = \begin{cases} 1, & \text{if } I(x_p, y_p) - I(x_m, y_m) \geq 0 \\ 0, & \text{otherwise,} \end{cases} \quad (13)$$

and  $I(x_p, y_p)$  is a pixel location intensity  $(x_p, y_p)$ . 8 neighborhood pixels at  $(x, y)$  are called in their simplest form, i.e.  $P = 8$ , and  $R = 1$ . Here, the MBP operator divides the image into two classes of pixels based on the median value of the configuration. Note that MBP compares two levels of severity, which influence the structure locally as well. Such patterns constitute the fundamental element of proposed texture description. The MBP operator produces binary patterns called MBP patterns, which range from 0 to  $2^P - 1$ . For each pixel of an image, MBP patterns are obtained under histograms with MBP patterns. These histograms define the texture of an image with a gray scale. When the value of  $P$  is higher, than scale of histogram patterns will be greater.

In MBP the threshold directly does not depend on amount of pixel intensity therefore it is invariant to monotonic gray scale changes in the image. Through the spatial interfaces in the given locality as a result the pattern information has been detected. Fig. 4 shows the calculation of MBP.

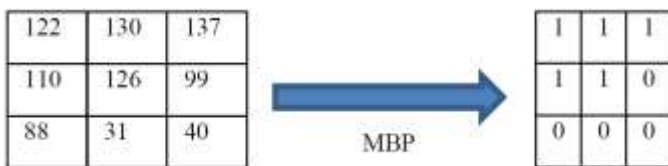


Fig. 4. Showing MBP, the Median Value is 110.

#### B. The Suggested MBPC Operator

The suggested MBPC operator uses a hyperplane in the  $m$  dimensions for partitioning or thresholding of color pixels. As  $m = 3$  in this case, 3D space is hyper-plan. Namely a hyper-ellipsoid or a hyper-cube, a hyper-sphere, are other threshold alternatives, they are, however, not as effective as a hyperplane. After extensive experiments with these threshold

alternatives, this work came to this conclusion. The future challenge is to construct a thresholding in 3D color that is performed as follows. The color pixel elements of RGB color space is to represented by the vector  $I(x, y) = (r(x, y), g(x, y), b(x, y))$  or, simply,  $I = (r, g, b)$ . Consequently, the color variable of  $m$  is 3. The local window  $(2R + 1) \times (2R + 1)$  of size is defined as the corresponding color vector  $(r_m, g_m, b_m)$ ,  $R \geq 1$ , median at pixel  $m$ . Let  $I_p = (r_p, g_p, b_p)$  be a pixel  $p$  of the neighborhood. Within the color space, the color plane  $L$  is established. Through specifying the normal and a reference point in the color space a plane is derived. The level is inferred by characterizing the typical to the plane and a reference point in the color space. Let  $n = (n_1, n_2, n_3)$  be denoted as normal to the plane and  $R_o = (r_o, g_o, b_o)$  be the reference point,

At the point the condition of the level with the normal vector  $n$  and the reference point  $R_o$  is presented as:

$$n \cdot (I_p - R_o) = 0. \quad (14)$$

or,

$$n_1(r - r_o) + n_2(g - g_o) + n_3(b - b_o) = 0. \quad (15)$$

The consequence of the dot product between the vector  $n$  and the vector generated by joining in two classes  $R_o = (r_o, g_o, b_o)$  to  $I_p = (r, g, b)$  space. The vector  $I_p - R_o$  produces every pixel on or above the plane. The color plane separates the color into unique class and all the other pixels below the plane into another class. There are many ways to choose normal vector  $n$  but a line that joins the dark pixel  $(0, 0, 0)$  and the pure white pixel  $(1, 1, 1)$  is an evident choice. It reflects the gray line and all primary colors  $R, G$  and  $B$  are equally present. The median pixel  $I = (r_m, g_m, b_m)$  of the square neighborhood would be an obvious choice for the reference point. Color plane is defined, which is normal for line connections  $(0, 0, 0)$  and  $(1, 1, 1)$ , with these values of the two parameter, and passes over  $R_o = I_m = (r_m, g_m, b_m)$ . It enables neighborhood pixel thresholding method  $I_p, p = 0, 1, \dots, P - 1$  into two classes: those above or on the plane, and some below the plane. Here,  $P$  reflects the total number of neighborhood pixels. The following expression can be assessed for a decision to this effect.

$$E_v(I_p) = E_v(r, g, b) = n_1(r_p - r_c) + n_2(g_p - g_c) + n_3(b_p - b_c), \quad (16)$$

The color pixel  $I_p = (r_p, g_p, b_p)$  is above or above the plane if  $E(I_p) \geq 0$  and below the plane if  $E(I_p) < 0$ . Therefore the color pixels are divided into two groups with a clearly defined process in the median pixel neighborhood. An extension of binary patterns on local grayscale, which is obtained by restricting gray values to the gray value for a median pixel  $(x_m, y_m)$  of the local window can be considered.

#### C. The Suggested MBPH Operator

MBP is a classical gray scale operator that cannot be further drawn-out to color images since a color pixel represented a vector quantity of  $R, G$ , and  $B$  components and a scalar quantity was the gray scale pixel. Consequently, for

color pixels to get a binary set, a comparison of the Eq. (18) type cannot be made.

The proposed Median binary patterns of the hue (H) component (MBPH) from HSI model.

$$MBPH_{p,r}(x_m^h, y_m^h) = \sum_{p=0}^{P-1} S(I(x_p^h, y_p^h) - I(x_m^h, y_m^h)) \times 2^p \quad (17)$$

Where

$$S(I(x_p^h, y_p^h) - I(x_m^h, y_m^h)) = \begin{cases} 1, & \text{if } I(x_p^h, y_p^h) - I(x_m^h, y_m^h) \geq 0 \\ 0, & \text{otherwise} \end{cases} \quad (18)$$

$I(x_p^h, y_p^h)$ , is an intensity of pixel at (x,y), the color component in H plane of HSI color image model and  $(x_m^h, y_m^h)$  is a median value of selected patch if an image.

#### IV. FUZZY COLOR HISTOGRAM (FCH) AND MEDIAN BINARY PATTERN OF THE HUE COMPONENT (MBPH) AND THEIR FUSION WITH MBPC FEATURES

Most of the time color features are dominant to identify objects as a whole from the image. The color texture may be obtained from colored image that represents in Hue (H) component of the HIS model. Therefore, the hue element is a normal option for segmenting an image dependent on color. To fuse with different features, suggested method uses fuzzy color histogram (FCH), MBPC of a color image driven by RGB color color space, and MBP is driven by the hue component of HSI image, which is then called MBPH. By using the following color conversion function from a colour image inside the color space of the RGB, the component hue is achieved.

$$H(i, j) = \begin{cases} \theta(i, j), & \text{if } b(i, j) \leq g(i, j) \\ 2\pi - \theta & \text{if } b(i, j) > g(i, j) \end{cases} \quad (19)$$

Where

$$\theta(i, j) = \cos^{-1} \left\{ \frac{\frac{1}{2}[(r(i,j)-g(i,j))+(r(i,j)-b(i,j))]}{[(r(i,j)-g(i,j))^2+(r(i,j)-b(i,j))(g(i,j)-b(i,j))]^{1/2}} \right\} \quad (20)$$

A gray scale image is treated as  $H(x, y)$  whose median binary patterns are in a manner identical to the median gray scale image binary patterns. Such features are called MBPH features. The features of MBPH is computed as follows.

The median of a window  $(2R + 1) * (2R + 1)$  be  $(x_m, y_m)$ . The features of MBPH are found as

$$MBPH(x_m, y_m) = \sum_{p=0}^{P-1} F_m(Q_p) * 2^p \quad (21)$$

Where

$$F_m(Q_p) = \begin{cases} 1 & \text{if } H(x_p, y_p) \geq H(x_m, y_m) \\ 0 & \text{otherwise.} \end{cases} \quad (22)$$

Here  $(x_p, y_p)$  is the coordinates of a neighborhood pixel.

Sometimes color histogram image descriptors especially useful. Fuzzy color histogram (FCH) is an extended version of CH descriptor.

Through separating them through image scale before their fusion, all histogram bins are normalized. The MBPC, MBPH and FCH are three feature vectors, of size  $s_c$ ,  $s_h$  and  $s$  correspondingly, are combined to create a single dimension vector of size  $s_c + s_h + s$  whose components are represented by.

$$\{MBPC[0], MBPC[1], \dots, MBPC[s_c - 1], MBPH[0], MBPH[1], \dots, MBPH[s_h - 1], FCH[0], FCH[1], \dots, FCH[s - 1]\}. \quad (23)$$

#### V. SIMILARITY MEASURES AND ESTIMATION METRICS

##### A. Similarity Measures

Various similarity measures have been suggested in the literature for image processing systems. Retrieval in the CBIR system, performance of retrieval not only depends on robust features but also measures through different similarity functions available in the literature. When work with histogram-based feature vectors, this dimension will underpin option of similarity tests. Four of these common parallels are Euclidean distance, chi-square, extended-Canberra and square-chord for histogram dependent function vectors. Certain distance measurements widely employed, including histogram cross section, L1-norm, L2-norm, Jeffrey gap, cos-correlation, etc., are not as successful as the previous ones [53]. With respect to such a similarity metric, evaluate the performance comparison and try to find an effective distance measure which provides the finest total recovery outcomes. It provides the following descriptions of the distance measures.

Let  $FV_i^q$  and  $FV_i^{db}$  represents the  $i^{th}$  feature components of image query 'q' and image database 'db', respectively. Feature vector size indicate by FVSize. Following is the formula of distance measures:

Euclidean Distance

$$Distance_{Euclidean} = \sqrt{\sum_{i=0}^n (|Q_i - D_i|)^2} \quad (24)$$

Extended-Canberra Distance:

$$Distance_{Extended-Canberra}(q, db) = \sum_{i=0}^{FVSize-1} \frac{|FV_i^q - FV_i^{db}|}{(FV_i^q + \mu^q) + (FV_i^{db} + \mu^{db})} \quad (25)$$

Chi-square

$$Distance_{chi}(q, db) = \sum_{i=0}^{FVSize-1} \frac{(FV_i^q - FV_i^{db})^2}{FV_i^q + FV_i^{db}} \quad (26)$$

Square – Chord distance

$$Distance_{Square-Chord}(q, db) = \sum_{i=0}^{Size-1} (\sqrt{FV_i^q} - \sqrt{FV_i^{db}})^2 \quad (27)$$

##### B. Estimation Metrics

All images in the experiments is used as a query image in the database. P(N) and R(N) precision can be used to measure the performance of the image retrieval. Top N images described in [54].

$$P(N) = \frac{I_o}{N}; R(N) = \frac{I_o}{M} \quad (28)$$

Where  $M$  is the total number of images that are identical to query images in the dataset and  $I_o$  is the total number of relevant images from higher index positions? The sum of all  $P(n)$  exact values is the average precision of the single query  $\bar{P}(q)$ ,  $n = 1, 2, \dots, N, i. e.$

$$\bar{P}(q) = \frac{1}{N} \sum_{n=1}^N P(n). \quad (29)$$

For each queries  $Q$ , the mean average precision (mAP) is the mean of the average scores:

$$mAP = \frac{1}{Q} \sum_{q=1}^Q \bar{P}(q). \quad (30)$$

If the number of the related images in each class differ, the graph  $P - R$  is not a satisfactory indicator. The  $mAP$  measure in [55].

## VI. RESULTS AND DISCUSSIONS

This section presents multiple experimental results which demonstrate the effectiveness and comparison [52] of the suggested methods with those of the closely linked operators of color textures, such as local bilateral component image patterns, LCVBP, MSLBP, MDLBP and QLRBP. For MDLBP operators, the decoder operator performs better than the adder operator. Therefore, decoder operator is considered for the performance comparison. The LBP image components are essentially an extension of the LBP to the R, G and B components of the color image. Another effective global descriptor for texture feature extraction is a Gabor filter which has been applied to gray scale images for texture image retrieval [44-49].

### A. Datasets

To analyzing the proposed method three datasets namely Wang, Corel-5K and Corel-10K are used. In the following, these datasets is briefly explained [51].

Wang [53]: It comprises 1000 color images separated into 10 groups of 100 images each. It includes one of Corel's image databases. Every class includes  $265 \times 384$  or  $384 \times 256$  pixel resolution images. The 10 classes of Wang image database are: African tribe people, Bus, Dinosaur, Flower, Beach, Elephant, Buildings, Food, Horse and Glacier.

Corel- 5K [55]: This dataset contains 50 groups of images and each group has image of size  $128 \times 192$  or  $192 \times 128$  pixels in JPEG format. Each group has 100 images, with different substance like mountain, tiger, fort, mushroom, car, ticket, ocean etc. in total 5000 images.

Corel- 10K [55]: This dataset contains 100 groups of images and each group has image of size  $128 \times 192$  or  $192 \times 128$  pixels in JPEG format. Each category has 100 images of various substances like rose, sunset, cat, train, duck, fish, judo-karate, etc. in total 10,000 images.

### B. Comparison with Existing Methods

In contrast with the following techniques, the findings of the suggested methods are compared: mean average precision (mAP) performance is based on LBP, ULBP, MSLBP, LCVBP, QLRBP, and MDLBP. The results also list the number of features used by all techniques. The methods

proposed are also applied separately by fusion of their characteristics in two combinations. MBPH, FCH are separate processes, while MBP+MBPH+CH while MBP+MBPH+FCH are two variations. Such strategies are tested to determine their relative efficiencies and to define the best solutions for high precision and low recovery times. The benefit of speed is therefore that the feature vector with low dimensional dimensions does not cause a great deal of recovery precision. There is a similar trend with MBPH and FCH operators, which is described in the following experimental analysis. The results are comparable with the following approaches: LBP, ULBP, LCVBP, MSLBP, QLRBP, GABOR, MDLBP, MBPC+MBPH+CH, and MBPC+MBPH+FCH.

Proposed method is compared the following approaches. In order to evaluate output of a value of  $N$ , a deeper examination of the amount of the images obtained was conducted, and the results for  $N=1$  to 12 were given from 100 values of  $N$ , in Table I.

Table II shows the mAP values, attained by the several method includes top 100 image ( $N = 100$ ) of the Wang dataset. The images in all databases are used as images for queries. It is noted from the table that, with the square chord distance, the MBPC+MBPH+FCH solution proposes the highest mAP of 63.9. The efficiency of the resulting method is significantly improved when FCH features fuse with MBPC+MBPH to obtain the MBPC+MBPH+FCH method. Finally, the efficiency of the estimated square chord interval is better seen in the mAP values displayed in Table II. The values of precision versus recall are shown in Fig. 5 for 9 methods (7 existing, 2 proposed) that produce the overall results.

Euclidean, Chi-sqaure, Extended Canberra and Square-Chord obtained average mAP value are 57.32%, 52.23%, 58.08%, and 63.9%, respectively.

The results shown by Table III is based on Corel-5K dataset. The proposed MBPC + MBPH + FCH method achieves a maximum mAP value of 42.47% followed by the MBPC + MBPH + CH method that results in a mAP of 41.88%. MDLBP, which uses Square-Chord, reaches the next largest mAP value. GABOR, which is 37.12 %, is the fourth largest mAP. The decreased mAP values of other existing methods performance is: LCVBP (36.4%), LBP (35.41), MSLBP (35.33), QLRBP (35.22), and ULBP (32.88%). The difference of mAP values from the proposed MBPC + MBPH + CH and MBPC + MBPH + FCH, which is insignificant, is only 0.59, whereas the size of feature vector differences in the proposed method is significant. The distance measured by Euclidean, Chi-sqaure, Extended Canberra and Square-Chord obtained average mAP value are 36.35%, 31.05%, 37.9%, and 42.47%, respectively for method MBPC + MBPH + FCH. Fig. 6 plots the precision and recall values of existing and suggested approaches. The suggested methods are seen to outperform the conventional methods for all recall values.

The findings for Corel-10 K are shown in Table IV. The development for mAP values is closed to the development for datasets Wang and Corel-5k. The suggested MBPC + MBPH + FCH solution reaches the 33.89% highest value of mAP. It should be noted here that although the mAPs of these two

proposed methods vary slight difference by 1 percent. Although the second largest value is obtained by MDLBP, which is 33.23% with 2048 features size is higher than compared with the proposed methods. Therefore, the third

largest value of mAP, is attained by MBPC+MBPH+CH that is 32.89% by using 1055 features. Precision versus recall values for N=100 is shown in Fig. 7.

TABLE I. NUMBER OF SIMILAR IMAGES (IN PERCENT) OBTAINED FOR EACH CATEGORIES FROM WANG DATASET VALUE OF  $N(N = 1 \text{ to } 12)$  BY THE SUGGESTED MBPC+MBPH+CH (METHOD X) AND BY MBPC+MBPH+FCH (METHOD Y)

| N  | Class   |     |       |     |          |     |     |     |          |     |          |     |        |     |       |     |         |     |      |     |
|----|---------|-----|-------|-----|----------|-----|-----|-----|----------|-----|----------|-----|--------|-----|-------|-----|---------|-----|------|-----|
|    | African |     | Beach |     | Building |     | Bus |     | Dinosaur |     | Elephant |     | Flower |     | Horse |     | Glacier |     | Food |     |
|    | X       | Y   | X     | Y   | X        | Y   | X   | Y   | X        | Y   | X        | Y   | X      | Y   | X     | Y   | X       | Y   | X    | Y   |
| 1  | 100     | 100 | 100   | 100 | 100      | 100 | 100 | 100 | 100      | 100 | 100      | 100 | 100    | 100 | 100   | 100 | 100     | 100 | 100  | 100 |
| 2  | 94      | 92  | 85    | 80  | 86       | 89  | 96  | 99  | 100      | 100 | 92       | 91  | 97     | 99  | 97    | 99  | 76      | 77  | 90   | 92  |
| 3  | 91      | 87  | 76    | 70  | 77       | 84  | 93  | 98  | 99       | 99  | 85       | 85  | 95     | 97  | 96    | 98  | 68      | 68  | 85   | 87  |
| 4  | 88      | 84  | 73    | 66  | 73       | 80  | 93  | 96  | 99       | 99  | 82       | 79  | 95     | 97  | 94    | 98  | 63      | 63  | 83   | 84  |
| 5  | 86      | 83  | 69    | 61  | 68       | 78  | 93  | 96  | 99       | 99  | 77       | 75  | 94     | 96  | 91    | 97  | 59      | 60  | 80   | 83  |
| 6  | 85      | 82  | 67    | 59  | 66       | 75  | 92  | 96  | 99       | 99  | 76       | 72  | 93     | 96  | 90    | 96  | 57      | 57  | 79   | 80  |
| 7  | 83      | 82  | 66    | 57  | 63       | 73  | 91  | 95  | 99       | 99  | 73       | 68  | 93     | 96  | 87    | 96  | 55      | 55  | 78   | 79  |
| 8  | 82      | 81  | 64    | 55  | 60       | 72  | 90  | 94  | 99       | 99  | 72       | 66  | 93     | 96  | 86    | 96  | 54      | 54  | 76   | 76  |
| 9  | 81      | 80  | 63    | 54  | 58       | 69  | 89  | 94  | 99       | 99  | 69       | 63  | 92     | 95  | 85    | 95  | 52      | 52  | 74   | 75  |
| 10 | 80      | 79  | 62    | 53  | 58       | 68  | 88  | 94  | 99       | 99  | 68       | 61  | 92     | 95  | 84    | 94  | 50      | 50  | 72   | 74  |
| 11 | 79      | 79  | 61    | 52  | 56       | 68  | 87  | 93  | 99       | 99  | 68       | 58  | 91     | 94  | 83    | 94  | 48      | 48  | 70   | 72  |
| 12 | 78      | 78  | 60    | 51  | 55       | 66  | 87  | 93  | 99       | 99  | 66       | 57  | 91     | 94  | 82    | 94  | 47      | 47  | 68   | 71  |

TABLE II. MEAN AVERAGE PRECISION (MAP) IN PERCENT FOR N = 100 OBTAINED USING VARIOUS APPROACHES ON WANG DATASET

|          | Method        | No. of features       | Euclidean | Chi-square | Extended-Canberra | Square chord |
|----------|---------------|-----------------------|-----------|------------|-------------------|--------------|
| Existing | LBP [15]      | $3 \times 256 = 768$  | 45.94     | 55.28      | 56.93             | 55.33        |
|          | ULBP[15]      | $3 \times 59 = 177$   | 48.11     | 53.34      | 54.19             | 53.37        |
|          | MSLBP[41]     | $9 \times 256 = 2304$ | 50.20     | 59.86      | 60.62             | 59.86        |
|          | LCVBP[22]     | $4 \times 59 = 236$   | 47.21     | 53.52      | 56.83             | 53.44        |
|          | QLRBP[42]     | $3 \times 256 = 768$  | 45.18     | 53.47      | 56.03             | 53.50        |
|          | GABOR[54]     | 96                    | 49.10     | 58.86      | 59.53             | 58.91        |
|          | MDLBP[43]     | $8 \times 256 = 2048$ | 50.10     | 59.58      | 60.82             | 59.58        |
| Proposed | MBPC_MBPH_CH  | =1055                 | 53.67     | 50.77      | 57.93             | 60.01        |
|          | MBPC_MBPH_FCH | =1631                 | 57.32     | 52.23      | 58.08             | 63.90        |

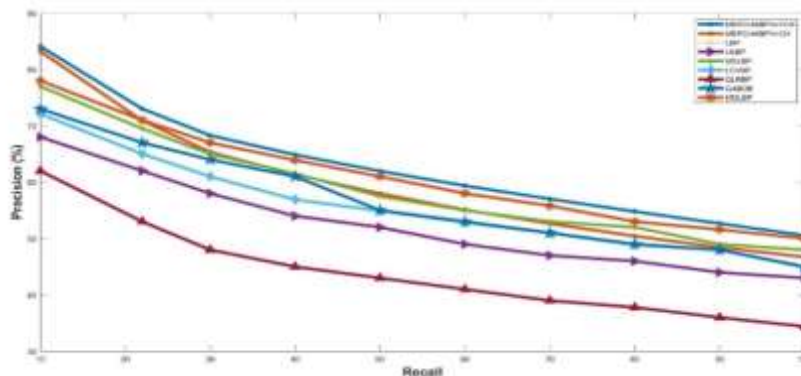


Fig. 5. Precision-Recall Curve for Existing and Proposed Methods (MBPC + MBPH + FCH and MBPC + MBPH + CH) for N = 100 on Wang Dataset.

TABLE III. MEAN AVERAGE PRECISION (MAP) IN PERCENT FOR N = 100 OBTAINED USING VARIOUS APPROACHES ON COREL-5K DATASET

|          | Method        | No. of features       | Euclidean | Chi-square | Extended-Canberra | Square chord |
|----------|---------------|-----------------------|-----------|------------|-------------------|--------------|
| Existing | LBP [15]      | $3 \times 256 = 768$  | 27.43     | 35.38      | 35.75             | 35.41        |
|          | ULBP[15]      | $3 \times 59 = 177$   | 27.21     | 32.87      | 34.26             | 32.88        |
|          | MSLBP[41]     | $9 \times 256 = 2304$ | 28.1      | 35.53      | 39.95             | 35.33        |
|          | LCVBP[22]     | $4 \times 59 = 236$   | 29.5      | 36.39      | 37.95             | 36.4         |
|          | QLRBP[42]     | $3 \times 256 = 768$  | 28.12     | 35.31      | 36.54             | 35.22        |
|          | GABOR[54]     | 96                    | 29.48     | 37.18      | 36.95             | 37.12        |
|          | MDLBP[43]     | $8 \times 256 = 2048$ | 29.56     | 37.89      | 39.99             | 38.08        |
| Proposed | MBPC_MBPH_CH  | =1055                 | 36.31     | 32.59      | 38.52             | 41.88        |
|          | MBPC_MBPH_FCH | =1631                 | 36.35     | 31.05      | 37.9              | 42.47        |

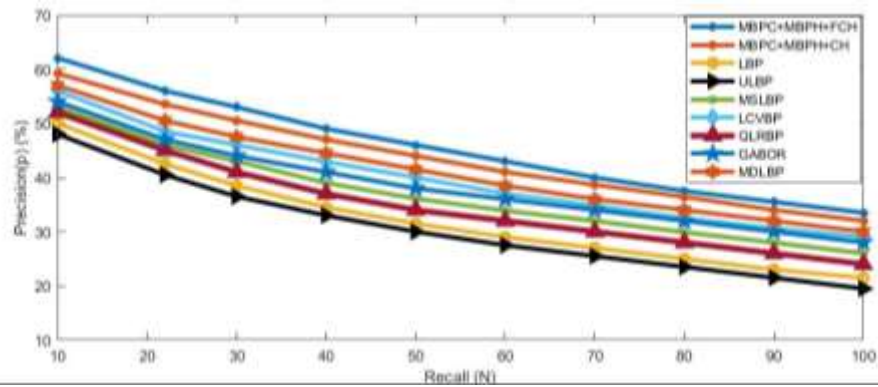


Fig. 6. Precision-Recall Curve for Existing and Proposed methods (MBPC + MBPH + FCH and MBPC + MBPH + CH) for  $N = 100$  on Corel-5K Dataset.

TABLE IV. MEAN AVERAGE PRECISION (MAP) IS PRESENT FOR N = 100 OBTAINED USING VARIOUS APPROACHES ON COREL-10 K DATASET

|          | Method        | No. of features       | Chi-square | Extended-Canberra | Square chord |
|----------|---------------|-----------------------|------------|-------------------|--------------|
| Existing | LBP [15]      | $3 \times 256 = 768$  | 28.97      | 29.33             | 28.98        |
|          | ULBP[15]      | $3 \times 59 = 177$   | 26.96      | 28.08             | 26.97        |
|          | MSLBP[41]     | $9 \times 256 = 2304$ | 28.99      | 31.69             | 28.86        |
|          | LCVBP[22]     | $4 \times 59 = 236$   | 29.25      | 29.72             | 29.26        |
|          | QLRBP[42]     | $3 \times 256 = 768$  | 26.4       | 27.64             | 26.39        |
|          | GABOR[54]     | 96                    | 28.06      | 29.68             | 28.06        |
|          | MDLBP[43]     | $8 \times 256 = 2048$ | 31.83      | 33.97             | 33.23        |
| Proposed | MBPC_MBPH_CH  | =1055                 | 25.16      | 31.19             | 32.89        |
|          | MBPC_MBPH_FCH | =1631                 | 23.89      | 31.30             | 33.89        |

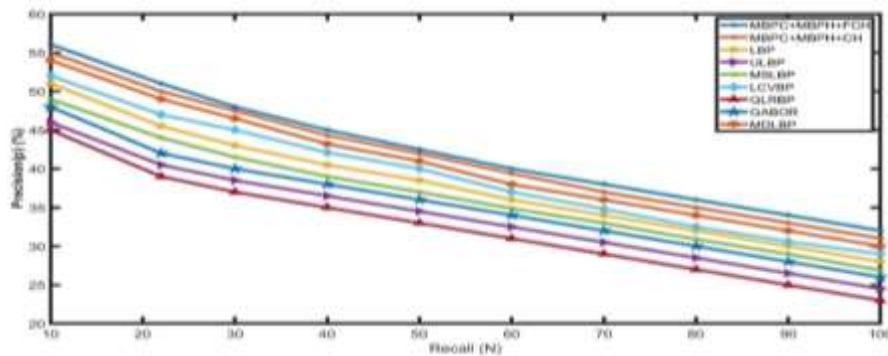


Fig. 7. Precision-Recall Curve for Existing and Proposed Methods (MBPC + MBPH + FCH and MBPC + MBPH + CH) for  $N = 100$  on Corel-10K Dataset.

## VII. CONCLUSION AND FUTURE WORK

In the Color Image Retrieval issue the suggested median binary patterns (MBPC) for color textures are relatively efficient. The precision of recovery is further improved by the derivation of local binary patterns of the HSI color space (H) variable, known as the MBPH. An effective color image descriptor is the Color Histogram (CH) and is used for improving recovery performance using the proposed methods. If such color descriptors have been combined, the output has been evaluated individually as well as in their mixture types. Since both methods are focused on standardized histograms it is a concatenation method that incorporates such functions. There are also strong variations in LBPC and LBPH. Exhaustive performance experimental analysis shows, compared with the significant available local texture based on color, the descriptor and multichannel decoded local binary pattern (MDLBP) of the suggested MBPC+MBPH+FCH method reaches the most elevated estimation of mAP over all datasets. Compared to the strongest current MDLBP system with a large function aspect (2048), it offers far more effective outcomes of recovery at low device costs. The experimental findings also demonstrate that Square chord distance calculation outperforms all the significant distance measurements to award certain approaches the maximum mAP values in all datasets. In future, other effective methods may also be combined to get efficient image retrieval. Feature selection method may also use to obtain a prominent feature subset to improve the retrieval performance.

### REFERENCES

- [1] Furht, B. "Encyclopedia of multimedia (2nd ed.). Springer Science & Business Media", 2008.
- [2] Lin, C.-H. , Chen, H.-Y. , & Wu, Y.-S. "Study of image retrieval and classification based on adaptive features using genetic algorithm feature", Expert Systems with Applications, Vol 41 (15), pp. 6611–6621, 2014.
- [3] Hiwale, S. S. & Dhotre, D. "Content-based image retrieval: Concept and current practices", Paper presented at the Electrical, Electronics, Signals, Communication and Optimization (EESCO), 2015 international conference on, Visakhapatnam, India. 2015.
- [4] M. Swain , D. Ballard , Color indexing, Int. J. Comput. Vis. Vol.7 (1), pp. 11–32, 2017.
- [5] Otávio A.B. Penatti, Eduardo Valle, Ricardo da S. Torres, Comparative study of global color and texture descriptors for web image retrieval, Journal of Visual Communication and Image Representation, Vol. 23, Issue 2, pp. 359-380, ISSN 1047-3203, 2012.
- [6] L. Liu, S. Lao, P. W. Fieguth, Y. Guo, X. Wang and M. Pietikäinen, "Median Robust Extended Local Binary Pattern for Texture Classification," in IEEE Transactions on Image Processing, Vol. 25, no. 3, pp. 1368-138, 2016.
- [7] C. Zhu , C.-H. Bichot , L. Chen , Image region description using orthogonal combination of local binary patterns enhanced with color information, Pattern Recognit. Vol. 46, pp. 1949–1963, 2013.
- [8] J. Li, N. Sang, C. Gao , Completed local similarity pattern for color image recognition, Neurocomputing Vol. 182, pp. 111–117, 2016.
- [9] Y. Zhao, D. Huang, W. Jia, Completed local binary count for rotation invariant texture classification, IEEE Trans. Image Process. Vol. 21 (10), pp. 492–4497, 2012.
- [10] S. Liao, M. Law, A. Chung, Dominant local binary patterns for texture classification, IEEE Trans. Image Process. Vol. 18 (5), pp. 1107–1118, 2009.
- [11] M. Heikkilä, M. Pietikäinen, C. Schmid, Description of interest regions with local binary patterns, Pattern Recognition. Vol. 42 (3), pp. 425–436, 2009.
- [12] L. Nanni, A. Lumini, S. Brahmam, Local binary patterns variants as texture descriptors for medical image analysis, Artif. Intell. Med. Vol. 49 (2), pp.117–125, 2010.
- [13] N. Dalal, B. Triggs , Histograms of oriented gradients for human detection, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 886–893, 2005.
- [14] D.G. Lowe, Distinctive image features from scale-invariant keypoints, Int. J. Comput. Vis. Vol. 60 (2),91–110, 2004.
- [15] Ojala T., M. Pietikäinen, T. Maenpää, Multiresolution gray-scale and rotation invariant texture classification with local binary patterns, IEEE Trans. Pattern Anal. Mach. Intell. Vol. 24 (7), pp. 971–987, 2002.
- [16] R.M. Haralick, K. Shangmugam, Textural feature for image classification, IEEE Trans. Syst. Man Cybern. Vol. 6, pp. 610–621 SMC-3, 1973.
- [17] H. Tamura, S. Mori, T. Yamawaki, Texture features corresponding to visual perception, IEEE Trans. Syst. Man Cybern. Vol. 8 (6) pp. 460–473, 1978.
- [18] T. Mäenpää, M. Pietikäinen , T. Ojala , Texture classification by multi-predicate local binary pattern operators, in: Proceedings of the International Conference on Pattern Recognition (ICPR), pp. 3951–3954, 2000.
- [19] T. Mäenpää, T. Ojala, M. Pietikäinen , M. Soriano , Robust texture classification by subsets of local binary patterns, in: Proceedings of the International Conference on Pattern Recognition (ICPR), pp. 3947–3950, 2000.
- [20] M. Pietikäinen, T. Mäenpää, V. Jaakko , Color texture classification with color histograms and local binary patterns, in: Workshop on Texture Analysis in Machine Vision, pp. 109–112, 2002.
- [21] Y. Guo, Z. Guoying , M. Pietikäinen , Discriminative features for texture description, Pattern Recognit. Vol. 45 (10), pp. 3834–3843, 2012.
- [22] S.H. Lee, J.Y. Choi, Y.M. Ro, K.N. Plataniotis , Local color vector binary patterns from multichannel face images for face recognition, IEEE Trans. Image Process. Vol. 21 (4), pp. 2347–2353, 2012.
- [23] T. Ahonen, A. Hadid , M. Pietikäinen , in: Face Recognition With Local Binary Patterns Lecture Notes in Computer Science, 3021, Springer-Verlag, Berlin Heidelberg New York, pp. 469–481, 2004.
- [24] T. Ahonen, A. Hadid, M. Pietikäinen, Face description with local binary patterns: application to face recognition, IEEE Trans. Pattern Anal. Mach. Intell. Vol. 28 (12), pp. 2037–2041, 2006.
- [25] C. Shan, S. Gong, P.W. McOwan, Robust facial expression recognition using local binary patterns, in: Proceedings of IEEE International Conference of Image Processing, Vol. 2, pp. 370–373, 2005.
- [26] Z. Guoying, M. Pietikäinen, Dynamic texture recognition using local binary patterns with an application to facial expressions, IEEE Trans. Pattern Anal. Mach. Intell. Vol. 29 (6), pp.915–928, 2007.
- [27] C. Shan, S. Gong, P.W. McOwan, Facial expression recognition based on local binary patterns: a comprehensive study, Image Vis. Comput. Vol. 27 (6), pp. 803–816, 2009.
- [28] S. Zhang , X. Zhao , B. Lei , "Facial expression recognition based on local binary patterns and local fisher discriminant analysis", WSEAS Trans. Sig. Process. Vol. 8 (1), pp. 21–31, 2012.
- [29] V. Takala, T. Ahonen, M. Pietikäinen, Block-based methods for image retrieval using local binary patterns, in: Proceedings of Scandinavian Conference on Image Analysis, pp. 882–891, 2005.
- [30] O.A.B. Penatti, E. Valle, R.S. Torres, Comparative study of global color and texture descriptors for web image retrieval, J. Vis. Commun. Image Represent. Vol. 23 (2), pp. 359–380, 2012.
- [31] Satpathy, X. Jiang, H. Eng, Lbp based edge texture features for object recognition, IEEE Trans. Image Process. Vol. 23 (5), pp. 1953–1964, 2014.
- [32] Fernández, M. Álvarez, F. Bianconi, Texture description through histograms of equivalent patterns, J. Math. Imaging Vis. Vol. 45 (1), pp. 76–102, 2013.
- [33] L. Nanni, A. Lumini, S. Brahmam, Survey on lbp based texture descriptors for image classification, Expert Syst. Appl. Vol. 39 (3), pp. 3634–3641, 2012.



- [34] Nguyen, P. Ogunbona, W. Li, A novel shape-based non-redundant local binary pattern descriptor for object detection, *Pattern Recognit.* Vol. 46 (5), pp. 1485–1500, 2013.
- [35] S. Manjunath, J.-Ohm, V. V. Vasudevan and A. Yamada, "Color and texture descriptors," in *IEEE Transactions on Circuits and Systems for Video Technology*, Vol. 11, no. 6, pp. 703-715, June 2001.
- [36] D. Huang, S. Caifeng, A. Mohsen, W. Yunhong, C. Liming, Local binary patterns and its application to facial image analysis: a survey, *IEEE Trans. Syst. Man Cybern. Part C* Vol. 41 (6), pp. 765–781, 2011.
- [37] Liu, G. H., Yang, J. Y., & Li, Z. Content-based image retrieval using computational visual attention model. *Pattern Recognition*, Vol. 48(8), pp. 2554–2566, 2015.
- [38] Shakoor, M. H., & Boostani, R. Radial mean local binary pattern for noisy texture classification. *Multimedia Tools and Applications*, Vol. 77(16), pp. 21481–21508, 2018.
- [39] M. Sotoodeh, M.R. Moosavi and R. Boostani, A novel adaptive LBP-based descriptor for color image retrieval, *Expert Systems With Applications*, Vol. 127, pp. 342–352, 2019.
- [40] Johari P.K., Gupta R.K., Retrieval of Content-Based Images by Fuzzified HSV and Local Textural Pattern. In: *Intelligent Computing Applications for Sustainable Real-World Systems. ICSISCET 2019. Proceedings in Adaptation, Learning and Optimization*, Vol. 13, pp. 219-229 Springer, Cham.
- [41] T. Mäenpää, M. Pietikäinen, J. Viertola, Separating color and pattern information for color texture discrimination, in: *Proceedings of 16th International Conference on Pattern Recognition*, Vol.1, pp. 668–671, 2002.
- [42] Rushi Lan, Yicong Zhou, and Yuan Yan Tang, Quaternionic Local Ranking Binary Pattern: A Local Descriptor of Color Images, *IEEE transactions on image processing*, Vol. 25, no. 2, pp. 566-579, 2015.
- [43] S.R. Dubey, S.K. Singh, R.K. Singh, Multichannel decoded local binary patterns for content-based image retrieval, *IEEE Trans. Image Process.* Vol. 25 (9), pp. 4018–4032, 2016.
- [44] G.J. Burghouts, J.-M. Geusebroek, Performance evaluation of local color invariants, *Comput. Vis. Image Understand.* Vol. 113 (1) pp. 48–62, 2009.
- [45] A. Yamada, "Color and texture descriptors," in *IEEE Transactions on Circuits and Systems for Video Technology*, Vol. 11, no. 6, pp. 703-715, 2001.
- [46] Li, C., Huang, Y., & Zhu, L. Color texture image retrieval based on Gaussian copula models of Gabor wavelets. *Pattern Recognition*, Vol. 64, pp. 118–129, 2017.
- [47] T. Mäenpää, M. Pietikäinen, Classification with color and texture: jointly or separately? *Pattern Recognit.* Vol. 37 (8), pp.1629–1640, 2004.
- [48] X. Qi, R. Xiao, C. Li, Y. Qiao, J. Guo, X. Tang, Pairwise rotation invariant cooccurrence local binary pattern, *IEEE Trans. Pattern Anal. Mach. Intell.* Vol. 36 (11), pp. 2199–2213, 2014.
- [49] Liu, Guang-Hai Jing-Yu Yang, Content-based image retrieval using color difference histogram, Vol. 46, Issue 1, pp. 188-198, 2013.
- [50] A. Hafiane, G. Seetharaman, B. Zavidovique, Median binary pattern for textures classification, in: *Proceedings of the 4th International Conference on Image Analysis and Recognition*, pp. 387–398, 2007.
- [51] G-H Liu, J-Y Yang, Content-based image retrieval using color difference histogram, *Pattern Recognit.* Vol. 46 pp. 188–198, 2013.
- [52] T. Deselaers, D. Keysers, H. Ney, Features for image retrieval: an experimental comparison, *Inf. Retrieval* Vol. 11 (2), pp. 77–107, 2008.
- [53] J.Z. Wang, J. Li, G. Wiederhold, Simplicity: semantics-sensitive integrated matching for picture libraries, *IEEE Trans. Pattern Anal. Mach. Intell.* Vol. 23 (9), pp. 947–963, 2001.
- [54] J. Han, K-K. Ma, Rotation-invariant and scale-invariant Gabor features for texture image retrieval, *Image Vis. Comput.* Vol. 25 (9), pp. 1474–1481, 2007.
- [55] G-H Liu, Image Rank Machine Based on Visual Attention Mechanism <<http://www.ci.gxnu.edu.cn/cbir/Dataset.aspx>>.

# A Review of Recommender Systems for Choosing Elective Courses

Mfowabo Maphosa<sup>1</sup>, Wesley Doorsamy<sup>2</sup>, Babu Paul<sup>3</sup>  
Institute for Intelligent Systems, University of Johannesburg  
Johannesburg, South Africa

**Abstract**—In higher education, students face challenges when choosing elective courses in their study programmes. Most higher education institutions employ advisors to assist with this task. Recommender systems have their origins in commerce and are used in other sectors such as education. Recommender systems offer an alternative to the use of human advisors. This paper aims to examine the scope of recommender systems that assist students in choosing elective courses. To achieve this, a systematic literature review (SLR) on recommender systems corpus for choosing elective courses published from 2010–2019 was conducted. Of the 16 981 research articles initially identified, only 24 addressed recommender systems for choosing elective courses and were included in the final analysis. These articles show that several recommender systems approaches and data mining algorithms are used to achieve the task of recommending elective courses. This study identified gaps in current research on the use of recommender systems for choosing elective courses. Further work in several unexplored areas could be examined to enhance the effectiveness of recommender systems for elective courses. This study contributes to the body of literature on recommender systems, in particular those applied for assisting students in choosing elective courses within higher education.

**Keywords**—Recommender systems; elective courses; data mining algorithms; systematic literature review; higher education

## I. INTRODUCTION

Looking through the current lens, in and post COVID-19, it is clear that higher education institutions (HEIs) have to change the way they engage with students from the traditional methods to an online or a blended approach. Popenici and Kerr [1] propose that it is time for HEIs to reimagine their function and pedagogical models in a new paradigm with technology at the centre. This calls for increased application and adaptation of artificial intelligence, machine learning and data mining tools to equip the education sector [2].

Many degree programmes offer elective courses in addition to compulsory ones. The courses that students fail to complete include both compulsory and elective courses. Students chose by a student elective courses based on their interests. Predicting student grades in the courses, they will enroll for is useful for guiding students and allowing them to make informed choices regarding compulsory, and elective courses [3].

In higher education, students are faced with difficulties when choosing elective courses. A survey of first-year students at the University College Dublin showed that almost half of the students selected elective courses outside their

major because they perceived the courses to be exciting. Some of the difficulties emanate from the limited capacity in some elective courses as well as timetable clashes with compulsory courses which make students choose other elective courses [4].

Finding the most suitable elective course from the available ones can be achieved by using recommender systems [5]. By analysing data on the courses that students completed, it is possible to categorize a student's interests. The ability to predict student enrolment patterns for courses provides an opportunity for HEIs to be effective in allocating resources and providing a high-quality learning experience [6]. Predicting student grades in future courses before they take them is an essential tool that can be used to assist students with choosing elective courses [3].

The purpose of recommender systems is to recommend a product to a user that would possibly interest them based on the user profile [7]. A typical recommender system uses three elements: a user, item and rating. The recommender system attempts to predict a rating that a particular user would provide for unrated items [8]. Recommender systems use different types of input data which are placed in a matrix with one dimension representing users and the other one items of interest [9].

The rest of this paper is organised as follows: Section 2 provides a brief literature review; Section 3 discusses the methodology that was followed for the study; Section 4 presents the findings of the study and proposes work that needs to be considered in the field of recommender systems for choosing elective courses; Section 5 discusses the implications of the findings and suggests new trends in the field that can enhance recommender systems and Section 6 summarises the paper.

## II. LITERATURE REVIEW

Laghari [10] warns that poor course selection can cause delays in completing a qualification because students have not completed prerequisite courses, or they have missed the minimum credit requirements for the qualification. Selecting the right elective course is vital for the student to complete their degree programme [11]. Choosing an elective course is influenced by several factors such as the student's personal and academic interest as well as institutional regulations that govern when a particular elective course can be enrolled for [6].

O'Mahony and Smyth [4] identified the following factors that influence students' choices of elective courses: interest and academic goals, career goals, course pre-requisites and co-requisites, ability to progress with their study, difficulty and format, of course, awareness options, availability of places, and timetable clashes. Other factors that influence the choice of elective courses include the number of compulsory and elective courses, the number of credits in a course and the maximum number of students that can be enrolled in a given course [6].

Machine learning tools and techniques have caused disruptive innovation in the way that most industries operate and education has not been spared. Artificial intelligence is defined as "computing systems that are able to engage in human-like processes such as learning, adapting, synthesizing, self-correction and the use of data for complex processing tasks" [1]. Machine learning is considered as a subfield of artificial intelligence. Machine learning involves the development of algorithms used to automatically make sense of data to adapt and learn from experience [12].

According to Gollapudi [7], there are many ways of grouping machine learning algorithms. One such method is the use of model-based grouping. In model-based grouping, machine learning algorithms can be classified into one of the following classes: association rule-based, Bayesian methods based, clustering methods based, deep learning-based, decision tree-based, dimensionality reduction based, ensemble method based, instance-based, kernel methods based and regression analysis based.

Data mining is an umbrella term for two separate processes: knowledge discovery and prediction. Knowledge discovery entails providing information in a form that can be understood by end-users and prediction allows the foretelling of future events [13]. Machine learning and data mining are different in that machine learning focus on using general knowledge, while data mining focuses on discovering new knowledge [7]. A sub-field of data mining with a focus on applying data mining tools and techniques is educational data mining (EDM). EDM can be defined as a process of applying computerised methods to identify patterns in educational data that are hard to detect because of the volume of data [14].

Recommender system techniques use different classifications based on the data that are used as input for the recommendation. There are four broad classifications of recommender systems: collaborative filtering (CF), content-based, knowledge-based and hybrid approaches [8] [12] [15].

#### A. Collaborative Filtering

CF is the most widely used recommendation technique because of its power and simplicity [12]. CF uses data about users and items. A recommendation is made by analysing relationships between users and interdependencies among items to identify new user-item associations [16]. CF techniques can be classified as memory-based and model-based. Memory-based refers to the use of user-based algorithms and item-based algorithms. User-based algorithms produce predictions for a given user by first identifying users with similar choices to the given user and then calculating the

most frequently rated items that the given user has not seen [17]. In model-based algorithms, models are used to predict the ratings of unrated items by learning intricate patterns based on training data and using these patterns to make predictions [12].

#### B. Content-based Filtering

Content-based filtering techniques are based on the idea that users will prefer items that are similar to items that previously offered them enjoyment [12]. Content-based filtering depends on data about users and categories that have been assigned to the available item descriptions [16]. Content-based filtering allows for the creation of a profile for each user to characterize its nature. This enables an association between a user and matching categories to be made by calculating a set of items that are most similar to items already known by the current user [9] [16]. Content-based filtering has an architecture that consists of components, such as item representations, user-profiles and the ability to learn a user model [12].

#### C. Knowledge-based Filtering

Knowledge-based filtering use domain knowledge to generate recommendations. This knowledge is made up of rules, metrics and items. Depending on the given user requirements, rules are used to describe the best approach to use to make a recommendation [16].

#### D. Hybrid Techniques

Hybrid techniques combine the above-discussed approaches to create a unified model that possesses characteristics of all approaches. The use of the unified model helps to mitigate certain limitations of the above approaches [12]. Hybrid techniques are a common feature because they provide opportunities to achieve better accuracy than the techniques mentioned above [16].

#### E. Matrix Factorisation

Matrix factorisation (MF) techniques can overcome the problem of data sparsity by employing dimensionality reduction to improve the model's ability to generalise [18]. MF can be used within CF recommender systems to achieve better levels of accuracy than those achieved by nearest neighbor techniques [9]. There are a variety of matrix factorization models and combinations in use today. These include singular value decomposition, PMF, non-negative MF, probabilistic sparse MF, Bayesian probabilistic MF and general probabilistic MF [18].

This paper aims to examine the scope of recommender systems for choosing elective courses. The study seeks to survey the landscape and determine the state of recommender systems for elective courses in higher education and to identify emerging technologies that could be explored to enhance recommender systems. To achieve this objective, a review of relevant literature on recommender systems for recommending elective courses was conducted.

### III. METHODOLOGY

Kitchenham and Charters [19] define an SLR as "a means of evaluating and interpreting all available research relevant to

a particular research question, topic area, or phenomenon of interest. Systematic reviews aim to present a fair evaluation of a research topic by using a trustworthy, rigorous, and auditable methodology". An SLR involves analysing relevant primary research studies by identifying, evaluating and interpreting corpus.

For this study, the SLR method proposed by guidelines of Kitchenham and Charters [19] is used together with the Preferred Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA) principles for the selection of the articles [20]. This section outlines the three steps, followed in performing the SLR study. In the first step, the need for performing the SLR study is identified. Then, a review protocol for conducting the SLR study is developed. In the third step, the process followed in conducting the SLR is described. The reporting of the findings of the study is described in Section 4: Results.

#### A. The Need for a Systematic Literature Review

Recommender systems are gaining prominence in the education sector. The search of the literature revealed that researchers had conducted systematic literature reviews (SLR) on recommender systems. Iatrellis, Kameas and Fitsilis [21] conducted an SLR study on academic advising systems and its impact on education. This study covered work published between 2008 and 2017. The study found that academic advising systems were used for choosing programs/majors, selecting courses and long-term academic planning.

Rivera, Tapia-Leon and Lujan-Mora [22] conducted a literature review on recommender systems in education. The study revealed that recommender systems are used to address different challenges in education with the majority of studies focusing on academic choice. During the review, no SLR studies of recommender systems for elective courses were found. Thus the motivation for this study is to identify, evaluate and analyze relevant literature on recommender systems that recommend elective courses in higher education.

#### B. Development of a Review Protocol

Kitchenham [23] states that a review protocol is essential as it defines the method that will be used to undertake the study. The following steps describe the development process for the review protocol for the study at hand.

1) *Identify the research goals and research questions:* The goal of this study was to conduct an SLR in selecting elective courses in higher education. The research sought to answer the following questions.

a) *Research question 1:* What is the state of recommender systems for elective courses?

b) *Research question 2:* What emerging trends in data mining should be explored to enhance recommender systems for elective courses?

2) *Identify keywords:* The literature search terms comprised of the following words and combinations:

"recommender systems", and "recommendation systems". Kitchenham and Charters [19] proposed the use of Boolean operators such as "AND" and "OR" for refining the keyword search string. For this study, the logical operator 'OR' is used to join the identified keywords, and the 'AND' operator is used to combine the keywords in the phrase. The study used the following search string: [{"recommender systems" OR "recommendation systems" AND "elective courses"}].

3) *Identify the sources:* Specific online databases and search engines for were searched for research articles related to recommender systems in higher education with a focus on assisting students in selecting elective courses. These included IEEE Xplore, ACM Digital Library, Science Direct, Emerald Insight, EBSCOhost and Google Scholar. The authors used these online databases and search engines because they assumed that these were the main sources for collecting relevant literature on recommender systems.

4) *Identify the inclusion criteria:* The inclusion criteria for this study are as follows:

a) Articles that satisfied the keyword conditions.

b) Articles that are written in English.

c) Articles published between 2010 and 2019.

d) The articles focused on selecting elective courses.

5) *Study quality assessment:* Kitchenham [23] asserts that it is critical to assess the quality of primary articles. The following questions were used to measure the quality of the articles to be included in the final list.

a) Is the research article focussed on recommender systems for recommending elective courses?

b) Is the research article a primary study?

6) *Identify the data extraction strategy:* The data extraction strategy involved extracting the following information from each research article: author and year, name of journal or conference proceeding, the objective of the study, the size of dataset used, the recommender system approach employed, the data mining algorithm used and the results of the study.

#### C. Conducting the Review

In the first phase, the search string was applied to the online databases and search engines. The search string was applied on all metadata and obtained 16 981 research articles, as shown in Fig. 1.

Next, refined the search was refined to only the article titles, which yielded to 3 021 papers. Filters were then applied on the online databases and search engines to exclude articles not written in English, non-peer-reviewed articles, articles not published in journals and conference proceedings and articles that are not full access. A total of 2 897 articles were excluded leaving 124 articles that were analyzed for the subject matter, leading to further exclusion of 28 articles.

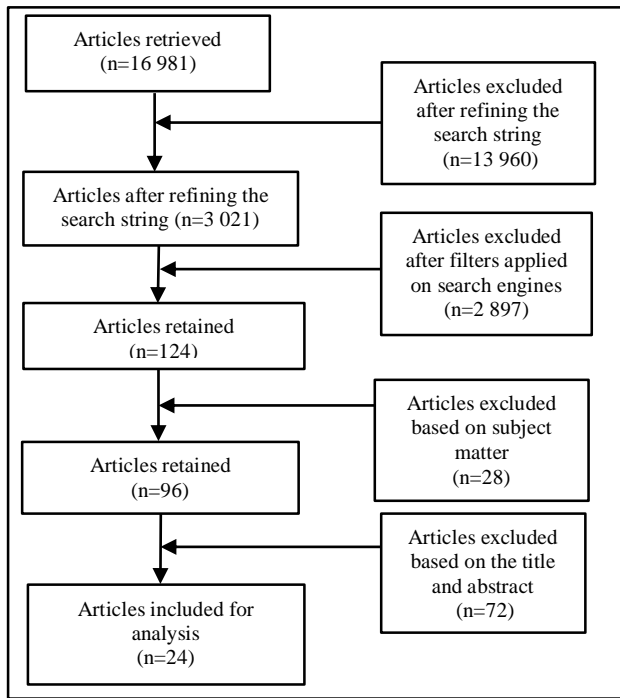


Fig. 1. Article Selection Process Adapted from the PRISMA Flowchart [20].

The references and abstracts for the remaining 96 articles were then uploaded to Rayyan (<https://rayyan.qcri.org>), a free web and mobile app for screening articles when conducting an SLR study. According to Ouzzani, Hammady, Fedorowicz and Elmagarmid [24], Rayyan is used to speed up the initial abstract/title screening of articles and also allows researchers to collaborate when performing SLR studies. It is easy to pick up duplicates in articles referenced in more than one database when using Rayyan. The titles and abstracts of the 96 articles were then analyzed on Rayyan and excluded 72 articles leaving 24 articles for analysis. These articles passed the quality assessment by having yes as the answer for both the quality assessment questions.

The 24 articles included for the final analysis were then uploaded into an online SLR software – SysRev (<https://sysrev.com>) to extract the necessary data. The following labels were created to extract data from each article – the year of publication, name of journal or conference, country, size of data set used for testing, recommender system technique used and data mining method used. Information about the objective of the article and the results obtained were also extracted. Table I contains the quality assessment results for the 24 reviewed articles.

TABLE I. LIST OF REVIEWED ARTICLES

| Ref  | Journal   | The objective of the study   | Size of the dataset used           | Recommender System Approaches used | Data mining Algorithms used                       | Results  |
|------|---|--|------------------------------------|------------------------------------|---|--|
| [25] | International Journal of Information Intelligence Systems, Technology and Management  | To develop a course recommendation system to help students choose elective courses.        | 255 student records, 25 courses    | Not provided                       | Item and user-based CF                            | Both item and user-based CF methods achieved high prediction accuracy.   |
| [26] | Athens Journal of Sciences  | To recommend elective courses.   | Not provided                       | Fuzzy clustering                   | Gustafson-Kessel clustering                       | The use of several clustering algorithms on students' data provided better results.  |
| [27] | 2015 IEEE Frontiers in Education Conference   | To develop a web-based recommender system that uses CF.                                    | 743 student records, 50 courses    | CF                                 | K-means algorithm clustering                      | The proposed system uses an intelligent advising component to provide a rough guideline to students on course selection and selection of majors. |
| [28] | 2016 IEEE International Conference on Big Data  | To recommend elective courses to students each semester based on courses taken previously. | 37 392 student records             | Markov-based CF                    | Not provided                                      | The Skip Markov Model performed better than the other recommender models.  |
| [8]  | Proceeding of 8th International Conference on Knowledge and Systems Engineering (KSE) | To create a framework for building a course recommendation system.                         | 4 017 student records, 353 courses | Biased MF                          | k-NN  | The system allows students to choose elective courses, and the system would recommend the best courses.  |
| [29] | Proceeding of 9th International Conference on Educational Data Mining                 | To recommend core and elective courses   | 1 444 student records              | Not provided                       | Four custom algorithms                            | The systems could warn students about challenging courses and recommend courses a student could benefit from taking.                             |
| [30] | International Journal of Advanced Computer Science and Applications                   | To show students the available elective courses and make recommendations.                  | 2 000 student records, 54 courses  | CF                                 | K-means Clustering, Association rule mining (ARM) | ARM can be used to recommend courses to a target student. The model achieved the highest precision rate of 90%.                                  |

|      |  |  |  |                                       |  |  |
|------|--|--|--|---------------------------------------|--|--|
| [31] | Procedia Computer Science  | To develop a hybrid recommender system to recommend courses.   | 300 courses, 5 programmes,               | Ontology modelling, Hybrid techniques | Classification   | The hybrid technique was much more effective in terms of accuracy.   |
| [32] | IEEE Transactions on Signal Processing   | To recommend courses adaptively based on the students' background.   | 1 444 student records                    | Multi-armed Bandits (contextual)      | Forward-Search Backward-Induction Algorithm                | The recommender system outperforms systems that ignore personalized context information.   |
| [11] | International Journal for Research in Applied Science and Engineering and Technology (IJRASET)       | To recommend elective subjects based on neural networks and association rule   | 250 student records, 22 courses          | CF                                    | Artificial neural networks (ANN) Multilayer Perceptron ARM | The recommender system predicts elective subjects for students based on their marks from the previous semester.  |
| [33] | International Journal of Information Technology and Computer Science                                 | To recommend courses to students based on their profile on Moodle  | 100 student records, 6 courses           | Not provided                          | K-means algorithm clustering                               | Based on the performance of one course, the study was able to recommend the most appropriate elective courses.   |
| [34] | Cybernetics and Information Technologies   | To predict student performance in courses using a recommender-based approach and a regression-based approach.          | 1268 student records                     | MF                                    | ANN, decision tree, SVM, logistic regression               | The regression-based approach performed better than the recommender-based approach.  |
| [35] | Procedia Computer Science  | To recommend elective courses based on previous grades.  | 1 000 student records                    | Not provided                          | ARM  | The system was tested on 100 students and achieved an efficiency of 90%.   |
| [5]  | Proceeding of 14th International Conference (Lecture Notes in Computer Science)                      | To broaden the range of elective courses that students are aware of by adding diversity to the recommendation process. | 100 student records                      | Hybrid approach                       | Vector Space Model Content-based                           | The system improves recommendation diversity than content-based and hierarchical taxonomy systems as module descriptions make recommendations more meaningful. |
| [36] | Proceeding of 11th International Conference on Educational Data Mining                               | To recommend elective courses based on course orderings and grade predictions.   | 1 700 student records, 72 courses        | Context-aware filtering               | Not provided   | Therefore, the course dependency graph seems to be more suitable for course recommendations.   |
| [37] | International Journal of Scientific Research in Computer Science Applications and Management Studies | To predict final grades for students and recommend elective subjects   | Not provided                             | Hybrid techniques                     | Ensemble (Pearson Algorithm and I to I)                    | The system recommends elective courses to the student to yield maximum grade.  |
| [38] | International Journal on Future Revolution in Computer Science and Communication Engineering         | To build a recommender system to recommend career paths for undergraduate students                                     | Not provided                             | Rule-based learning system            | k-NN   | The new dataset was used to evaluate the system with an accuracy of 75%.   |
| [39] | 11th International Conference on Educational Data Mining   | To determine the most relevant criteria for recommending courses.  | 1700 course ratings (survey), 63 courses | CF                                    | Not provided   | The study used different weights for each criterion to use the combination of multiple criteria which provided better results.                                 |
| [40] | International Journal of AI and Data Mining  | To design a course recommender model to assist decision-making for elective course selection.                          | 798 student records                      | CF                                    | Clustering, Fuzzy Association Rule                         | The system could recommend appropriate elective courses and predict the likely students' grade.  |
| [41] | International Journal of Data Science and Analysis   | To assist students in choosing the most appropriate elective courses for better performance dynamically.               | 10 601 student records                   | Knowledge-based                       | k-NN   | The results of these calculations prove that the model has a high level of accuracy. The model achieved an accuracy rate of 95.6%.                             |
| [42] | International Journal for Research Trends and Innovation   | To predict student performance and to recommend elective courses   | 16 features                              | CF                                    | MF, probabilistic MF (PMF) Gene Fuzzy model                | The system classifies students into one of three categories – theory, testing and practical so that student can know what to focus on the following semester.  |

|      |   |   |                                     |            |                           |  |
|------|---|---|-------------------------------------|------------|---------------------------|--|
| [43] | Proceedings of 9th International Conference on Learning Analytics and Knowledge | To personalize course pre-requisite inference for a goal-based recommendation based on adaptations of a recurrent neural network. | 164 196 student records, 10 courses | Goal-based | Recurrent Neural Networks | The study shows that recommendation was set up to recommend course preparation for a single semester.  |
| [44] | International Journal of Computer Science and Information Technology            | To recommend courses based on what other students have taken after applying the association rules algorithm on course data.       | 384 student records                 | ARM        | k-NN                      | ARM can be used to recommend elective courses for students. The system did not provide recommendations when the matching rule was less than 50%. |
| [45] | Journal of Theoretical and Applied Information Technology                       | To use a context-aware recommender system which recommends undergraduate programs to students based on academic performance.      | 3 421 student records               | CF         | Naïve Bayes, J48          | The use of a rating matrix is useful when using contextual information. The effectiveness of the recommender system was 98%.                     |

#### IV. RESULTS

In this study, an SLR is performed to ascertain the state of recommender systems for choosing elective courses and to identify emerging technologies for recommender systems. After the screening process using Rayyan, 24 primary studies were selected. Fig. 2 shows the topics and themes covered in the 24 articles reviewed. The size and the frequency of the topic or theme show its prevalence in the articles. The results of this SLR are structured according to the two research questions.

##### A. The State of Recommender Systems for Elective Courses

1) *Publications per year:* None of the reviewed articles was published in the years 2010 and between 2012 and 2014. The majority of the articles were published in the years 2016 to 2019, as shown in Table I. The steady increase in the number of articles in the years 2016–2019 shows that there is more interest in recommender systems for recommending elective courses from researchers. Fig. 3 shows the number of articles published each year.

2) *Publications type:* All 24 reviewed articles were published in different journals and conference proceedings. 67% of these articles were published in journals with the remainder being published in conference proceedings, as shown in Fig. 4. Interestingly, all the 24 reviewed articles were published in different conference proceedings, with only two articles being published by the same journal [31] [35]. The fact that research on recommender systems on elective courses is published in both journals and conference proceedings highlights the interest this subject has among researchers.

3) *Objectives of the study:* The analysis of the articles shows that there are a variety of ways of performing the task of recommending elective courses. These include making recommendations based on the student’s background and marks, broadening the range of elective courses available to students, and providing descriptions of the elective courses other than just the name of the course. The variety of ways of performing the task of recommending elective courses has increased over the years showing that this field is growing.

4) *Dataset used:* In terms of the datasets used, studies that used datasets that consisted of students and courses accounting for 37.5%. Another 37.5% of the articles used a dataset consisting of students and course data. Articles that used a

dataset consisting, of course data without providing student details accounted for 8.3%, and one article (4.2%) used 16 features which were not specified to be either relating to students or courses. Lastly, three articles, representing 12.5%, did not specify the size of the dataset used.

The analysis of the articles revealed that recommending elective courses can be done by using recommender system techniques, data mining techniques or both. 70.8% of the articles reviewed used a combination of recommender system techniques and data mining techniques. 16.7% used data mining techniques to make a recommendation, and 12.5% used recommender system techniques to make a recommendation. It was also interesting to note that all the reviewed articles published in 2019 utilised both recommender system techniques and data mining techniques.

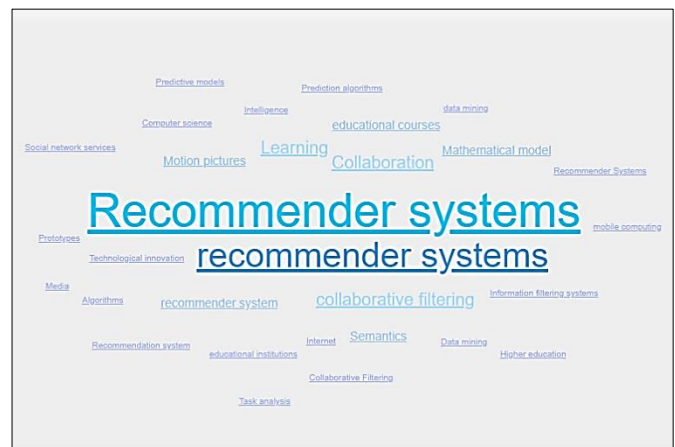


Fig. 2. Topics and Themes Covered in the Reviewed Articles.

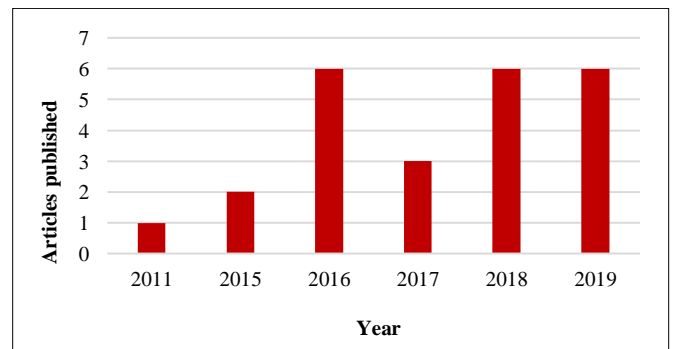


Fig. 3. Articles Published Per Year.

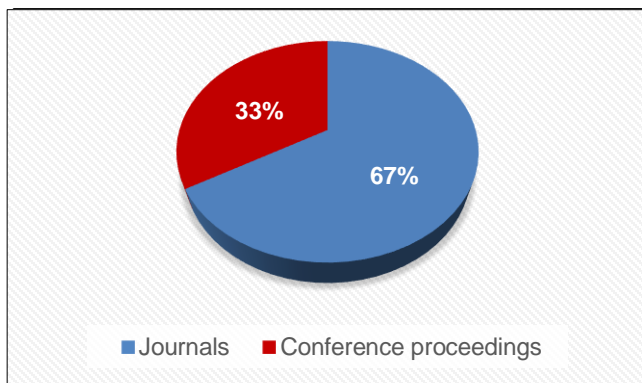


Fig. 4. Publication Type.

5) *Recommender system approaches*: The analysis of the results shows that the most widely used recommender system technique is CF. of the 20 papers that utilised recommender systems, 45% used CF, 10% used the hybrid techniques, and 10% used MF. The rest of the articles used, ARM, fuzzy clustering, knowledge-based, rule-based and multi-armed bandits (contextual). Also, the results show that different and new techniques are being applied to recommender systems. These techniques include ontology modelling, context-aware filtering and goal-based.

6) *Data mining algorithms*: This review revealed that several data mining approaches were applied in the reviewed articles. These approaches include k-nearest neighbors (k-NN), k-means algorithm, MF, ARM, recurrent neural networks, multilayer perceptron, item and user-based CF, and ANN. Interestingly, some articles used custom algorithms.

The first research question was to show the state of recommender systems for recommending elective courses. This section has highlighted the state of recommender systems for recommending elective courses using these sections: publications per year, publications type, objectives of the study, size of dataset used, recommender system approaches and data mining algorithms used. The analysis shows that recommender systems are useful for recommending elective courses. In the research articles where empirical evidence is given, the efficiency, precision or accuracy rates are greater than 90%, indicating the effectiveness of these recommender systems.

#### B. Emerging Technologies for Recommender Systems

The second research question for the study sought to establish the emerging trends in data mining that can be explored to enhance recommender systems for elective courses. The following paragraphs discuss the emerging trends based on the analysis of the reviewed articles.

1) Most studies on recommending elective courses suffer from limitations related to the use of structured data. This SLR study has shown that several unexplored areas could enhance the effectiveness of recommender systems for elective courses.

2) Several studies reported on the fact that some students did not take the recommended elective courses. This could be to do with acceptance or lack-of for recommendations drawn

from recommender systems. Researchers could attempt to incorporate models to improve the acceptance of recommender systems by students, which may result in them choosing the recommended elective courses.

3) Another critical area that needs to be considered is the selection of the right programme or discipline before recommending the correct elective course. Recommender systems should start with guiding students on the qualification path based on their interests.

4) Schnabel, Bennett and Joachims [46] suggest that user feedback can increase the learning accuracy of recommender systems. Using the information foraging theory, the authors prove that foraging interventions are complementary to improving algorithms and result in more effective recommender systems. There is a need to consider allowing students to provide feedback on the recommendations made by the recommender systems. Such feedback could be solicited at the beginning of the semester and retrospectively at the end of the semester.

#### V. DISCUSSION

This study reports an SLR regarding recommender systems for elective courses. This study aimed to ascertain the state of recommender systems for elective courses. It also sought to establish the emerging trends in data mining that can be explored to enhance recommender systems that assist students in choosing elective courses. This study reviewed 24 articles on the corpus and reported the results using different themes. The results showed that research on recommender systems has been increasing with the majority of the articles published in journals. Also, recommender systems are used to address a myriad of challenges faced by HEIs. This finding is not surprising given the impact recommender systems are having in other fields such as commerce, medicine and entertainment. Currently, HEIs have at their disposal vast amounts of data, both structured data and unstructured data obtained from social media [12]. There are possibilities for the task of recommending elective courses to incorporate structured data and unstructured data. Another challenge that most recommender systems reviewed face is the assumption that past student performance is a determinant of future performance. This is not always the case as social factors could have influenced past performance. It could be useful if recommender systems could incorporate data from social media.

The results showed that a variety of datasets were used – datasets on courses, students and a combination of both. The size of the datasets ranges from small to large, with data sourced from institutional repositories and other data sourced from survey questionnaires. The analysis of these articles shows that there are papers focused on recommending elective courses and others focus on recommending courses and predicting grades that the student would likely obtain should they choose the recommended elective course. Thus, this development is considered vital as it addresses one of the challenges grappling higher education–student’s poor performance.



CF was the most widely used recommender system technique. Furthermore, several data mining algorithms have been employed to address the challenge of recommending elective courses. The variety of data mining algorithms may, in part, be attributable to the breakthroughs in terms of big data, data analysis and data science. This is in line with the findings of Jembere, Rawatlal and Pillay [47] that EDM provides numerous prediction tools that can be used to guide students on choosing courses.

Recommender systems must not only focus on broad outcomes such as courses but also on recommending learning resources and activities that will assist students in passing the recommended course. Such recommendations can take into consideration the student's needs, interests, preferences and past activities [14]. Fourthly, none of the reviewed articles touched on the aspect of Big Data. It is crucial for recommender systems to incorporate Big Data techniques as HEIs are producing vast amounts of Big Data. Higher education is changing post-COVID-19, and part of the changes will include online learning, which will produce more Big Data for analysis. This is in line with predictions by Popenici and Kerr [1] that HEIs need to reimagine their function and pedagogical models in a new paradigm with technology at the centre.

The analysis of the recommendations of the articles highlighted the direction that recommender systems for elective courses are taking. Only a handful of the reviewed articles utilised ensemble methods. Ensemble methods combine many models that are built independently to use the combined models to make predictions. The individual models that are combined are referred to as weaker models because the results of these models cannot make the required task on their own [7]. Two standard ensemble methods are bagging and boosting. Bagging (bootstrap aggregating) as a concept was introduced by Breiman [48] to reduce the variance of a predictor. Bagging is a simple ensemble method in which many independent models are built and combined using some model averaging techniques. Bagging is akin to improving existing methods by adding a loop in front that selects the bootstrap sample [48]. Boosting is an ensemble method in which the models are not made independently, but sequentially [12]. Ensemble methods provide an opportunity to build algorithms to recommend elective courses and algorithms to predict the grades that a student is likely to achieve for the recommended course.

Deep learning is defined as "a class of machine learning techniques that exploit many layers of nonlinear information processing for supervised or unsupervised feature extraction and transformation, and for pattern analysis and classification" [12]. Deep learning is an area of machine learning that deals with building more complex neural networks to solve problems classified under semi-supervised learning and operates on datasets that have little labelled data. Some of the widely used deep learning techniques include convolutional networks, restricted Boltzmann machine, deep belief networks and stacked auto-encoders [7]. Deep learning can be employed in a bid to create more accurate user profiles. These user profiles are central to the problem of recommending elective courses.

## VI. CONCLUSION

Recommender systems employed to recommend elective courses to students are gaining traction. This growth can be attributed to the rise in the effectiveness of recommender systems that recommend products and services in sectors such as commerce and entertainment. In this paper, 24 articles on recommender systems aimed at recommending elective courses to students in higher education are reviewed. This review offers some insight into the state of recommender systems in this domain. Through this SLR, the recommender systems techniques and the data mining methods used in these papers to make recommendations were identified. The review revealed that several recommender systems approach and data mining algorithms are used to achieve the task of recommending elective courses. More importantly, this study has suggested emerging trends in the field that need to be explored by recommender systems to improve their effectiveness. These include the incorporation of acceptance models to increase the acceptance of recommendations, the effectiveness of user feedback. There is also a need to consider recommendation systems that begin with recommending the qualification path. This review is useful as it summarises current trends and makes suggestions on the future of this field of recommender systems for recommending elective models.

## ACKNOWLEDGMENT

MM conceived the concept, performed the literature review and drafted the original article. WD performed critical revision and approval of the submitted version. BP performed critical revision and approval of the submitted version. All authors read and approved the final manuscript.

## REFERENCES

- [1] S. A. D. Popenici and S. Kerr, "Exploring the impact of artificial intelligence on teaching and learning in higher education," *Research and Practice in Technology Enhanced Learning*, vol. 12, no. 22, pp. 1-13, 2017.
- [2] S. A. Becker, B. Malcolm, E. Dahlstrom, A. Davis, K. DePaul, V. Diaz and J. Pomerantz, "NMC Horizon Report: 2018 Higher Education Edition," EDUCAUSE, Louisville, 2018.
- [3] A. Cakmak, "Predicting student success in courses via collaborative filtering," *International Journal of Intelligent Systems and Applications in Engineering*, vol. 5, no. 1, pp. 10-17, 2017.
- [4] M. P. O'Mahony and B. Smyth, "A Recommender system for on-line course enrolment: An initial study," in 2007 ACM Conference on Recommender Systems, Minneapolis, 2007.
- [5] N. Hagemann, M. P. O'Mahony and B. Smyth, "Module advisor: Guiding students with recommendations," in *International Conference on Intelligent Tutoring Systems*, Cham, 2018.
- [6] I. Ognjanovic, D. Gasevic and S. Dawson, "Using institutional data to predict course selections in higher education," *The Internet and Higher Education*, vol. 29, pp. 49-62, 2016.
- [7] S. Gollapudi, *Practical Machine Learning*, Birmingham: Packt Publishing, 2016.
- [8] H. Thanh-Nhan, H. Nguyen and N. Thai-Nghe, "Methods for building recommender systems," in 2016 Eighth International Conference on Knowledge and Systems Engineering (KSE), Hanoi, 2016.
- [9] Y. Koren, R. Bell and C. Volinsky, "Matrix factorization techniques for recommender systems," *Computer*, vol. 42, no. 8, pp. 30-37, 2009.
- [10] M. S. Laghari, "Automated course advising system," *International Journal of Machine Learning and Computing*, vol. 4, no. 1, pp. 47-51, 2014.

- [11] N. B. Samrit and A. Thomas, "A recommendation system for prediction of elective subjects," *International Journal of Research in Applied Science and Engineering Technology (IJRAET)*, vol. 5, no. 4, pp. 36-43, 2017.
- [12] R. Akerkar, *Artificial Intelligence for Business*, eBook: Springer, 2019.
- [13] A. Kunjir, P. Pardeshi, S. Doshi and K. Naik, "Recommendation of data mining technique in higher education," *International Journal of Computational Engineering Research (IJCER)*, vol. 5, no. 3, pp. 29-34, 2015.
- [14] O. Scheuer and B. M. McLaren, "Educational Data Mining," in *Encyclopedia of the Sciences of Learning*, Springer, 2011.
- [15] S. Dwivedi and K. Roshni, "Recommender systems for big data in education," in *5th National Conference on E-Learning & E-Learning Technologies (ELELTECH)*, Hyderabad, 2017.
- [16] A. Felfernig, M. Jeran, G. Ninaus, F. Reinfrank, S. Reiterer and M. Stettinger, "Basic approaches in recommendation systems," in *Recommendation Systems in Software Engineering*, M. P. Robillard, W. Maalej, R. J. Walker and T. Zimmermann, Eds., Berlin, Springer, 2014, pp. 15-37.
- [17] M. Skilton and F. Hovsepian, *The 4th Industrial Revolution: Responding to the Impact of Artificial Intelligence on Business*, Cham: Springer, 2018.
- [18] Z. Zhang, Y. Lin and Z. Zhang, "Field-aware matrix factorization for recommender systems," *Open Access Journal - Special Section on Recent Computational Methods in Knowledge Engineering and Intelligence Computation*, vol. 6, pp. 45690-45698, 2018.
- [19] B. Kitchenham and S. Charters, "Guidelines for performing systematic literature reviews in software engineering version 2.3," *Engineering*, vol. 45, no. 4ve, p. 1051, 2007.
- [20] A. Liberati, D. G. Altman, J. Tetzlaff, C. Mulrow, P. C. Gøtzsche, J. P. Ioannidis, M. Clarke, P. J. Devereaux, J. Kleijnen and D. Moher, "The PRISMA statement for reporting systematic reviews and meta-analyses of studies that evaluate health care interventions: Explanation and elaboration," *Journal of Clinical Epidemiology*, vol. 62, no. 10, pp. 1-34, 2009.
- [21] O. Iatrellis, A. Kameas and P. Fitsilis, "Academic advising systems: A systematic literature review of empirical evidence," *Education Sciences*, vol. 7, no. 4, pp. 90-107, 2017.
- [22] A. C. Rivera, M. Tapia-Leon and S. Lujan-Mora, "Recommendation systems in education: A systematic mapping study," in *Proceedings of International Conference on Information Theoretic Security*, Cham, 2018.
- [23] B. Kitchenham, "Procedures for performing systematic reviews," *Keele University*, Keele, 2004.
- [24] M. Ouzzani, H. Hammady, Z. Fedorowicz and A. Elmagarmid, "Rayyan — a web and mobile app for systematic reviews. *Systematic Reviews*, vol. 5, no. 1, pp. 210-220, 2016.
- [25] S. Ray and A. Sharma, "A collaborative filtering based approach for recommending elective courses," in *Proceedings of International Conference on Information Intelligence, Systems, Technology and Management*, Berlin, 2011.
- [26] E. Bedalli and I. Ninka, "Exploring an educational system's data through fuzzy cluster analysis," *Athens Journal of Sciences*, vol. 2, no. 1, pp. 33-44, 2015.
- [27] K. Ganeshan and X. Li, "An intelligent student advising system using collaborative filtering," in *2015 IEEE Frontiers in Education Conference*, El Paso, 2015.
- [28] E. S. Khorasani, Z. Zhenge and J. Champaign, "A Markov chain collaborative filtering model for course enrollment recommendations," in *2016 IEEE International Conference on Big Data*, Washington, 2016.
- [29] H. Bydzovska, "Course enrollment recommender system," in *The 9th International Conference on Educational Data Mining EDM 2016*, Raleigh, 2016.
- [30] A. Al-Badarenah and J. Alsakran, "An automated recommender system for course selection," *International Journal of Advanced Computer Science and Applications*, vol. 7, no. 3, pp. 166-175, 2016.
- [31] D. Upendran, S. Chatterjee, S. Sindhumol and K. Bijlani, "Application of predictive analytics in intelligent course recommendation," *Procedia Computer Science*, vol. 93, pp. 917-923, 2016.
- [32] J. Xu, T. Xing and M. van der Schaar, "Personalized course sequence recommendations," *IEEE Transactions on Signal Processing*, vol. 64, no. 20, pp. 5340-5352, 2016.
- [33] B. Rawat and S. K. Dwivedi, "An architecture for recommendation of courses in e-learning systems," *International Journal of Information Technology and Computer Science (IJITCS)*, vol. 9, no. 4, pp. 39-47, 2017.
- [34] T. O. Tran, H. T. Dang, V. T. Dinh, T. M. Truong, T. P. Vuong and X. H. Phan, "Performance prediction for students: A multi-strategy approach," *Cybernetics and Information Technologies*, vol. 17, no. 2, pp. 164-182, 2017.
- [35] Z. Gulzar, A. A. Leema and G. Deepak, "PRCS: Personalized course recommender system based on hybrid approach," *Procedia Computer Science*, vol. 125, pp. 518-524, 2018.
- [36] M. Backenkohler, F. Scherzinger, A. Singla and V. Wolf, "Data-driven approach towards a personalised curriculum," *Buffalo*, 2018.
- [37] S. D. Tupe, "A student performance prediction and course recommendation system: A survey," *International Journal of Scientific Research in Computer Science Applications and Management Studies*, vol. 7, no. 5, 2018.
- [38] N. Sawarkar, M. M. Raghuvanshi and K. R. Singh, "Intelligent recommendation system for higher education," *International Journal on Future Revolution in Computer Science & Communication Engineering*, vol. 4, no. 4, pp. 311-320, 2018.
- [39] A. Esteban, A. Zafra and C. Romero, "A hybrid multi-criteria approach using a genetic algorithm for recommending courses to university students," in *11th International Conference on Educational Data Mining*, New York, 2018.
- [40] S. Asadi, S. M. Jafari and Z. Shokrollahi, "Developing a course recommendation by combining clustering and fuzzy association rule," *Journal of AI and Data Mining*, vol. 7, no. 2, pp. 249-262, 2019.
- [41] A. O. Ogunde and E. Ajibade, "A k-nearest neighbour algorithm-based recommender system for the dynamic selection of elective undergraduate courses," *International Journal of Data Science and Analysis*, vol. 5, no. 6, pp. 128-135, 2019.
- [42] T. S. Deepak, "A student performance prediction and course recommendation system," *International Journal for Research Trends and Innovation*, vol. 4, no. 8, pp. 67-71, 2019.
- [43] W. Jiang, Z. A. Pardos and Q. Wei, "Goal-based course recommendation," in *Proceedings of 9th International Conference on Learning Analytics and Knowledge (LAK'19)*, Tempe, 2019.
- [44] W. A. AlZoubi, "Cluster based association rule mining for courses recommendation system," *International Journal of Computer Science & Information Technology (IJCSIT)*, vol. 11, no. 6, pp. 13-19, 2019.
- [45] V. Vaidhehi and R. Suchithra, "Design of a context-aware recommender systems for undergraduate program recommendations," *Journal of Theoretical and Applied Information Technology*, vol. 97, no. 23, pp. 3583-3596, 2019.
- [46] T. Schnabel, P. N. Bennett and T. Joachims, "Improving recommender systems beyond the algorithm," *arXiv preprint arXiv*, 2018.
- [47] E. Jembere, R. Rawatlal and A. Pillay, "Matrix factorisation for predicting student performance," in *Proceedings of 7th World Engineering Education Forum (WEEF)*, Kuala Lumpur, 2018.
- [48] L. Breiman, "Bagging predictors," *Machine Learning*, vol. 24, pp. 123-140, 1996.

# Susceptible, Infectious and Recovered (SIR Model) Predictive Model to Understand the Key Factors of COVID-19 Transmission

DeepaRani Gopagoni<sup>1</sup>, P V Lakshmi<sup>2</sup>  
Computer Science and Engineering  
GITAM Institute of Technology  
Vishakhapatnam, Andhra Pradesh  
India

**Abstract**—On 31 December 2019, WHO was alerted to several cases of pneumonia in Wuhan City, Hubei Province of China. The virus did not match any other known virus. This raised concern because when a virus is new, general behavior and how it affects, people do not know. Initial few cases reportedly had some link to a large seafood and animal market, suggesting animal-to-person spread. However, a growing number of patients reportedly have not had exposure to animal markets, indicating person-to-person spread is occurring. At this time, it's unclear how easily or sustainably this virus is spreading between people. At any given time during a flu epidemic, firstly, should know the number of people who are infected. Second, to know the numbers who have been infected and have recovered, because these people now have immunity to the disease. Well established SIR modeling methodology is used to develop a predictive model in order to understand the key factors that impact the COVID-19 transmission.

**Keywords**—COVID-19; SIR modeling; WHO; disease spread

## I. INTRODUCTION

In December 2019 World Health Organization alerted to several cases of pneumonia in Wuhan City, Hubei Province of China [1]. The virus did not match any other known virus. Novel Corona virus (2019-nCoV) is a virus (more specifically, a corona virus) identified as the cause of an outbreak of respiratory illness. Early on, many of the patients in the outbreak in Wuhan, China reportedly had some link to a large seafood and animal market, suggesting animal-to-person spread. However, a growing number of patients reportedly have not had exposure to animal markets, indicating person-to-person spread is occurring [2-8]. This raised concern because when a virus is new, does not know how it affects people. At this time, it's unclear how easily or sustainably this virus is spreading between people [9]. Moreover, according to one study, presumed hospital-related transmission of SARS-CoV-2 was suspected in 41% of patients [8]. Based on the evidence of a rapidly increasing incidence of infections [11] and the possibility of transmission by asymptomatic carriers [12], SARS-CoV-2 can be transmitted effectively among humans and exhibits high potential for a pandemic [5, 10, 13]. It is very important to stay informed during this outbreak. Moreover, this novel virus is new to the scientific world and many features of the virus are still not understandable due to

its new strains [14]. Hence, the worldwide researchers are now very active to explore the new insights of the virus in order to understand its biological character and mode of spreading. This real boost of research interest on the virus has actually started after the emergence of SARS and MARS, and subsequent COVID-19.

## II. MATERIALS AND METHODS

In response to the COVID-19 pandemic, the White House and a coalition of leading research groups have prepared the COVID-19 Open Research Dataset (CORD-19). CORD-19 is a resource of over 200,000 scholarly articles, including over 90,000 with full text, about COVID-19, SARS-CoV-2, and related corona viruses. This freely available dataset is provided to the global research community to apply recent advances in natural language processing and other AI techniques to generate new insights in support of the ongoing fight against this infectious disease. [15].

## III. EXPLORATORY DATA ANALYSIS (EDA)

The dataset covers 163 countries and almost 2 full months from 2020, which is enough data to get some clues about the pandemic. Let's see a few plots of the worldwide tendency in Fig. 1 to extract some insights:

### Observations:

The global curve shows a rich fine structure, but these numbers are strongly affected by the vector zero country, China. Given that COVID-19 started there, during the initial expansion of the virus there was no reliable information about the real infected cases. In fact, the criteria to consider infection cases was modified around 2020-02-11, which strongly perturbed the curve as you can see from Fig. 1.

### A. COVID-19 Behavior

Since China was the initial infected country, the COVID-19 behavior is different from the rest of the world. The medical system was not prepared for the pandemic; in fact no one was aware of the virus until several cases were reported.

Moreover, China government took strong contention measures in a considerable short period of time and, while the virus is widely spread, they have been able to control the increasing of the infections.

**Observations:**

a) *Smoothness:* Both plots are less smooth than theoretical simulations or the curve from the rest of the world cumulative.

b) *Infected criteria:* The moment in which the criteria to consider an infected case was changed is directly spotted.

c) *Irregularities:* There are some irregularities. I should check the literature in depth to look for evidences, but the reasons may be that both the resources spent to monitor the epidemic and the security measures that have been changing over time.

d) *Plateaux:* It looks like the curve has reached a plateau, which would imply that China is on their maximum

of contagion, which strongly perturbed the curve as you can see from Fig. 2.

**B. Italy, Spain, UK and Singapore**

Both Italy and Spain are experiencing the larger increase in COVID-19 positives in Europe. At the same time, UK is a unique case given that it's one of the most important countries in Europe but recently has left the European Union, which has create an effective barrier to human mobility from other countries. The fourth country studied in this section is Singapore, since it's an Asiatic island, is closer to China and its socio-economic conditions is different from the other three countries, which strongly perturbed the curve as you can see from Fig. 3.

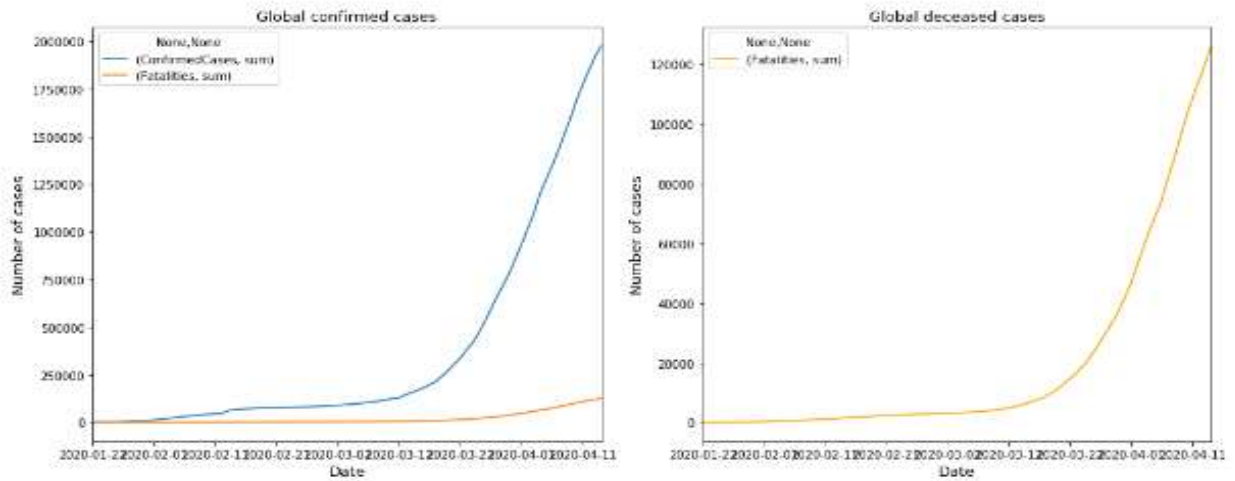


Fig. 1. Global Confirmed Cases Date-Wise.

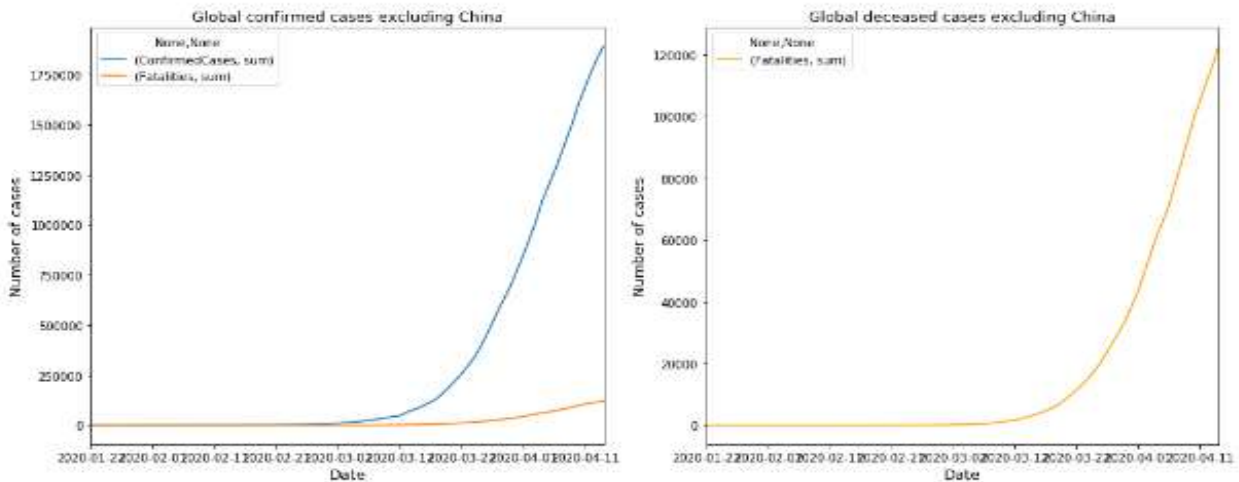


Fig. 2. Global Confirmed Cases Excluding China Date-Wise.

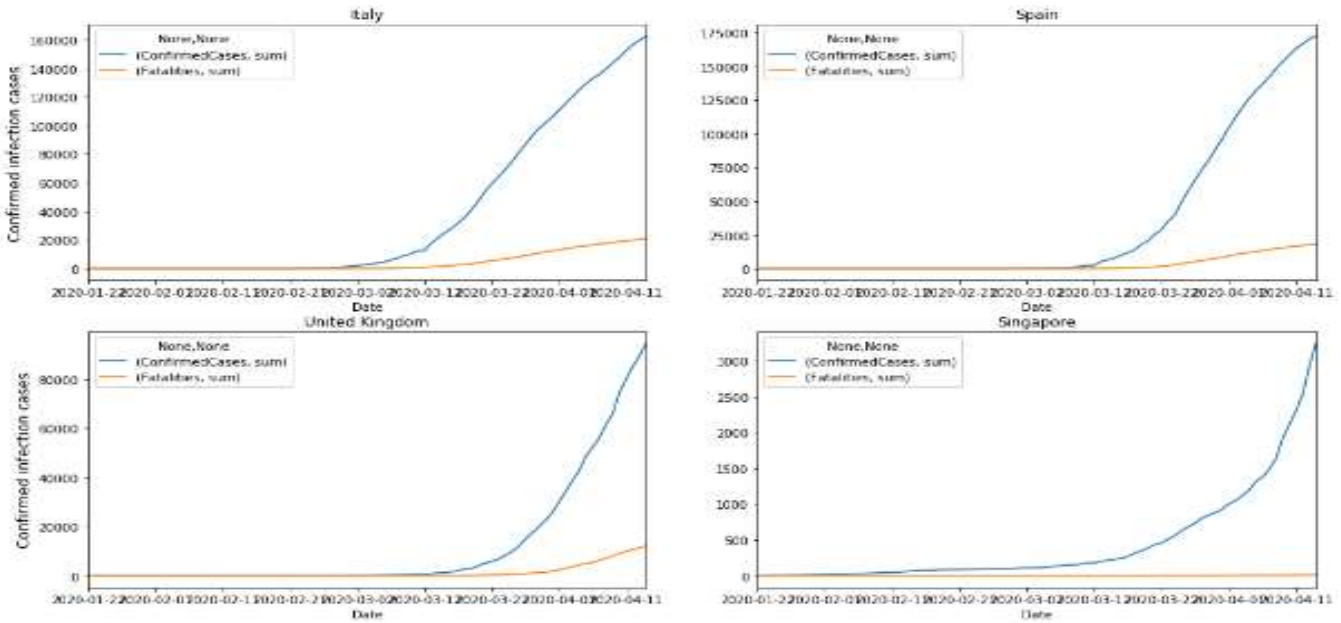


Fig. 3. Confirmed Infection Cases in different Countries.

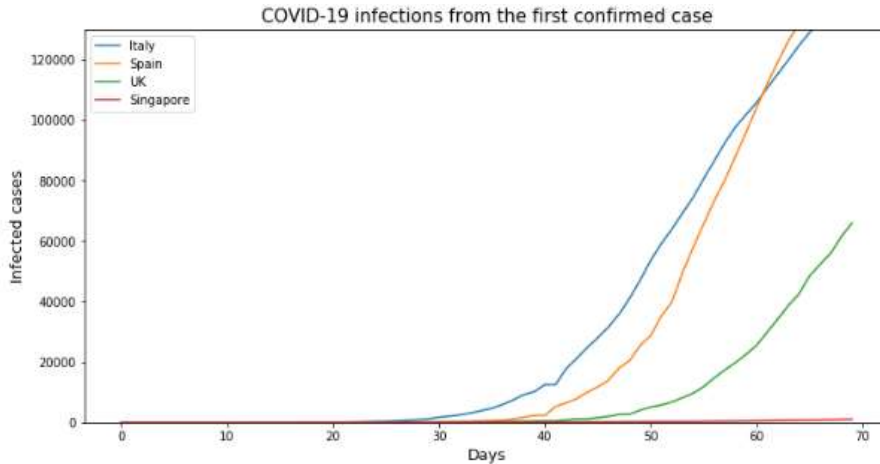


Fig. 4. Infections from the First Confirmed Case Spread Over different Countries.

As a fraction of the total population of each country, in order to compare the four countries, it's also interesting to see the evolution of the infections from the first confirmed case, which is plotted in Fig. 4.

*Observations:*

a) *Italy.* With almost 120.000 confirmed cases, Italy shows one of the most alarming scenarios of COVID-19. The infections curve is very steep, and more than 2% of population has been infected.

b) *Spain.* Spain has the same number of cumulative infected cases than Italy, near 120.000. However, Spain's total population is lower (around 42 millions) and hence the percentage of population that has been infected rises up to 3%.

c) *United Kingdom.* Despite not being very far from them, the UK shows less cases. This may be due to the number of tests performed, but it's soon to know for sure. The

number of cases is around 40.000, this is, a 0.6 of the total population.

d) *Singapore.* Singapore is relatively isolated given that is an island, and the number of international travels is lower than for the other 3 countries. The number of cases is still very low (>1000), despite the general tendency is to increase. However, the infections started faster in the beginning, but the slope of the infections curve hasn't increased very much in the past weeks. A 0.2% of the population was infected.

IV. SIR MODEL

Some general behavior of the virus in aggregated data, for the country where the corona virus was originated and for four other interesting countries. The purpose of this study is to develop a predictive model in order to understand the key factors that impact the COVID-19 transmission, Let's move on to one of the most famous epidemiologic models: SIR, the workflow shown in Fig. 5.

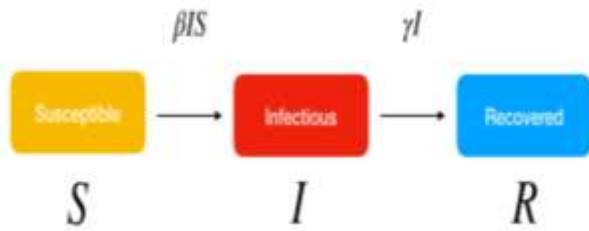


Fig. 5. SIR Workflow.

SIR is a simple model that considers a population that belongs to one of the following states:

- *Susceptible (S)*. The individual hasn't contracted the disease, but she can be infected due to transmission from infected people.
- *Infected (I)*. This person has contracted the disease.
- *Recovered/Deceased (R)*. The disease may lead to one of two destinies: either the person survives, hence developing immunity to the disease, or the person is deceased.

There are many versions of this model, considering birth and death (SIRD with demography), with intermediate states, etc. However, since world is in the early stages of the COVID-19 expansion and interest is focused in the short term, will consider that people develops immunity (in the long term, immunity may be lost and the COVID-19 may come back within a certain seasonality like the common flu) and there is no transition from recovered to the remaining two states.

#### A. Implementing the SIR Model

SIR model can be implemented in many ways: from the differential equations governing the system, within a mean field approximation or running the dynamics in a social network (graph). For the sake of simplicity run a numerical method (Runge-Kutta) to solve the differential equations system.

In order to solve the differential equations system, a 4th order Runge-Kutta method is developed.

And finally, to obtain the evolution of the disease, simply define the initial conditions and call the Runge-Kutta method.

The number of infected cases increases for a certain time period, and then eventually decreases given that individuals recover/decease from the disease. The susceptible fraction of population decreases as the virus is transmitted, to eventually drop to the absorbent state 0, which is predicted in Fig. 6. The opposite happens for the recovered/deceased case. Notice that different initial conditions and parameter values will lead to other scenarios, feel free to play with these numbers to study the system.

#### B. Fit SIR Parameters to Real Data

The SIR model is purely theoretical, and interested in a real approximation of the COVID-19 expansion in order to extract insights and understand the transmission of the virus. Model needs to extract the  $\beta$  and  $\gamma$  parameters for each case to predict the evolution of the system.

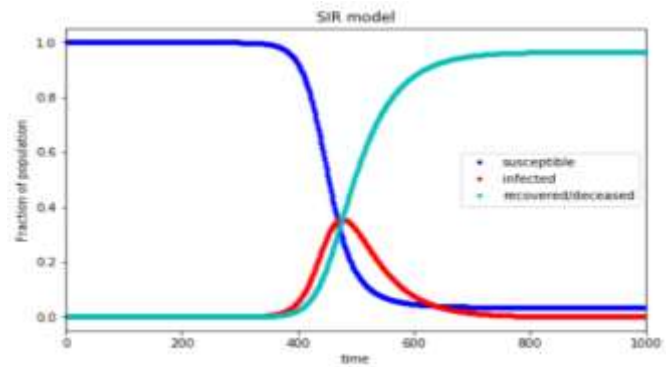


Fig. 6. Susceptible Fraction of Population Decreases as the Virus is Transmitted.

#### C. Data Enrichment

Analyzing SIR simulations was meant to understand a model that approximately resembles the transmission mechanism of many viruses, including the COVID-19. However, there are alternative methods that may prove being equally useful both to predict and to understand the pandemic evolution. Many of these methods rely on having rich data to extract conclusions and allow algorithms to extrapolate patterns in data, and that is exactly what is going to be implemented.

#### D. Main Workflow of this Section

- Join data, filter dates and clean missing.
- Compute lags and trends.
- Add country details.

*Disclaimer:* This data enrichment is not mandatory and could end up without using all of the new features in the model. However, this is consider as a didactical step that will surely add some value, for example in an in-depth exploratory analysis.

1) *Join data, filter dates and clean missing:* First of all, let's perform some pre-processing to prepare the dataset, consisting on:

- *Join data.* Join train/test to facilitate data transformations.
- *Filter dates.* According to the challenge conditions, remove Confirmed Cases and Fatalities post 2020-03-12. Create additional date columns.
- *Missing.* Analyze and fix missing values.

#### Observations:

- a) "Confirmed Cases" and "Fatalities" are now only informed for dates previous to 2020-03-12.
- b) The dataset includes all countries and dates, which is required for the lag/trend step.
- c) Missing values for "Confirmed Cases" and "Fatalities" have been replaced by 0, which may be dangerous if it is not remembered at the end of the process. However,

since training is done only on dates previous to 2020-03-12, this won't impact the prediction algorithm.

d) A new column "Day" has been created, as a day counter starting from the first date.

2) *Compute lags and trends:* Enriching a dataset is a key to obtain good results. In this case, two different transformations are applied:

a) *Lag.* Lags are a way to compute the previous value of a column, so that the lag 1 for Confirmed Cases would inform the column from the previous day.

b) *Trend.* Transforming a column into its trend gives the natural tendency of this column, which is different from the raw value.

The backlog of lags is applied for 14 days, while for trends is for seven days.

3) *Add country details:* Variables like the total population of a country, the average age of citizens or the fraction of people living in cities may strongly impact on the COVID-19 transmission behavior. Hence, it's important to consider these factors. The dataset is based on Web Scrapping for this purpose.

4) *Predictions for the early stages of the transmission:* The objective in this section consists of predicting the evolution of the expansion from a data-centric perspective, like any other regression problem. To do so, remember that the challenge specifies that submissions on the public LB should only contain data previous to 2020-03-26.

a) *Tools utilized:* Previously published automated machine learning tool (<https://automated-machinelearning-gitamcse.shinyapps.io/MLPv3/>) [16, 17] is utilized here for building multiple models on the imputed dataset. The natural advantage of the AMLT tool is to choose multiple train and test sets coupled with a suitable statistical algorithm to build

the best models out of the available data. AMLT tool also does the test validation automatically, which will be helpful to understand the accuracy of each model.

b) *Models to apply:*

- 1) Linear Regression for one country
- 2) Linear Regression for all countries

## V. LINEAR REGRESSION FOR ONE COUNTRY

Since we are interested into predicting the future time evolution of the pandemic, the first approach consists on a simple Linear Regression. However, remind that the evolution is not linear but exponential (only in the beginning of the infection), so that a preliminary log transformation is needed.

Visual comparison of both cases for Spain and with data from last 10 days informed, starting on March 1<sup>st</sup> is depicted in Fig. 7.

As you see, the log transformation results in a fancy straight-like line, which is awesome for Linear Regression. However, let me clarify two important points:

- This "roughly exponential behavior" is only true for the initial infection stages of the pandemic (the initial increasing of infections on the SIR model), but that's exactly the point where most countries are at the moment.
- Why do I only extract the last 10 days of data? For three reasons:

- 1) In order to capture exactly the very short term component of the evolution
- 2) To prevent the effects of certain variables that have been impacting the transmission speed (quarantine vs. free circulation)
- 3) To prevent differences on criteria when confirming cases (remember that weird slope on the China plot?).

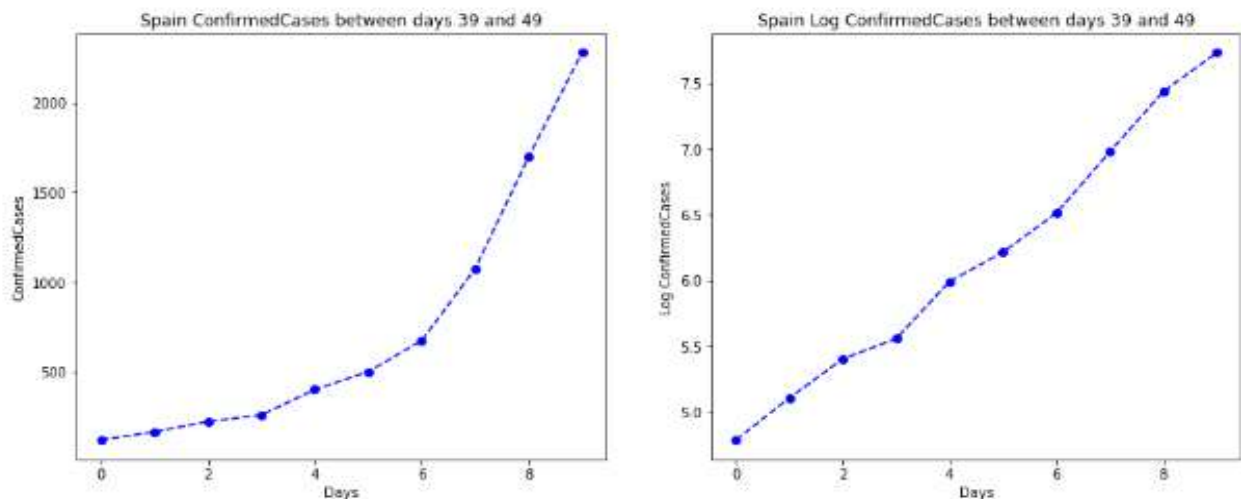


Fig. 7. Spain Confirmed Cases Day-Wise.

This first model is very simple, and only elemental features will be considered: Country/Region, date information, Long and Lat. Lags. Engineered columns like lags, trends and country details are not introduced as an input. Finally, the workflow for the Basic Linear Regression model is:

- 1) *Features*. Select features.
- 2) *Dates*. Filter train data from 2020-03-01 to 2020-03-18.
- 3) *Log transformation*. Apply log transformation to Confirmed Cases and Fatalities.
- 4) *Infinities*. Replace infinities from the logarithm with 0. Given the asymptotic behavior of the logarithm for  $\log(0)$ , this implies that when applying the inverse transformation (exponential) a 1 will be returned instead of a 0. This problem does not impact many countries, but still needs to be tackled sooner or later in order to obtain a clean solution.
- 5) *Train/test split*. Split into train/valid/test.
- 6) *Prediction*. Linear Regression, training country by country and joining data.
- 7) *Submit*. Submit results in the correct format, and applying exponential to reverse log transformation.

#### A. Linear Regression for All Countries

An alternative method to setting the number of days for the training step is to simply keep all data for each country since the first case was confirmed. However, since there are certain countries where the initial outbreak was very smooth (i.e. in Spain there was only one confirmed case for 7 days in a row), predictions may be biased by these initial periods.

Final LMSE score for week 2, with training data prior to 2020-03-19 and measures on date 2020-04-01: 1.19681.

## VI. Conclusion

### A. Results

1) *Parameters*. Two full weeks of training used (from February 26th to March 11th), with their previous 30 lags.

2) *Enough data*. (Spain, Italy, Germany). For countries with several Confirmed Cases  $\neq 0$  in the train dataset (prior to March 11th), predictions are very precise and similar to actual confirmed data.

3) *Poor data*. Countries with a small number of data points in the train dataset show a potentially disastrous prediction. Given the small number of cases, the log transformation followed by a Linear Regression is not able to capture the future behavior.

4) *No data*. When the number of confirmed cases in the train dataset is 0 or negligible, the model predicts always no infections.

### B. Discussion

1) The objective of this work is to provide some insights about the COVID-19 transmission from a data-centric perspective in a didactical and simple way. Predicted results should not be considered in any way an affirmation of what will happen in the future. Observations obtained from data exploration are personal opinions.

2) Models tailored specifically for epidemic spreading (i.e. SIR and its versions) are designed to reproduce a certain phenomenology, in order to understand the underlying mechanics of a contagion process. On the other hand, the simple machine learning approaches I used aim to predict the short term evolution of the infection in the current regime. They might eventually help to find some features or parameters that are particularly important for the model's fitting, but by no means should they be confused with scientific epidemic models.

3) The success of the current predictions is strongly dependent on the current spreading regime, in which the number of infections is still increasing exponentially for many countries. However, they cannot provide a reliable expected day by which the maximum contagion peak will be reached. Epidemic models are closer to obtaining such estimations, but there's a large number of variables that need to be considered for this (quarantines, quality of the medical resources deployed, environmental measures...).

4) In order to achieve such results, a considerable amount of tuning is required. Filter how many previous dates should be used for the fitting step, when to use lags or not, and even missing replacements were very rough due to the log transformation.

### C. Declaration

Predictive models can be used for several purposes, but they never (try to) substitute recommendations from experts.

#### REFERENCES

- [1] Lu H, Stratton CW, Tang YW. Outbreak of pneumonia of unknown etiology in Wuhan China: the mystery and the miracle. *J Med Virol* 2020 Jan 16 [Epub ahead of print]. doi: 10.1002/jmv.25678.
- [2] Li Q, Guan X, Wu P, Wang X, Zhou L, Tong Y, et al. Early transmission dynamics in Wuhan, China, of novel coronavirus-infected pneumonia. *N Engl J Med* 2020 Jan 29 [Epub ahead of print]. doi: 10.1056/NEJMoa2001316.
- [3] Gorbalenya AE, Baker SC, Baric RS, de Groot RJ, Drosten C, Gulyaeva AA, et al. Severe acute respiratory syndrome-related coronavirus: the species and its viruses—a statement of the Coronavirus Study Group. *bioRxiv* 2020 Feb 11. doi: 10.1101/2020.02.07.937862.
- [4] Chen N, Zhou M, Dong X, Qu J, Gong F, Han Y, et al. Epidemiological and clinical characteristics of 99 cases of 2019 novel coronavirus pneumonia in Wuhan, China: a descriptive study. *Lancet* 2020;395:507–13. doi: 10.1016/S0140-6736(20)30211-7.
- [5] Huang C, Wang Y, Li X, Ren L, Zhao J, Hu Y, et al. Clinical features of patients infected with 2019 novel coronavirus in Wuhan, China. *Lancet* 2020;395:497–506. doi: 10.1016/S0140-6736(20)30183-5.
- [6] Wang C, Horby PW, Hayden FG, Gao GF. A novel coronavirus outbreak of global health concern. *Lancet* 2020;395:470–3. doi: 10.1016/S0140-6736(20)30185-9.
- [7] Holshue ML, DeBolt C, Lindquist S, Lofy KH, Wiesman J, Bruce H, et al. First case of 2019 novel coronavirus in the United States. *N Engl J Med* 2020 Jan 31 [Epub ahead of print]. doi: 10.1056/NEJMoa2001191.
- [8] Wang D, Hu B, Hu C, Zhu F, Liu X, Zhang J, et al. Clinical characteristics of 138 hospitalized patients with 2019 novel coronavirus-infected pneumonia in Wuhan, China. *JAMA* 2020 Feb 7 [Epub ahead of print]. doi: 10.1001/jama.2020.1585.
- [9] Chang D, Lin M, Wei L, Xie L, Zhu G, Dela Cruz CS, et al. Epidemiologic and clinical characteristics of novel coronavirus infections involving 13 patients outside Wuhan, China. *JAMA* 2020 Feb 7 [Epub ahead of print]. doi: 10.1001/jama.2020.1623.



- [10] Carlos WG, Dela Cruz CS, Cao B, Pasnick S, Jamil S. Novel Wuhan (2019- nCoV) coronavirus. *Am J Respir Crit Care Med* 2020;201:P7–8. doi: 10.1164/rccm.2014P7.
- [11] Zhao S, Lin Q, Ran J, Musa SS, Yang G, Wang W, et al. Preliminary estimation of the basic reproduction number of novel coronavirus (2019-nCoV) in China, from 2019 to 2020: a data-driven analysis in the early phase of the outbreak. *Int J Infect Dis* 2020;92:214–17. doi: 10.1016/j.ijid.2020.01.050.
- [12] Biscayart C, Angeleri P, Lloveras S, Chaves T, Schlagenhaut P, Rodriguez- Morales AJ. The next big threat to global health? 2019 novel coronavirus (2019-nCoV): What advice can we give to travellers? — Interim recommendations January 2020, from the Latin-American Society for Travel Medicine (SLAMVI). *Travel Med Infect Dis* 2020;101567. doi: 10.1016/j.tmaid.2020.101567.
- [13] Munster VJ, Koopmans M, van Doremalen N, van Riel D, de Wit E. A novel coronavirus emerging in China—key questions for impact assessment. *N Engl J Med* 2020 Jan 24 [Epub ahead of print]. doi: 10.1056/NEJMp20 0 0929.
- [14] Qing, E., Gallagher, T., 2020. SARS coronavirus redux. *Trends Immunol.* 41, 271–273. <https://doi.org/10.1016/j.it.2020.02.007>.
- [15] <https://www.semanticscholar.org/cord19:arXiv:2004.10706>.
- [16] DeepaRani Gopagoni, P V Lakshmi, 2020. Automated machine learning tool, the first stop for data science and statistical model building. *IJACSA* 2020, 11(2),410-418,DOI: 10.14569/IJACSA.2020.0110253. <https://automatedmachinelearning-gitamcse.shinyapps.io/MLPv3/>.
- [17] DeepaRani Gopagoni, P V Lakshmi, An Application of Machine Learning Strategies to Predict Alzheimer’s Illness Progression in Patients *International Journal of Advanced Research in Engineering and Technology (IJARET)* Volume 11, Issue 6, June 2020, pp. 1056-1063, Article ID: IJARET\_11\_06\_095 DOI: 10.34218/IJARET.11.6.2020.95.

# Automated Estrus Detection for Dairy Cattle through Neural Networks and Bounding Box Corner Analysis

Nilo M. Arago<sup>1</sup>, Chris I. Alvarez<sup>2</sup>, Angelita G. Mabale<sup>3</sup>, Charl G. Legista<sup>4</sup>, Nicole E. Repiso<sup>5</sup>  
Rodney Rafael A. Robles<sup>6</sup>, Timothy M. Amado<sup>7</sup>, Romeo Jr. L. Jorda<sup>8</sup>, August C. Thio-ac<sup>9</sup>  
Jessica S. Velasco<sup>10</sup>, Lean Karlo S. Tolentino<sup>11</sup>  
Department of Electronics Engineering, Technological University of the Philippines<sup>1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11</sup>  
University Extension Services Office, Technological University of the Philippines<sup>11</sup>  
Ermita, Manila 1000, Philippines<sup>1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11</sup>

**Abstract**—Thorough and precise estrus detection plays a crucial role in the fertility of dairy cows. Farmers commonly used direct visual monitoring in recognizing estrus signs which demands time and effort and causes misinterpretations. The primary sign of estrus is the standing heat, where the dairy cows stand to be mounted by other cows for a few seconds. Through the years, researchers developed various detection methods, yet most of these methods involve contact and invasive approaches that affect the estrus behaviors of cows. So, the proponents developed a non-invasive and non-contact estrus detection system using image processing to detect standing heat behaviors. Through the TensorFlow Object Detection API, the proponents trained two custom neural network models capable of visualizing bounding boxes of the predicted cow objects on image frames. The proponents also developed an object overlapping algorithm that utilizes the bounding box corners to detect estrus activities. Based on the conducted tests, an estrus event occurs when the centroids of the detected objects measure a distance of less than 360px and have two interior angles with another fixed point of less than 25° and greater than 65° for Y and X axes, respectively. If the conditions are met, the program will save the image frame and will declare an estrus activity. Otherwise, it will restart its estrus detection and counting. The system observed 17 cows, a carabao, and a bull through the cameras installed atop of a cowshed, and detects the estrus events with an efficiency of 50%.

**Keywords**—Dairy cows; estrus detection; image processing; TensorFlow Object Detection API; custom neural network; object overlapping

## I. INTRODUCTION

The estrus cycle of mammals, such as dairy cattle and water buffaloes, is the period from one estrus to the next. On a typical basis, the cycle has an average duration of 21 days. In the Philippines, farmers observe a period of between 18 and 24 days. Research shows that estrus usually lasts between 10 and 18 hours. Even so, recent studies show that modern dairy cows' cycles are about 8 hours shorter [1] [2]. A livestock requires thorough heat detection, and correct timing of artificial insemination. So, not being able to detect in-heat signatures of cattle may lead to low fertility. If the producers could not detect and differ the in-heat and non-heat signs of the cattle, the farm may suffer. Also, the extended calving intervals and semen expenses affect the farm's economic status.

Farmers and researchers have introduced various methods to determine in-heat signatures in livestock. Today, farmers

commonly use visual observation of estrus signs of cows. But doing so may lead to misinterpretations as well. Meanwhile, some farmers track the roaming activities of the cows through a motion sensor on the cows' neck or leg. This method still varies depending on the efficiency and accuracy of the devices [3].

Several companies in America and Europe developed electronic products and services such as the AfiACT, the HeatWatch system, the MountCount, etc. to identify the cows' estrus behaviors [4]. But in Asia, there are few companies known to offer such products. And in the Philippines, companies offering these types of products and services are non-existent. These show how underdeveloped the cattle industry in the Philippines is. According to the Philippine Statistics Authority, the fourth reading of the total cattle production in 2018 is 0.33 percent lower than in 2017. The stock of cattle is also decreased by 0.73 percent, and the rate of slaughter is high and rising [5]. These statistics proved that the Philippines' performance in cattle production is slower than in other ASEAN countries. That is why farmers and researchers should develop new methods to meet the demands of the country.

As a solution to the problem, in this paper, the researchers proposed a non-invasive and non-contact estrus detection system that uses image processing and artificial intelligence through TensorFlow Object detection API to identify standing heat behaviors of Holstein-Friesian and Sahiwal crosses. The research specifically aims to: (1) develop an automated estrus detection system which visualizes bounding boxes of the cattle objects, and verifies if the overlapping instances are estrus activities through the surveillance system; and (2) conduct an evaluation and assessment on the system's functionality and reliability of detection in comparison with the manual visual inspection methods of the farmers.

The findings of the study will benefit small and large farms in the cattle industry, given the current lack of commercially available products and services, and advanced breeding methods. The implementation of the estrus detection system minimizes the workload of farmers through the real-time monitoring capabilities of the system and increases the dairy production and fertility rate of cows through immediate insemination. Such benefits consequently contribute to the economic growth of the farms.

This research paper is structured as follows: Section II pertains to the gaps and limitations of the related researches, Section III defines the materials and methods used by the researchers, Section IV explains the detection and database results of the study, Section V declares the conclusion and Section VI enumerates possible future works of the research.

## II. RELATED WORKS

Researchers develop high-tech devices that helps farmers track the estrus signs of cows. Such technologies based its efficiency on the detection of physical activities, mounting behaviors, body temperature, etc. [6].

In [7], the researchers developed an estrus detection system based on the following behavior of the cows for a short time using IP cameras. The system implements a motion detection technique to identify probable mounting regions, and blob analysis on the said regions to detect changes on the image frames. By incorporating both methods, the proponents were able to accurately identify true estrus events on the surveillance feed.

Talukder *et al.* tested the effectivity of implementing infrared thermography (IRT) in detecting estrus behaviors of dairy cattle. The proponents also incorporated a breeding indicator with IRT which resulted in a sensitive heat detector with false-positive results. The technology can only yield true estrus events only when the IRT was implemented during the ovulation phase of the subjects [8].

In [9], the researchers devised a cattle identifier based on Region Based Convolutional Neural Networks (R-CNN) in an open field setup using unmanned aerial vehicles (UAV) drones. The study has shown great results in detecting unique individual cow patterns through deep learning frameworks and end-to-end training of image datasets. However, false-positive results still occur due to the similarity of structures and features of some cows.

Yang *et al.* also proposed an estrus detection system based on the following and restlessness behaviors of the cows using infrared technology. The infrared cameras were able to monitor and detect estrus events at both daytime and nighttime with the aid of artificial lighting. Despite that, their experimentations showed that the efficiency for detecting objects was greater in contrast to the visual observation considering good illumination in the area [10].

Meanwhile, Xia *et al.* constructed an estrus detection system based on the activities of the cows using pedometers and readers. Through the pedometers and the readers, the system was able to gather and analyze cow information to declare estrus and notify the end-users via text messages. The results proved the system's accuracy, in which it can replace the conventional rectum identification of cows in detecting estrus [11].

In [12], the researchers also proposed an estrus detection system through geometric region analysis using fixed IP cameras. This system's operability is similar to the aforementioned studies that filter the collected image frames and extracts the relevant features of the cows from the images to perform analysis and identification of estrus. Still, the

proposed techniques in this research accurately recognized the mounting behaviors of the cows with minimal false-positive detection rates.

Table I shows the comparison framework of the related works in this research. Unlike with the aforementioned studies, this research performs estrus detection by detecting Holstein-Friesian and Sahiwal Crosses, a bull, and a water buffalo from the surveillance feed of three pan-tilt-zoom (PTZ) cameras (DH-SD22404T-GN Lite Series, 4 MP). The researchers also customized two neural network models using pre-trained frameworks from the TensorFlow Zoo for the object detection and utilized bounding box corners for the analysis of overlapping instances in the image sequences and declaration of estrus events.

TABLE I. BRIEF COMPARISON FRAMEWORK OF THE RELATED WORKS

| Authors                                     | Breed of Cattle to be monitored | Materials and Methods   |   |
|---|---------------------------------|-------------------------|---|
|   |                                 | Sensors used            | Techniques and Algorithms used  |
| Tsai and Huang (2014) [7]                   | Holstein                        | IP Dome Camera          | Motion Detection, Region Segmentation, Foreground Segmentation, and Blob Analysis   |
| Talukder <i>et al.</i> (2014) [8]           | Holstein-Friesian               | Thermal Infrared Camera | Infrared Thermography   |
| Andrew, Greatwood, and Burghardt (2017) [9] | Holstein-Friesian               | UAV integrated camera   | R-CNN Localization, and Tracking  |
| Yang, Lin, and Peng (2017) [10]             | Holstein                        | Infrared Camera         | Motion Detection, Region Segmentation, Foreground Segmentation.   |
| Xia <i>et al.</i> (2017) [12]               | Holstein                        | Pedometer and reader    | Motion Analysis   |
| Guo, Zhang, He, Niu, and Tan (2019) [13]    | Holstein                        | Fixed IP Camera         | Background Subtraction with Color and Texture Features (BSCTF), Geometric and Optical Flow feature extraction, Support Vector Machine (SVM) |

## III. METHODOLOGY

### A. Research Locale - Barn

In this research, the estrus detection system is deployed in a small-scale commercial farm in the province of San Ildefonso, Bulacan, in the Philippines. The barn houses 17 Holstein-Friesian and Sahiwal crosses, a bull, and a water buffalo. Similarly with the research of Porto *et al.* [13], they have observed some delimiting factors in the barn that may affect the automated detection system, such as: high variation in lamination in areas near the open side of the barn; metal surfaces of stable crossbars; color indifferences of cows; and surface reflection caused by manure or dirt. The panoramic top-viewed images of the barn are crucial in to capture image frames which shows the true shape of cow's body [13]. To

capture the panoramic top-view images, three 4 Megapixel Pan-tilt-zoom (PTZ) Network Cameras, as in [13], were installed at a height of 3.78m. Each camera monitors an area for about 4.87 m x 3.97 m with a separation distance of approximately 2.98 m apart atop the cowshed, as shown in Fig. 1 and Fig. 2.

### B. TensorFlow Object Detection API

TensorFlow Object Detection API is a framework that is currently being utilized today to resolve object detection problems. With this, deploying accurate machine learning models that can localize and identify multiple objects in an image frame is easier, as in. Within the models, the feature extraction and the classification processes play vital roles in the cow pattern recognition, as in [14].

According to Huang *et al.*, there will be trade-offs between speed and accuracy in constructing an object detection architecture that depends on the application and platform [15]. In their repository, the user can modify the model to satisfy his/her requirements and platform. The TensorFlow Object Detection API library comprises of object detection structures, such as Single Shot Detector (SSD), Faster Region-based Convolutional Neural Network (Faster R-CNN), etc.

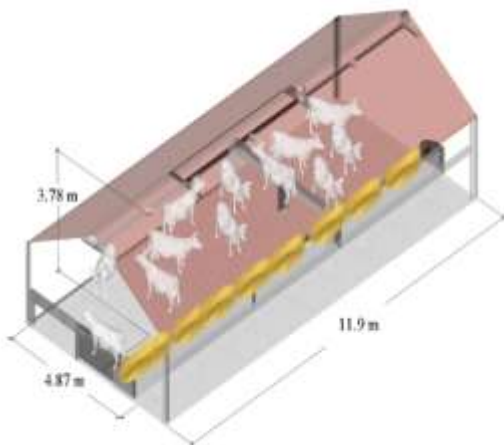


Fig. 1. The Isometric view of the Experimental Setup which Displays the Position of the Cameras.

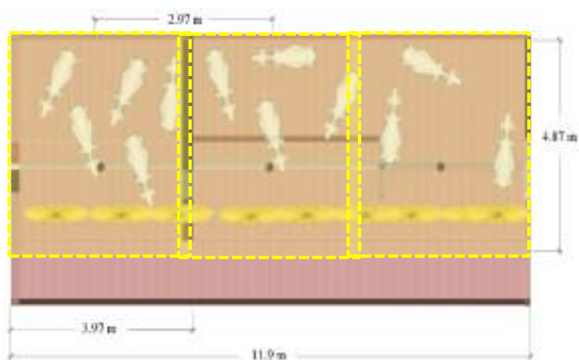


Fig. 2. The Bottom view of the Experimental Setup which Depicts the Camera Separation and the Field of views.

Feature extractors such as Inception, MobileNet [16] and Resnet play critical roles in the speed and accuracy trade-off of

the framework. Even with the recent studies of various researchers, constructing convolutional networks from scratch requires a great volume of image datasets and a long period of training and testing time. That is why transfer learning is more applicable with pre-trained models like the TensorFlow API [17]. Transfer learning is a technique in which a model is reprocessed as a starting point for a second function model [18] [19].

In this research, two (custom) object detection frameworks using TensorFlow CPU and the pre-trained Faster R-CNN [20] and SSD [21] models were developed and integrated as its core architectures from the TensorFlow Zoo.

### C. Data Acquisition and Pre-Processing

In this research, all of the cows, including the bull and the water buffalo, are pre-identified with a corresponding ID. In building the dataset, a total of 1400 images for each defining class for the Faster R-CNN model, and a total of 21,912 images of cows for the SSD model were used. By accessing the playback videos from the Network Video Recorder, and using image processing techniques through OpenCV, the image frames were obtained at a rate of 1 frame per second.

To provide the necessary supervised learning for the detection system, the researchers used a label annotator tool, as in [19]. For the Faster R-CNN model, each cow object on every image frame were annotated as: "BULL"; "CARACOW"; "COW A"; "COW B"; "COW C"; "COW D"; "COW E"; "COW F"; "COW G"; "COW H"; "COW I"; "COW J"; "COW K"; "COW L"; "COW M"; "COW N"; "COW O"; "COW P"; and "COW Q" in accordance to its COW ID whereas, for the SSD model, all objects were labeled as "COW". The annotations will be saved as Extensible Markup Language data files (XML) and will be processed after the data slicing. Next, the image datasets were divided into the training and the testing data. The partition used for data slicing is 90:10 wherein 90% is for the training data while the 10% is for the testing data, as in [9] [17] [22] [23].

Afterwards, two label maps for each model were created, in which 19 labels were listed for the Faster R-CNN model but only 1 label for the SSD model. From the XML data files, TensorFlow Records in "RECORD" format will be generated. These records contain the filename, the labels (classes), the height and width of the images, and the bounding box corners (xmin, ymin, xmax, and ymax), as in [9] [24].

### D. Configuring the Pipeline

In selecting a pre-trained model, the performance, speed, and mean Average Precision (mAP) that define the accuracy of the detector were considered, as in [16] [18]. According to the analysis of Huang *et al.* [15], the Faster R-CNN model with Inception V2 and SSD model with Inception V2 yields a mAP of 28 and 24, respectively, which requires a speed of at least 58 ms and 42 ms per image, respectively. To configure the pipeline, the researchers utilized two of the pre-trained models provided by TensorFlow Zoo. The speed and mAP of the given pre-trained models were considered, and the Faster R-CNN and the SSD with Inception V2 models will be implemented.

The pipeline configurations given in Fig. 3 and Fig. 4 only show the changes made from the pre-configured models.

Adjusting some of the parameters does not necessarily give similar results on other applications.

#### E. Training the Networks

In training the custom neural network models, it is expected to obtain a minimum TotalLoss value of 1.0 or less. The training job for both the Faster R-CNN and SSD with Inception V2 models can be monitored using the TensorBoard. Once the optimal range of TotalLoss is observed, the training job can be interrupted. Also, checkpoints that represent the training steps are being saved in the system unit as the training progresses. These checkpoints will be used in visualizing the training performance. The training for both networks took approximately 387 hours.

Fig. 5 depicts the TotalLoss graph obtained from training the Faster R-CNN with Inception V2 model while Table II shows the model's training metrics having TotalLoss between approximately 0.04 and 0.14.

Fig. 6 depicts the TotalLoss graph obtained from training the SSD with Inception V2 model while Table III shows the model's training metrics having TotalLoss between approximately 1.7 and 2.0.

Once the training jobs are complete, trained inference graphs will be generated to be integrated into the object detection program.

#### F. Estrus Detection Criteria

According to the research done by Tsai *et al.* an estrus event in images projects an object with a size of about 2-cows which will change into roughly 1.5-cows during the activity. Furthermore, based on the blob analysis and segmentation approach, if the distance between two centroids of the cows exhibiting "following" behavior is equal to or less than the distance threshold for more than 2 seconds or exactly equal to 4 seconds, the system will declare an estrus activity [7]. By adapting this research with the abovementioned study, the researchers were able to construct a similar detection rule for identifying the standing-heat activities of cows. The researchers initially hypothesized that in a panoramic top-viewed image depicting a standing-heat activity, the mounting (top) cattle's head and half body overlaps the other (bottom) cattle's half body. Consequently, having both objects stand very close to each other, an estrus activity can be declared.

In the numerical and photographic perspective, if the cattle's head and half of its body is treated as 0.5-cow while it mounts the other cattle's body (1.0-cow) on the prescribed time, the total length will eventually be equivalent to roughly 1.5-cows, giving the idea that the cow's features in pixels will be in the same range of value with the latter. Also, if the distance and the angles between their centroids meets a certain threshold, an estrus activity can be declared while taking all into account that the objects are highlighted by bounding boxes through the TensorFlow Object Detection API.

| Faster R-CNN Configuration  |
|---|
| <b>Input:</b> training and testing TFRecord files, and a label map file of the subjects |
| <b>Output:</b> configuration file for training  |
| 1: num_classes = {19};  |
| 2: feature_extractor:   |
| 3: type = {faster_rcnn_inception_v2};   |
| 4: second_stage_post_processing:  |
| 5: batch_non_max_suppression:   |
| 6: scales = {0.25, 0.5, 1.0, 2.0}   |
| 7: aspect_ratios = {0.5, 1.0, 2.0};   |
| 8: first_stage_nms_iou_threshold = {0.7};   |
| 9: second_stage_post_processing:  |
| 10: batch_non_max_suppression:  |
| 11: score_threshold = {0.0}   |
| 12: iou_threshold = {0.75}  |
| 13: max_detections_per_class = {1}  |
| 14: max_total_detections = {300};   |
| 15: train_config:   |
| 16: batch_size = {1}  |
| 17: learning_rate = {0.0002, 0.00002, 0.000002};  |
| 18: num_steps = {200,000};  |
| 19: data_augmentation_options = {autoaugment_image};                                    |
| 20: eval_config:  |
| 21: num_examples = {26,600};  |

Fig. 3. Pipeline Configuration for the Faster R-CNN Model.

| SSD Configuration  |
|--|
| <b>Input:</b> training and testing TFRecord files, and a label map file of the subjects                            |
| <b>Output:</b> configuration file for training   |
| 1: num_classes = {1};  |
| 2: feature_extractor:  |
| 3: type = {ssd_inception_v2};  |
| 4: anchor_generator:   |
| 5: ssd_anchor_generator:   |
| 6: aspect_ratios = {0.333, 0.5, 1.0, 2.0};   |
| 8: loss:   |
| 9: hard_example_miner:   |
| 10: iou_threshold = {0.99};  |
| 13: post_processing:   |
| 14: batch_non_max_suppression:   |
| 15: score_threshold = {0.0}  |
| 16: iou_threshold = {0.75}   |
| 17: max_detection_per_class = {19}   |
| 18: max_total_detections = {19};   |
| 19: train_config:  |
| 20: batch_size = {4};  |
| 21: learning_rate = {0.0002, 0.00002, 0.000002};   |
| 22: num_steps = {200,000};   |
| 23: data_augmentation_options = {random_rotate_90, random_horizontal_flip, random_vertical_flip, ssd_random_crop}; |
| 24: eval_config:   |
| 25: num_examples = {21 912};   |

Fig. 4. Pipeline Configuration for the SSD Model.

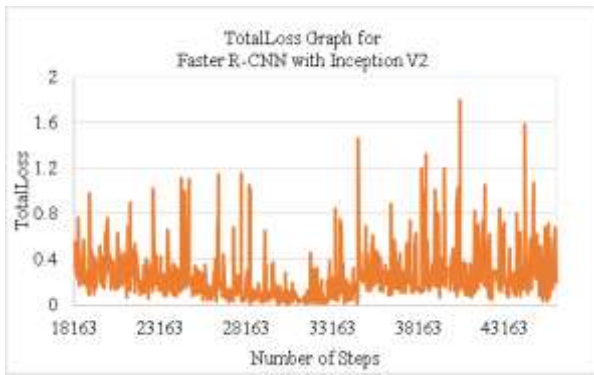


Fig. 5. The Line Graph Representation of the TotalLoss for the Faster R-CNN with Inception V2 Model.

TABLE II. TRAINING METRICS OF THE FASTER R-CNN MODEL

| Steps | Value       |
|-------|-------------|
| 45616 | 0.040272284 |
| 45672 | 0.431099415 |
| 45699 | 0.083715603 |
| 45727 | 0.708893716 |
| 45755 | 0.06934201  |
| 45783 | 0.584971786 |
| 45810 | 0.254837424 |
| 45866 | 0.302880734 |
| 45894 | 0.229811206 |
| 45977 | 0.143239096 |

The formula for the Euclidean distance, as in [25], (1) and the interior angles between centroid (2 and 3) are as follows:

$$D = \sqrt{[(x_2 - x_1)^2 + (y_2 - y_1)^2]} \quad (1)$$

$$\theta_y = \sin^{-1}[(y_2 - y_1)/D] \cdot (180^\circ/\pi) \quad (2)$$

$$\theta_x = \sin^{-1}[(x_2 - x_1)/D] \cdot (180^\circ/\pi) \quad (3)$$



Fig. 6. The Line Graph Representation of the TotalLoss for the SSD Model.

TABLE III. TRAINING METRICS OF THE SSD MODEL

| Steps  | Value       |
|--------|-------------|
| 140129 | 1.688523769 |
| 140171 | 1.409463882 |
| 140213 | 2.607795238 |
| 140467 | 1.906678438 |
| 140764 | 1.774537325 |
| 140978 | 3.171329737 |
| 141275 | 3.082624435 |
| 141445 | 1.949746    |

where  $D$  is the Euclidean distance between two centroids in pixels,  $x_1$  is the centroid of the first object in x-axis,  $x_2$  is the centroid of the second object in x-axis,  $y_1$  is the centroid of the first object in y-axis,  $y_2$  is the centroid of the Second object in y-axis,  $\theta_y$  is the interior angle between centroids in y-axis, and  $\theta_x$  is the interior angle between centroids in x-axis.

### G. Overall Structure of the System

In the input section, the program will load the necessary packages, the label map, and the frozen inference graph that is generated and trained. Consequently, the camera will process the image frames through the VideoCapture objects of the program. In the image processing section, the SSD-based neural network will visualize “COW” predictions and identify object overlapping activities through bounding box corner analysis in real-time, as in [24], if the prediction score exceeds seventy percent. The program will also generate data frames [23] to contain information such as the Cow Name, ID, box coordinates and angles, and date and time of detection, considering there is only one class to be predicted in the image. If the data frames contain more than one detection, the program will filter out the prediction and will calculate the distances between two centroids of object instances and the interior angles between the two centroids and a point connecting it. After meeting the criteria, the program will iteratively count for the overlapping of object instances from 2 to 8 frames per second. If an overlapping of object instances occurred, as in [9] [24], then a copy of the frame will be directed to the Faster RCNN model, which will be initialized to perform image classification and object detection. The model will also be generating data frames to contain the Cow Names, IDs, box coordinates and angles, and date and time of detection of the nineteen classes predicted in the image. If the similar conditions are met in the Faster R-CNN model, an object overlapping or estrus activity will be declared, and the current image frame and record will be locally saved. Subsequently, the program will restart its counter and will continue to perform object detection. A flowchart represents the program flow of the automated estrus detection system using TensorFlow object detection API is illustrated in Fig. 7.

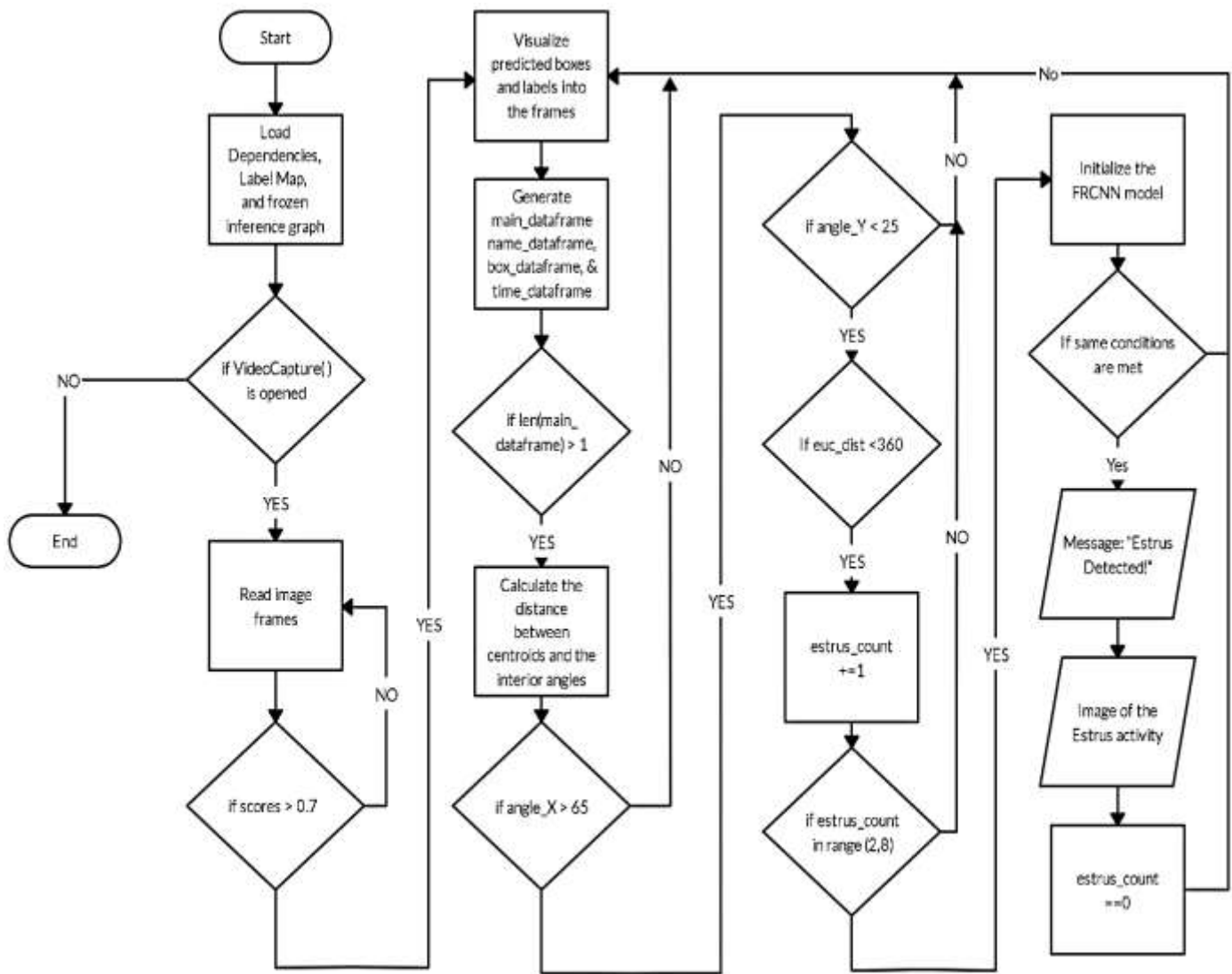


Fig. 7. The Line Graph Representation of the TotalLoss for the SSD Model.

#### IV. RESULTS AND DISCUSSION

##### A. Object Detection Results

The researchers deployed the system and operated locally in the barn for 4 months, with 10 hours of daylight and artificial light exposure in the barn. The system unit can execute the program at 30 fps and 1fps for image frame recognition with the SSD and the Faster R-CNN models, respectively.

Based on the results obtained, the system reported only two confirmed estrus events for 19 subjects in the trials, as shown in Fig. 8 and Fig. 9. Even after attaining acceptable and low TotalLoss values from the training of the Faster R-CNN model, the system still produced inaccurate cow predictions with 50% detection efficiency. According to the cow caretaker, the estrus activity depicted in Fig. 8 between “COW H” and “CARACOW” is validated. But in Fig. 9, the event is misidentified since it should be in-between the “BULL” and “COW Q”, but not in-between “COW P” and “COW N”, respectively.

Moreover, the confidence scores of the model for “COW N” and “COW P” are 71% and 75%, whereas, the confidence scores for “COW H” and “CARACOW” are 96% and 97%, respectively. Nevertheless, the SSD model effectively visualized “COW” objects with confidence scores of 94%, as shown in Fig. 10. These results suggest additional training time, dataset acquisition, and data cleaning to attain higher prediction scores for both models.

##### B. Database Results

Table III represents the validity of the results in monitoring the standing-heat of cattle. Based on the verification of the cow caretaker from the locally saved dataframes and images, the detected event in-between “CARACOW” and “COW H” is “TRUE” while the detected event in-between “COW N” and “COW P” is “FALSE” due to the misidentification of the Faster R-CNN model which led to the 50% detection efficiency.

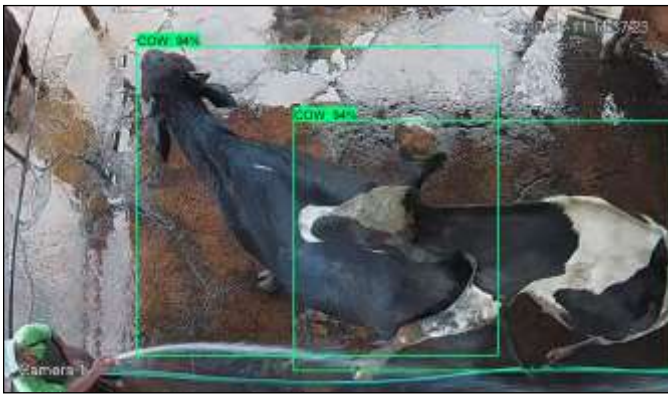


Fig. 8. An Image Frame of an Estrus Activity between the Cow (“COW H”) and the Water Buffalo (“CARACOW”).



Fig. 9. An Image Frame of an Estrus Activity between the Cow (“COW P”) and the Bull (“COW N”).

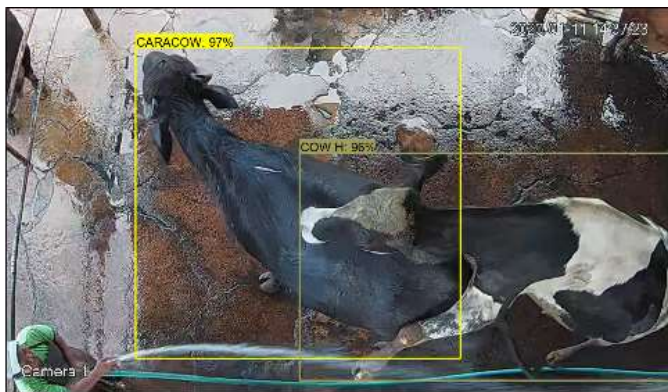


Fig. 10. An Image Frame of an Estrus activity between the Cow (“COW H”) and the Water Buffalo (“COW”).

As shown in Fig. 11, there are a total of 4 app-detections of standing-heat, 4 manually detected standing-heat signs, and 2 “True Positive” and “False Positive” detections from the program. As represented in Table IV, the end-user stated “FALSE” due to the incorrect detection of the system with “COW N” and “COW P” as in-heat cows, which instead should be the “BULL” and “COW Q”. Still, the system initially and correctly detected 4 standing-heat signs, but with 2 false predictions and identifications leading to 2 “True Positive” and 2 “False Positive” results, attaining a 50% detection efficiency.

### C. Performance Assessment with other Related Works

Table V represents the summarized comparison framework between the proposed method and other relevant researchers in estrus detection. As abovementioned, this research deals with the detection of mounting behaviors of Holstein-Friesian and Sahiwal crosses, a bull, and a water buffalo. In contrast with the papers [7], [10]- [12], the proponents integrated a cattle identifier using customized neural network frameworks with a detection efficiency of approximately 90% and 50% for the Faster R-CNN and SSD models, respectively. Besides, most of the formulated methods do not include cattle identifiers since the researchers and the cow caretakers employ manual inspection of the cow tags after the process of standing-heat detection, by which, in this case, the system automatically identifies the cows and declares the estrus event at the same time.

The system also calculated a detection efficiency of 50% as a subsequent effect from the system’s image classifier or cattle identifier. These results suggest the integration of other machine learning algorithms such as Foreground segmentation, background subtraction, support vector machine, and more within the deep learning framework, or the application of unsupervised learning in the detection system.

TABLE IV. TABULAR REPRESENTATION FOR THE VERIFICATION OF THE ACTUAL PROJECT TESTING

| Cow ID | Date of Detection | Time of Detection | Estrus validity | Inseminated |
|--------|-------------------|-------------------|-----------------|-------------|
| 1      | January 11, 2020  | 2:37:22 PM        | TRUE            | NO          |
| 1229-1 | January. 11, 2020 | 2:37:22 PM        | TRUE            | YES         |
| 257-2b | April 2, 2020     | 06:24:13 AM       | FALSE           | NO          |
| 67     | April 2, 2020     | 06:24:13 AM       | FALSE           | NO          |



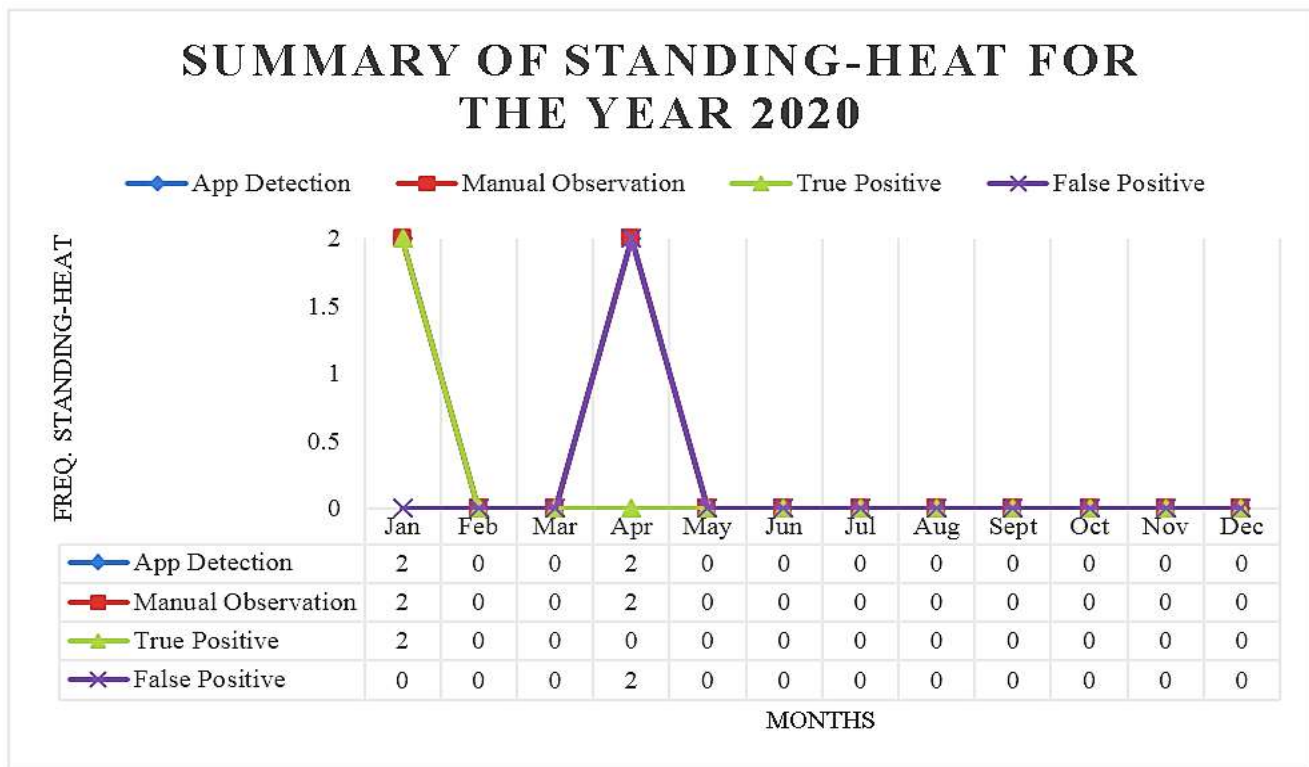


Fig. 11. Summarized Graphical Representation of Database Result for the Frequency of Standing-heat from January to December.

TABLE V. COMPARISON FRAMEWORK OF THE PROPOSED METHOD WITH OTHER RELEVANT WORKS IN TERMS OF PERFORMANCE AND ACCURACY

| Author                                | Tsai and Huang [7]  | Andrew, Greatwood, and Burghardt [9]                 | Yang, Lin, and Peng [10]  | Xia <i>et al.</i> [11]                        | Guo, Zhang, He, Niu, and Tan [12]                         | Proposed work   |
|---------------------------------------|---|--|---|---|---|---|
| Cattle Breed                          | Holstein  | Holstein-Friesian                                    | Holstein  | Holstein                                      | Holstein  | Holstein-Friesian, Sahiwal  |
| Sensors used                          | IP Dome Camera  | UAV integrated camera                                | Infrared Camera   | Pedometer and reader                          | Fixed IP Camera   | PTZ Camera  |
| Algorithm                             | Motion Detection, Region Segmentation, Foreground Segmentation, and Blob Analysis | R-CNN Localization, and Tracking                     | Motion Detection, Region Segmentation, Foreground Segmentation. | Motion Analysis                               | BSCTF, SVM, Geometric and Optical Flow Feature extraction | Faster R-CNN and SSD Localization and Tracking  |
| Output                                | Image   | Image  | Image   | Steps   | Image   | Image, Printed message  |
| Accuracy in object (cattle) detection | -   | 98.13%   | -   | -   | 98.3%   | 94% (SSD), 50% (FRCNN)  |
| Accuracy in estrus detection          | 100% (TP)<br>0.333% (FP)  | -  | -   | 91.86%  | 90.9% (TP)<br>4.2% (FP)                                   | 50% (TP)<br>50% (FP)  |
| Limitations                           | Only at daytime, Indoor setup   | Only at daytime, Not suitable for marker-less cattle | Shadow appearances  | Indoor setup<br>Lack of cattle identification | Indoor setup<br>Lack of cattle identification             | Indoor setup, Not suitable for marker-less cattle   |
| Recommendations                       | Lameness detection  | Larger herd identification                           | Changing type of camera, Cattle identification                  | None  | None  | Integrate other machine learning algorithms, implement unsupervised learning in the framework |

## V. CONCLUSION

In this study, the researchers presented a novel way of detecting estrus for dairy cattle, specifically the Holstein-Friesian and Sahiwal crosses, using the TensorFlow Object Detection API and its pre-trained models such as the Faster R-CNN and Single Shot Detector models with the Inception V2 as the feature extractor. Based from the obtained results, it can be concluded that (1) the Single Shot Detector (SSD) with Inception V2 proved to be effective in visualizing bounding boxes on the single class objects (e.g. "COW") with confidence scores of more than 90%, and (2) the Faster R-CNN with Inception V2 proved to be inaccurate in identifying objects with color indifferences between the subjects and the surface area of the barn obtaining a detection efficiency of 50%. Despite the inaccuracy, the proposed system can detect mounting behaviors of dairy cattle, given that the system will classify only one class (e.g. "COW") as shown in Fig. 8.

## VI. FUTURE WORK

This research aims to report the preliminary attempt and provide learnings for other researchers to devise a system which classifies the dairy cattle subjects as well. The researchers also recommend to: (1) implement unsupervised learning techniques and machine learning algorithms within the deep learning framework that enhances the efficiency of cow classification and estrus detection without the aid of cowhide patterns, (2) develop a system that can monitor and detect mounting behaviors of cows on an outdoor setup, (3) define other suitable estrus detection criteria that maximize the camera's performance and line-of-sight, and (4) integrate a notification subsystem to immediately inform the end-users of the estrus events and initiate insemination on the cows.

## ACKNOWLEDGMENT

The authors would like to express their deepest gratitude to Mr. Arcadio Francisco De Belen Jr., Farm Foreman, and the personnel of De Belen Dairy Farm for allowing the researchers to use the facilities, and to implement and deploy the project on the barn; and to the University Research and Development Services (URDS) of the Technological University of the Philippines, Ermita, Manila for the project funding.

## REFERENCES

- [1] T. J. Parkinson, "Infertility in the Cow Due to Functional and Management Deficiencies," *Veterinary Reproduction and Obstetrics*, vol. 10, pp. 361-407, 2019.
- [2] M. Mária, P. Strapák, I. Szenczióvá, E. Strapáková and O. Hanušovský, "Several Methods of Estrus Detection in Cattle Dams: A Review," *Acta Universitatis Agriculturae et Silviculturae Mendelianae Brunensis*, vol. 66, no. 2, p. 619 – 625, 2018.
- [3] J. B. Roelofs, C. Krijnen and E. v. E.-v. d. Kooij, "The effect of housing condition on the performance of two types of activity meters to detect estrus in dairy cows," *Theriogenology*, vol. 93, pp. 12-15, 2017.
- [4] R. W. Rorie, T. R. Bilby and T. D. Lester, "Application of electronic estrus detection technologies to reproductive management of cattle," *Theriogenology*, vol. 57, pp. 137-148, 2002.
- [5] L. G. S. Bersales, "Philippine Statistics Authority," January-December 2018. [Online]. Available: <https://psa.gov.ph/livestock-poultry-ips/cattle/inventory>. [Accessed 2019].
- [6] S. Reith and S. Hoy, "Review: Behavioral signs of estrus and the potential of fully automated systems for detection of estrus in dairy cattle," *The Animal Consortium* 2017, vol. 12, no. 2, pp. 398-407, 2018.
- [7] D.-M. Tsai and C.-Y. Huang, "A Motion and Image Analysis Method for Automatic Detection of Estrus and Mating Behavior in Cattle," *Computers and Electronics in Agriculture*, vol. 104, pp. 25-31, 2014.
- [8] S. Talukder, K. L. Kerrisk, L. Ingenhoff, P. C. Thomson, S. C. Garcia and P. Celi, "Infrared technology for estrus detection and as a predictor of time of ovulation in dairy cows in a pasture-based system," *Theriogenology*, vol. 81, no. 7, pp. 925-935, 2014.
- [9] W. Andrew, C. Greatwood and T. Burghardt, "Visual Localisation and Individual Identification of Holstein Friesian Cattle via Deep Learning," in 2017 IEEE International Conference on Computer Vision Workshops (ICCVW), Venice, Italy, 2017.
- [10] C.-J. Yang, Y.-H. Lin and S.-Y. Peng, "Develop a Video Monitoring System for Dairy Estrus Detection at Night," in 2017 International Conference on Applied System Innovation (ICASI), Sapporo, Japan, 2017.
- [11] T. Xia, C. Song, J. Li, C. Li, G. Xu, F. Xu, J. Liu, G. M. P. O'Hare and Q. Zhou, "Research and Application of Cow Estrus Detection Based on the Internet of Things," in 2017 IEEE International Conference on Computational Science and Engineering (CSE) and IEEE International Conference, Guangzhou, China, 2017.
- [12] Y. Guo, Z. Zhang, D. He, J. Niu and Y. Tan, "Detection of cow mounting behavior using region geometry and optical flow characteristics," *Computers and Electronics in Agriculture*, vol. 163, 2019.
- [13] Porto, S. MC, C. Arcidiacono, U. Anguzza and G. Cascone, "A computer vision-based system for the automatic detection of lying behaviour of dairy cows in free-stall barns," *Biosystems Engineering*, vol. 115, pp. 184-194, 2013.
- [14] J. O'Rourke and G. T. Toussaint, "Pattern Recognition," in *Handbook of Discrete and Computational Geometry*, Florida, CRC Press LLC, 2017, pp. 1421-1450.
- [15] J. Huang, V. Rathod, C. Sun, M. Zhu, A. Korattikara, A. Fathi, I. Fischer, Z. Wojna, Y. Song, S. Guadarrama and K. Murphy, "Speed/Accuracy Trade-Offs for Modern Convolutional Object Detectors," in 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Honolulu, HI, USA, 2017.
- [16] J. S. Velasco, C. G. Pascion, J. W. Alberio, J. Apuang, J. S. Cruz, M. A. Gomez, B. J. Molina, L. Tuala, A. Thio-ac and R. L. J. Jr., "A Smartphone-Based Skin Disease Classification Using MobileNet CNN," *International Journal of Advanced Trends in Computer Science and Engineering*, vol. 8, pp. 2632-2637, 2019.
- [17] F. Al-Azooa, A. M. Taqia and M. Milanova, "Human Related-Health Actions Detection using Android Camera based on TensorFlow Object Detection API," (IJACSA) International Journal of Advanced Computer Science and Applications, vol. 9, no. 10, pp. 9-23, 2018.
- [18] J. S. Velasco, N. M. Arago, R. M. Mamba, M. V. C. Padilla, J. P. M. Ramos and G. C. Virrey, "Cattle Sperm Classification Using Transfer Learning Models," *International Journal of Emerging Trends in Engineering Research*, vol. 8, pp. 4325-4331, 2020.
- [19] X. Ma, S. Pan, Y. Li, C. Feng and A. Wang, "Intelligent welding robot system based on deep learning," in 2019 Chinese Automation Congress (CAC), Hangzhou, China, 2019.
- [20] S. Ren, K. He, R. Girshick and J. Sun, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 6, pp. 1137-1149, 2016.
- [21] D. A. D. E. Wei Liu, C. Szegedy, S. E. Reed, C.-Y. Fu and A. C. Berg, "SSD: Single Shot MultiBox Detector," *European Conference on Computer Vision* 2016, vol. 9905, pp. 21-37, 2016.
- [22] V. N. and S. A., "Pre-processing," in *Studies in Systems, Decision and Control*, New York City, Springer, 2020, pp. 89-120.
- [23] Porcu, Python for Data Mining Quick Syntax Reference, New York City: Appress, 2018.
- [24] Y. He, C. Zhu, J. Wang, M. Savvides and X. Zhang, "Bounding box regression with uncertainty for accurate object detection," 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 2883-2892, 2019.
- [25] O. B. Elementary Differential Geometry, Amsterdam, Netherlands: Elsevier Inc., 2006.

# An Efficient Cluster-Based Approach to Thwart Wormhole Attack in Adhoc Networks

Kollu Spurthi<sup>1</sup>

Research Scholar, Department. of Computer Science and Engineering, KLEF, AP, India

Dr.T N Shankar<sup>2</sup>

Associate Professor, Department. of Computer Science and Engineering, KLEF, AP, India

**Abstract**—Mobile Ad-hoc networks is an ascertaining domain with promising advancements, attracting researchers with a scope of enhancements and evolutions. These networks lack a definite structure and are autonomous with dynamic nature. The strength of the Ad-hoc network lies in the routing protocols making it an apt choice for transmission. With several types of routing protocols available our focus is on LGF (Location-based Geo-casting and Forwarding) protocol that falls in Position based category. LGF assures to grab the attention with its feature of low bandwidth consumption and routing overhead at the cost of unvolunteered attacks resulting in compromising the security of data. In our approach, we present a technique to overcome the profound attacks like Wormhole and Blackhole by aggregating LGF with k++ Means Clustering aiming at route optimization and promoting security services. The proposed mechanism is evaluated against QoS factors like End to End delay, Delivery Ratio, Load balancing of LGF using Simulator NS3.2 which envisioned drastic performance acceleration in the aforementioned model.

**Keywords**—MANET; LGF; K++ Means clustering; network security; wormhole attack; blackhole attack; secured algorithm

## I. INTRODUCTION

MANETs, a group of nodes that provide communication through wireless links without a predefined infrastructure and exhibiting dynamic nature has been the choice of practitioners and researchers for two long decades. The ad-hoc feature of MANETS makes it a favorable choice in several applications like Vehicular communication handling various disaster scenarios, defense, Security, and Online meetings. These applications depend on information exchange between nodes that play a vital in the process of Communication. The Crucial component acting as the backbone for node elucidation and improving the strength of MANETS ate routing protocols. It is a syntactic rule for defining a methodology to be undertaken by routers for transmission of data.

Based on the consideration surveyed by various researchers these protocols are classified as Topology based and Position-based. The former protocols rely on the respective structure of the network, whereas the latter originate with the location information of nodes. Topology based routing protocols fall into three well-known models like Proactive, Reactive, and Hybrid. Proactive protocols as the name reflect works based on prior information stored in the table in contrast reactive build a route on-demand when a request triggers. Bridging the gaps among both Hybrid protocols intersects the characteristics of Proactive and

Reactive [24]. Few well-known protocols falling on the Proactive side are DSDV, FSR, OLSR, and reactive are AODV, DSR, and TORA. ZRP, ZHLS, CEDAR are occupied under hybrid Class [23]. These protocols fail to outperform when the network turns to be densely populated with a huge number of nodes resulting in large network sizes, thereby lowering performance [20]. To leverage and sustain the network efficiency even with dense networks, MANETS impend on position-based routing protocols with urging requirement of security features [15]. Few protocols of interest are namely LAR, LGF, and Landmark. SLAR is also proposed to provide security against different attacks [16]. These protocols ensure efficient performance when clustered into zones [4]. Position based class mainly emphasizes the position of the node in the network and their performance is analyzed based on qualitative characteristics like Loop free, Decentralized operations, Path strategy, Performance metrics, Scalability, Reliable Delivery service, and Robustness. These Protocols support few strategies in packet forwarding namely Greedy Forwarding, Constrained directional flooding, and additionally Hierarchical or multilevel methods [2]. The greedy method of forwarding works by using optimization criteria for selecting the next node for the transmission of messages [22]. With directional flooding sender floods, packets to nodes toward the direction of destination satisfying predefined constraints and the Hierarchical method works when huge network scaling is on-demand [26].

Among several position-based routing protocols, our focus is on LGF that targets the reduction of routing overhead and bandwidth. In LGF the neighboring nodes in the forwarding zone perform rebroadcasting route request packet and acknowledge the source with route reply. Unlike all routing protocols, LGF is also vulnerable to serious attacks like Wormhole and Blackhole attack. These attacks exhibit an adverse impact on the performance of location-based ad-hoc networks [8].

Our paper enlightens all the fore-mentioned issues and proposes an enhanced approach based on the K++ Clustering technique to overcome attacks in LGF.

## II. RELATED WORK

Ahmad, Hameed, & Ikram, 2019, analyzed Ad-hoc networks and came up with a unique cluster-based algorithm for a reduction in the size of the routing algorithm. AI-Shrugan, Ghazali, & Hassan, 2012, gave a qualitative comparison of the position-based protocol in the context of the

greedy forwarding strategy. Alinci, Spaho, Lala, & Koli, 2015, reviewed MANETS related to clustering schemes like mobility-based, Energy-based connectivity with their pros and cons. Amouris, Papavasilliou, Maaloi, 1999, designed a protocol based on location routing zones which is efficient in the utilization of bandwidth for the huge size of networks. Dyabi, Hajami, & Allali, 2014, proposed the MANET clustering algorithm based on node density for cluster head selection. This approach promises improved results. Farjamnia, Gasimov, & Cavanshir, 2019, contributed a detailed review of handling the wormhole and analyzed its effect on the wireless network. Gayatri, et al., 2019 discussed in detail the wormhole attack in AODV and analyzed the framework by tracking the high transmission node. Giordano, Stojmenovic, 2004, presented a clear taxonomy on position-based routing models in Adhoc networks. Gupta, Singh, 2016, contributed a detailed study on wormhole attacks in wireless networks. Hamad, Kang, Jeon, & Nam, 2008 contributed to K-Means clustering in RDMAR protocol using the distance between nodes. Hossian et al., 2019 put a ray of light on the detection of a black hole in AODV and AOMDV adopting fusion of SHA-3 and Diffie-Hellman. Joo-Han song, et al., 2007 contributed the Secure Geographic Forwarding technique and SGLS with LRS (Location Reputation System) and comparison of their performance analysis. Kulkarni, Bukate, & Nanaware, 2018, provided an immense study on different attacks in Manets. Lattif, Ali, Ooi, & Fisal, 2005, proposed a detailed description of LGF implementation in MANETS with the GPS-FREE mechanism. Mahmood & Manivannan, 2018, discussed on Greedy Routing Protocol related to backtracking and compared performance issues with AODV and DSR protocol. Moudini, Er-Rouidi, Mouncif, 2016 evaluated secure Adhoc routing protocols categorized them into three types, and analyzed different protocols for secured and efficient routing in MANETS. Muthupriya, Revathi, & Rahman, 2017, designed a new algorithm SLAR enhancing security in LAR protocol against various types of malicious nodes. Patel A, Patel N, Patel R, 2015, proposed a Hash-based compression function based on a hash function for the RREQ packet with promising results. Priya Maidamour, Nekita chavan, 2012 surveyed and analyzed the vulnerable security threats like the Wormhole attack. Mishra, Gandhi, & Singh proposed a weighted forward method that is a fusion of forwarding, selection schemes of a node within a predefined area. Rajkumar Kapur, & Sunil Kumar Khatri, 2015, provided a detailed analysis of several vulnerabilities on routing protocols. Razaee, Yaghmaee, 2014, analyzed on cluster stability and proposed a Weight based algorithm for nodes with enhanced results. Royer, Toh Chai-Keong, 1999, reviewed about eight routing protocols, their functions, advantages, disadvantages, and provided a detailed comparison between these protocols, which helped our work in getting deeper. Sumit, Mitra, & Gupta, 2014, proposed an effective K-Means clustering and implemented IDS in MANETS using ZRP to avoid malicious activity. Srivastava, Daniel, Singh, & Saini, 2012, proposed a protocol for Energy Efficient Position-based routing with two new methods for route maintenance in Ad-hoc networks. Teotil, Dhurandher, Woungang, & Obaidat, 2015, proposed the COTA Approach in Position-Based Routing Protocol LAR1, which showed

efficiency in terms of security against the Wormhole attack. Yih-Chun Hu, Perrig, & Johnson, 2006, came up with a TIK protocol to handle Wormhole Attack in MANETS.

### III. EXISTING APPROACH

LGF: LGF with its beneficiary factors like lowering bandwidth and packet dropping rises to be the best choice for leveraging performance about measuring concerns like efficient packet forwarding in MANET's. Steps involved in LGF include path discovery and message forwarding [13, 22].

- The process initiates from a source with a multicast PREQ packet to all neighboring nodes based on the IP address of the destination. The protocol limits its range within a predefined distance.
- RREQ packet is forwarded to all the neighboring nodes with a distance less than the source node to the destination.
- The process repeats until the RREQ packet reaches the destination that further acknowledges the path to the source node from various intermediate nodes.
- Finally, the optimal shortest path is captured, and intended packets are transferred among source and destination.

Despite limitations with LGF when the range increases, it also suffers from a Wormhole attack that targets the shortest path with an illusion perspective. To handle this perturbation our algorithm fusions LGF with a clustering approach resulting in a secure, reliable transmission.

#### A. Attacks on ad-hoc Networks

An attack aims at compromising the security of transmission innumerable ways like Interruption, Interception, modification, Fabrication, or denial of service. A Wireless network is mainly prone to such type of attacks due to their dynamic nature [22]. Based on the method of disrupting security services [11], attacks are characterized by direct manipulation to the transmitted data, conversely passive as eavesdropping the communication between the parties [12]. Many attacks are figuring out, of which Wormhole attack and Black hole are considered [13,7].

#### B. Impact of Wormhole Attack

Limited availability of resources dispenses Ad-hoc network to attacks. An unauthorized entity with high power supply, memory, and computational capability is successful in introducing malicious attacks over MANET's [6].

A Wormhole attack is one worth enough to affect the network without revealing the cryptographic mechanisms embedded [9]. This attack has two variations that are hidden and exposed [29, 18].

1) *Hidden wormhole attack*: In this scenario the attacker succeeds in hiding the identity of the nodes between source and destination, creating an illusion of source and destination as one-hop neighbors [25].

2) *Exposed wormhole attack*: Here the attackers introduce themselves into the network with route discovery technique,

thereby exposing Wormhole nodes and hiding liable nodes between source and destination [28]. Based on these nodes different forms of Wormhole attacks are-encapsulated packet-based, Wormhole attack, out-of-band path, relaying of packets, and Protocol manipulation wormhole attack [21].

### C. Black Hole Attack

A serious problem endeavoring wireless sensor networks is the Blackhole attack characterized to absorb everything that comes on the way thereby decreasing the performance of the network [17]. In this attack, a malicious node or attacker node announces that it indexes the shortest path to the destination resulting in packet loss. It succeeds in communication failure among wireless networks and base stations. This attack results in topology modification including packet damage with forged routing information [23,10].

## IV. PROPOSED APPROACH

### A. Lgf with Clustering Approach

As LGF is prone to several attacks discussed and even restricted with range constraints. We propose a variation of LGF in combination with the clustering technique to handle the demerits of LGF. K++ Clustering is embedded in our proposed mechanism to strengthen the LGF for overriding the deficiencies [11].

**K++ Means:** This algorithm aims at clustering the given dataset into clusters and mainly focuses on seed or initial value selection, as an input to k-means. It overcomes the poor Clustering results of K-means which is an NP-hard problem. K-means gives the worst results for super polynomials in input and bad approximation of objective function in comparison with optimal clustering [5] that is overridden in k++ by the defined procedure to initial Clusters [1,3].

Step by Step Procedure for k++ is given as:

- 1) Select a data point C randomly.
- 2) For every data point P, Calculate the distance  $d(p)$  between P and C.
- 3) Pick a new data point based on weighted probability distribution with n proportional to  $d(p)^2$ .
- 4) Iterate steps 2 and 3 until optimal k centers are selected.
- 5) Continue with k-Means Clustering after the initial seed value is selected [27].

Position Based Protocols by virtue depends on the location of the node which ascertains performance assurance. These Protocols rely on three main sources to identify the exact location of the nodes, namely, based on signal strength, coordinates between nodes, and GPS based node location. Our approach uses the coordinates to locate the point of the node which helps in identifying malicious node location that promotes traffic bypassing within the network.

### B. Phases in Proposed Approach:

**Phase1:** The network is divided into clusters using the k++ Clustering technique, resulting in k value. The source node initiates the RREQ packet for route discovery when a

communication link is needed. This RREQ packet is forwarded to the nearest neighbors with the shortest Euclidean distance.

**Phase 2:** In this Phase Destination Node acknowledges the RREQ packet with RREP through the shortest path opted. The legitimacy of the RREP packet is judged or evaluated based on the predefined threshold value of RREP packets permitted.

**Phase 3:** Here a node is considered malicious once it violates the threshold value constraint of the RREP packet. The path the malicious node resides is excluded from the transmission path for forwarding packets. This phase also checks the hop count in the routing table to identify the path established by the hidden malicious node.

Our procedure succeeds in overcoming the Range constraint using k++ and Node Authentication by considering the location of the node as criteria using coordinates for calculating the Euclidean distance. Euclidean distance also adds to overcome the shortest path illusion injected by Wormhole and Black attacker nodes with the help of the RREP packet threshold value.

### C. Algorithm for Clustering enhanced LGF

There are some descriptions as given below to recover a safe RREP Packet through intermediate nodes.

Assume source node value = 1 and 2.

Intermediate nodes value = 1, 2.

Hop count node value = 3.

If (source node value == intermediate nodes value 1 and 2)

{  
Accept the RREP packet to the source node.

}

Else

{

The Source node discards the RREP packet because it is malicious node paths.

}

3. With this condition as benchmark source node waits and checks for safe route reply RREP packet through the intermediate nodes.

#### // location selection of node using k++ means

Require: k++means Function(cluster :list)

for k++meansFunction(cluster : list)

do  $X_k \leftarrow X_k + \text{cluster}.X$

$Y_k \leftarrow y_k + \text{cluster}.y$

end for

$XK \leftarrow XK / \text{long}(\text{cluster})$

$YK \leftarrow YK / \text{long}(\text{cluster})$

$N \leftarrow N - 1$  return (XK,YK)

#### // selection of cluster (based on the location of a node)

Require: k++means Function(cluster :list)

for (ifrom1tolong(Cluster))

do  $X \leftarrow \text{Cluster}.X$

$Y \leftarrow \text{Cluster}.Y$

$\text{distance} = \text{sqrt}((X_k - X) * (X_k - X) + (Y_k - Y) * (Y_k - Y))$

if (distance < distancemin) then

distance min=distance end if end for return (distancemin)

V. PERFORMANCE EVALUATION AND ANALYSIS

To Implement the Proposed approach using NS3.2, the Simulation parameters are initialized as shown in Table I. The performance of LGF with K++ Means is evaluated by considering the parameters like Load Balancing, End to End Delay, and Delivery ratio.

- 1) *End to End delay*: End to End delay is defined as the time incurred to travel from source to destination [19].
- 2) *Packet delivery ratio*: This is the ratio of the number of packets delivered to the number of packets sent by the source node [14].
- 3) *Load Balancing*: A parameter that defines an efficient distribution of transmission load during transmission in a network.

Fig. 1 shows a clear comparison of the reduced End to End delay factor proposed in comparison to the existing system. Here proposed system is indicated as lgfc-delay with a green spike.

Fig. 2 shows a comparison of the Delivery ratio which shows the packet loss of the Proposed and Existing system. Here proposed system Promises a reduced packet loss with a green spike.

TABLE I. SIMULATION PARAMETERS

| PARAMETERS              | VALUES                 |
|-------------------------|------------------------|
| Simulator               | NS 3.2                 |
| Simulator Time          | 120 s                  |
| Simulation Area         | 1000*1000 m            |
| Proposed Protocol       | LGF-C(hybrid approach) |
| Initial Energy of nodes | 1J                     |
| Number of Nodes         | 22                     |
| Bit Rate                | 1Mb/sec                |
| Packet Length           | 600 byte               |



Fig. 2. Graph Showing Delivery Ratio.

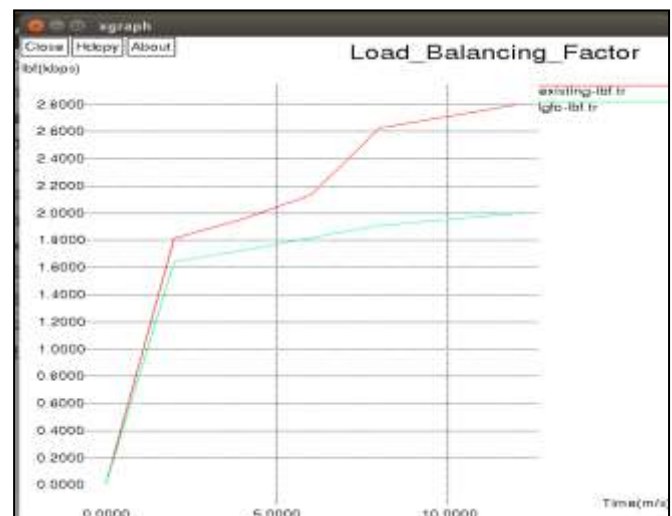


Fig. 3. Graph Showing Load Balancing Factor.

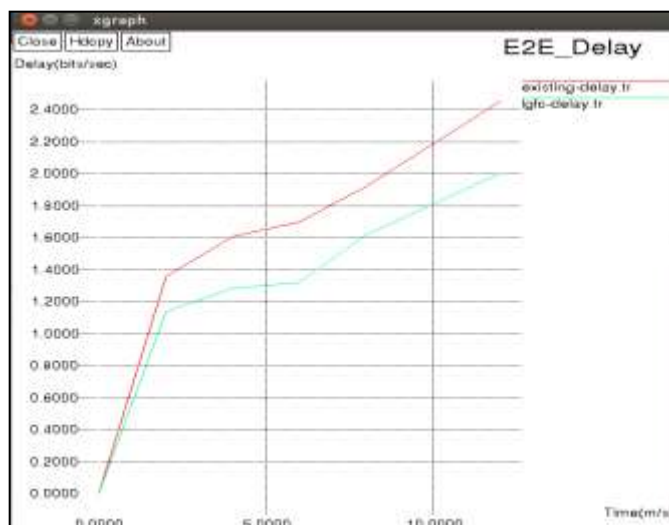


Fig. 1. Graph Showing End-To-End Delay.

Fig. 3 shows a comparison of the load balancing factor which projects a better performance by the pproposed system of lgfc indicated with a green spike.

Hence the simulation results of the proposed protocol outperform in terms of End to end delay, load Balancing, and in the reduction of Packet Loss by providing node authentication, Reliability and stability thereby leveraging the performance of the network in Position-Based Routing Protocols.

VI. CONCLUSION

This paper aimed to discuss the importance of MANETS and concentrated on the adverse effects of Wormhole and Blackhole attacks in the position-based routing protocol. LGF Protocol is studied for various setbacks related to delivery, avoiding Attacks, and providing Authentication of nodes with Location as a constraint. Our approach is considered a clustering-based method to overcome the Prior mentioned issues in LGF with enhanced K++ Mean’s supporting attack free and secure packet transmission in Wireless Ad-hoc Networks.

REFERENCES

- [1] M. Ahmad, A.Hameed, A. Ikram, & I. Wahid, "State of the art clustering schemes in mobile ad hoc networks: objectives, challenges, and future direction", *IEEE Access*, DOI:10.1109/access.2018.2885120, 2018.
- [2] M. Al-Shugran, O. Ghazali, & S. Hassan, "Performance Comparison of Position-Based Routing Protocols in the Context of Solving Greedy Failure", *International Conference on Advanced Computer Science Applications and Technologies (ACSAT)*. DOI:10.1109/acsat.2012.20, 2012.
- [3] M. Alinci, E. A. Spaho, A. Lala, & V. Kolici, "Clustering Algorithms in MANETs: A Review", *Ninth International Conference on Complex, Intelligent, and Software Intensive Systems*. DOI:10.1109/cisis.2015.47, 2015.
- [4] K. N. Amouris, S. Papavassiliou, & Li. Miao (n.d.), "A position-based multi-zone routing protocol for wide-area mobile ad-hoc networks", *IEEE 49th Vehicular Technology Conference (Cat. No.99CH36363)*. DOI:10.1109/vetec.1999.780570, 1999.
- [5] M. Dyabi, A. Hajami, & H. Allali, "A new MANETs clustering algorithm based on nodes performances", *International Conference on Next Generation Networks and Services (NGNS)*. DOI:10.1109/ngns.2014.6990222, 2014.
- [6] G.Farjamnia, Y.Gasimov & K.Cavanshir "Review of the Techniques Against the Wormhole Attacks on Wireless Sensor Networks", *Wireless Personal Communications*, volume 105, pages1561–1584,
- [7] S. Hossain, M. S. Hussain, R. R. Ema, S. Dutta, S. Sarkar, & T. Islam, "Detecting Blackhole attack by selecting appropriate routes for authentic message passing using SHA-3 and Diffie-Hellman algorithm in AODV and AOMDV routing protocols in MANET", *10th International Conference on Computing, Communication and Networking Technologies (ICCCNT)*, DOI:10.1109/icccnt45670.2019.8944395, 2019.
- [8] S. Gayathri, R. Seetharaman, L. H. Subramanian, L. H., Premkumar, S., Viswanathan, S., & Chandru, S., "Wormhole Attack Detection using Energy Model in MANETS", *2nd International Conference on Power and Embedded Drive Control (ICPEDC)*. DOI:10.1109/icpedc47771.2019.9036536, 2019.
- [9] N. Gupta, & S. N. Singh, "Wormhole attacks in MANET", *6th International Conference - Cloud System and Big Data Engineering (Confluence)*. DOI:10.1109/confluence.2016.7508120, 2016.
- [10] O. F. Hamad, M. Y. Kang, J. H. Jeon, & J. S. Nam, "Neural Network's k-means Distance-Based Nodes-Clustering for Enhanced RDMAR Protocol in a MANET", *IEEE International Symposium on Signal Processing and Information Technology*. DOI:10.1109/isspit.2008.4775666, 2008.
- [11] R. K. Kapur, & S. K. Khatri, "Analysis of attacks on routing protocols in MANETs", *International Conference on Advances in Computer Engineering and Applications*. DOI:10.1109/icacea.2015.7164811, 2015.
- [12] A. N. Kulkarni, R. R. Bukate, & S. D. Nanaware, "Study of Various Attacks and Routing Protocols in MANETS", *International Conference on Information, Communication, Engineering and Technology (ICICET)*. DOI:10.1109/icicet.2018.8533696, 2018.
- [13] L. A. Latiff, A. Ali, Chia-Ching Ooi, & N. Fisal, "Location-based geo casting and forwarding (LGF) routing protocol in mobile ad hoc network", *Advanced Industrial Conference on Telecommunications/Service Assurance with Partial and Intermittent Resources Conference/E-Learning on Telecommunications Workshop (AICT/SAPIR/ELETE'05)*. DOI:10.1109/aict.2005.55, 2005.
- [14] B. A. Mahmood, & D. Manivannan, "GRB: Greedy Routing Protocol with Backtracking for Mobile Ad-Hoc Networks", *IEEE 12th International Conference on Mobile Ad Hoc and Sensor Systems*. DOI:10.1109/mass.2015.49, 2015.
- [15] H. Moudni, M. Er-round, H. Mouncif, & B. E. Hadadi, "Secure routing protocols for mobile ad hoc networks", *International Conference on Information Technology for Organizations Development (IT4OD)*. DOI:10.1109/it4od.2016.7479295, 2016.
- [16] V. Muthupriya, S. Revathi, & B. S. A. Rahman, "Secure Location Aided Routing (SLAR) for mobile ad hoc networks", *IEEE International Conference on Power, Control, Signals and Instrumentation Engineering (ICPCSI)*. DOI:10.1109/icpcsi.2017.8392279, 2017.
- [17] A. Nabou, M. D. Laanaoui, & M. Ouzzif, "Evaluation of MANET Routing Protocols under Black Hole Attack Using AODV and OLSR in NS3", *6th International Conference on Wireless Networks and Mobile Communications (WINCOM)*. DOI:10.1109/wincom.2018.
- [18] A. Patel, N. Patel, & R. Patel, "Defending against Wormhole Attack in MANET", *Fifth International Conference on Communication Systems and Network Technologies*, DOI:10.1109/cnsnt.2015.253, 2015.
- [19] PanelJoo-HanSong, Vincent.W.S.Wong, & Victor C.M.Leung "Secure position-based routing protocol for mobile ad hoc networks Ad Hoc Networks", Volume 5, Issue 1, January 2007, Pages 76-86
- [20] Priya Maidamwar, Nekita Chavan, "A Survey on Security Issues to Detect Wormhole Attack in Wireless Sensor Network International Journal on AdHoc Networking Systems, volume 2 issue(4):37-50 · October 2012 DOI: 10.5121/ijans.2012.2404
- [21] Priya Mishra, Charu Gandhi, & Buddha Singh, "An Improved Greedy Forwarding Scheme in MANETs Technology", *JIIT, Noida, U.P., Ind, journal of telecommunications and information technology*, pp. 50-55, 2017.
- [22] Rajinder Singh, Parvinder Singh, Manoj Chavan "An effective implementation of the security-based algorithmic approach in mobile Adhoc networks", *Human-centric Computing and Information Sciences*, volume 4, Article number: 7 (2014)
- [23] M. Rezaee & M. H. Yaghmaee, "A new clustering protocol for Mobile Ad-Hoc Networks", *International Symposium on Telecommunications*, DOI:10.1109/istel.2008.4651331, 2008.
- [24] E. M. Royer & Toh Chai-Keong, "A review of current routing protocols for ad hoc mobile wireless networks", *IEEE Personal Communications*, 6(2), 46–55. DOI:10.1109/98.760423, 1999.
- [25] S. Giordano, I. Stojmenovic, "Position-Based Routing Algorithms For Ad hoc Networks: A taxonomy", *Ad Hoc Wireless Networking*, pp 103-136, volume 14, 2004.
- [26] S. Srivastava, A. K. Daniel, R. Singh, & J. P. Saini, Energy-efficient position-based routing protocol for mobile ad hoc networks. *2012 International Conference on Radar, Communication, and Computing (ICRC)*, DOI:10.1109/icrc.2012.6450540, 2012.
- [27] S. Sumit, D. Mitra, & D. Gupta, "Proposed Intrusion Detection on ZRP based MANET by effective k-means clustering method of data mining", *International Conference on Reliability Optimization and Information Technology (ICROIT)*, DOI:10.1109/icroit.2014.6798303, 2014.
- [28] V. Teotia, S. K. Dhurandher, I. Woungang, & M. S. Obaidat, "Wormhole prevention using COTA mechanism in position based environment over MANETs", *IEEE International Conference on Communications (ICC)*, DOI:10.1109/icc.2015.7249448, 2015.
- [29] Yih-Chun Hu, A. Perrig, & D. B. Johnson, "Wormhole attacks in wireless networks", *IEEE Journal on Selected Areas in Communications*, 24(2), 370–380, DOI:10.1109/jsac.2005.861394, 2006.

# A Proposed User Requirements Document for Children's Learning Application

Mira Kania Sabariah<sup>1</sup>

Department of Electrical and Information Engineering  
Universitas Gadjah Mada, Yogyakarta, Indonesia  
School of Computing, Telkom University  
Bandung Indonesia

Paulus Insap Santosa<sup>2</sup>, Ridi Ferdiana<sup>3</sup>

Department of Electrical and Information Engineering  
Universitas Gadjah Mada  
Yogyakarta, Indonesia

**Abstract**—User requirements are the highest level of requirements. Flawed user requirements document can cause defects in the software being built—aspects of applications that were not presented in the user requirements document to cause a defect. In learning applications for children, there are aspects of pedagogy that need to be well documented. This aspect is not available in the general user requirements document, so it is often not well presented. The learning style and thinking skills level is crucial to be well presented in the user requirements document. That was because the children's persona cannot be compared at every range criteria of developmental age. That factor will undoubtedly affect the specifications of the software to be built. Users' viewpoints about different requirements can also make developers wrong in determining requirements. Applying requirements prioritization in the user requirements document can help resolve the problem. Measurement of document quality was also performed using parameters in measuring the quality of the user requirements document. The results of measuring the quality of the user requirements document found that it is reliable for use.

**Keywords**—User requirements; user requirements document; learning application; aspect of pedagogy children

## I. INTRODUCTION

Requirement documents are needed to verify and validate software requirements [1], [2]. A requirements document was also made to facilitate a series of requirements engineering activities. Changing requirements in the requirements engineering process often cause chaos [3]. The required documents' existence is a medium for communication between the team and stakeholders [4]. Software requirements specifications (SRS) need to be known by application developers and users through a document [5]. Therefore, requirements documents that could present in detail and as required. That matter was needed to achieve the objectives and quality of the application being built.

Software requirements were divided into three parts: user requirements, business requirements, and software requirements specifications [6], [1]. User requirements are the highest level in the requirements and were obtained from the results of the user's point of view. Failure to document requirements may occur presented in natural language [7]. Natural language is often used in user requirements because it can be the primary means of communication between stakeholders and developers [1]. However, problems such as

misunderstanding, inaccuracy, ambiguity, and inconsistency are the causes of failure [8]. Activities in requirements elicitation need documents to make requirements prioritization easier for the development team [8]. The process of requirements elicitation to produce valid software requirements specifications is not easy [9]. Requirements document are often misinterpreted, misunderstood, and not well documented [10]. That can happen due to the unavailability of components from specific aspects of the type of application that will be built on the general requirements document template.

User requirements for certain types of applications will certainly be different. These differences can occur when the application domain to be built specific aspect [11]. In the learning application, the particular aspects were a pedagogy aspect. Pedagogical aspects need to be present, and each attribute value must be written clearly. These aspects will have different values because they were influenced by the user's persona, such as in adult and child users. The two types of users have different characteristics that can be expressed in persona. Persona influences the value of pedagogical aspects, especially on learning style. Children have more diverse characteristics because they were influenced by the range criteria of their development age. While for adults, there are no range criteria of age. These differences will affect the learning style. In addition to learning styles, children's level thinking skills also need to be considered based on the range criteria of the age of development. That is because children in each range criteria of developmental age have different cognitive abilities.

The different pedagogical aspects need to be well presented in the user requirements document [11]. These aspects were needed so that learning outcomes from learning were achieved. In the application of children's learning, positive, psychomotor, and emotional aspects need to be well defined, and the range criteria of the children's development [12]. Various forms of learning applications will undoubtedly affect the document structure of user requirements. The problem is how to present a user requirements document that matches the characteristics and type of a children's learning application. Compilation of user requirements documents for children's learning applications was expected to guide the elicitation team in exploring the needed aspects. The quality of user document requirements needs to be measured so that the document's legibility was fulfilled. There are often defects in



the presentation of the user requirements document [13]. That condition caused the resulting application does not match the expectations of the application user. The next problem is how to measure the quality of the user requirements documents created. Quality measurements were carried out so that documents can be understood by application developers clearly and correctly.

Based on these problems, the research will focus on how to provide a guide regarding user requirements documents (URD) for children's learning applications. The URD was expected to make it easier for the elicitation of the child learning apps team to define a set of user requirements. The URD also presents a collection of aspects that need to be determined when building children's learning applications. The URD gives requirements prioritization so that it can reduce conflicts when requirements were made. The validity of the URD also needs to be measured to determine the legibility and clarity documents.

## II. RELATED WORK

### A. User requirements Document (URD)

User requirements were often referred to as user needs. Describe what the user does with the system. User requirements were proper if they were obtained directly from the user and state the domain's properties generated by introducing a new system [13]. The user requirements document is an artifact that contains a set of requirements obtained based on the views of the user [14]. In the requirements engineering phase, the requirements document preparation process is carried out [15]. The goal is that the developer gets clear information regarding the system requirements to be built. The user requirements document contains specifications of the application software requirements to be built. Software Specification Requirements are possible to develop due to the type of software project. The change occurred because of the inaccuracies and shortcomings of the SRS [16].

User requirements documents were generally written using natural language[8]. This language often makes documents present ambiguous, inaccurate, and unclear information [8]. These conditions cause differences in understanding between the requirements engineering team and the application development team [14]. Other problems, it is crucial to consider the presentation of the requirements prioritization in the user requirements document. Requirements prioritization can be taking into several variables, including time, staff, and costs [17].

### B. Children's Learning Application (CLA) vs. Adult's Learning Application (ALA)

In learning applications, pedagogical aspects need to be well defined, such as learning outcomes and learning styles [18]. How to learn in each individual has a difference. That difference occurs because it was influenced by the personality of each individual and influences the learning process. In adults and children, the difference is noticeable. Children who have this range of criteria of age development cannot be equated at every age level. Different range criteria positively affect aspects of pedagogy, such as learning styles and

thinking skills level. According to experts in children's learning and literature review, learning styles for children need to be in the form of visuals, audio, read/write, and kinesthetic (VARK) [18], [19]. While in adults, there is no type of age. Learning style differences are formed based on their experience in learning. The concept of andragogy was often used as a reference for determining adult learning styles [20], [21].

The differences between a children's learning application and an adult's learning application can be distinguished based on the persona. According to Piaget's, children have four range of criteria of developmental age [22]. The Psychomotor, cognitive, and emotional development of children who are different in each range criteria of age development becomes something to consider in building learning applications[23]. The pedagogical aspects that significantly influence children's learning are learning style and thinking skills level. The reason is that the child is in developmental age and does not have experience in the learning process. In addition to learning style, thinking skills level in children's learning needs to be considered. Limitations of cognitive abilities at every range of criteria of development affect children's level of thinking skills. In the learning process, children also need to be given an award. Appreciation is the basis for children's motivation to learn[24]. While in adults, the range of age is not a measure to determine of learning style. Learning experiences that affect learning behavior in adults. It also affects the learning style of adults. Andragogy is a learning style that is suitable for adults [20] because adult learning aims to enrich their knowledge to solve their problems.

The differences in the persona, which is influenced by the value of pedagogical aspects in children and adults, is the reason for differences in learning application. The difference in learning styles will affect the implementation aspects. Children's cognitive limitations also affect the way children can quickly receive information. The presentation of objects in the application needs to be adjusted by the range criteria of the development age. Children's learning application was made to help children in the learning process [25]. There are two types of applications that tend to be made for children's learning applications based on interviews with five child education application developers. That type of application is in the form of a game and simulation (non-game). Both types of applications have different characteristics and approaches to the development process. That has an impact on the aspect requirements that need to be controlled.

In adult learning applications, a feature of material selection and material source selection needs to be provided. That is because adults do learn to solve problems. Although there are differences in the three aspects' value, there are slices in the two learning applications. Aspects of generic environment issues need to be defined, for example, Platform applications. That is because children still have limitations in psychomotor. In children's and adult learning applications, learning outcomes need to be determined. That attribute also needs to be in the user requirements document. That is because the learning outcome is exposure to the form of the content presented in the application. Fig. 1 is a Venn diagram of the differences between CLA and ALA.

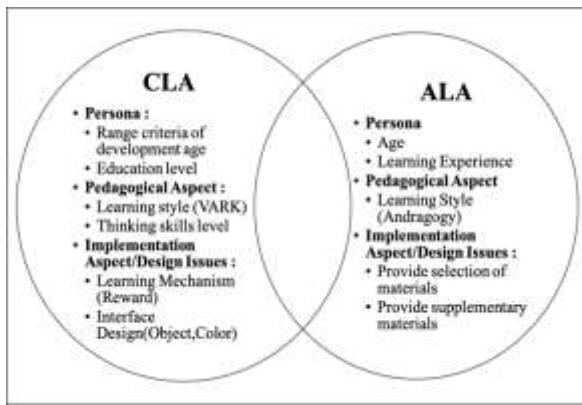


Fig. 1. Venn Diagram of the differences between CLA vs. ALA.

### III. RESEARCH METHOD

The methodology used in proposing the URD was divided into two stages. The first stage is structuring the URD for children's learning applications, and the second stage is to measure the quality of the proposed URD.

- In the first stage, the preparation of the URD was carried out. The first stage is the stage carried out to answer RQ1. An analysis of the URD for general application and characteristics of learning applications for children. The analysis is done by looking at aspects of the requirements that need to be present from the application of children's learning using literature review and interviews with an expert. Then after that was found, the URD structure is made. URD Structure was made based on aspects of user requirements and learning applications for children. Then, the URD was implemented in an elicitation application. The app is an application to assist the elicitation team in gathering needs. The formed URD also presents requirements prioritization for each aspect using ranking methods.
- The second stage is a stage to answer RQ2. They made appropriate measuring tools to carry out URD quality measurement in measuring the quality of user requirements. In compiling the measuring instrument, an approach was made using the user requirements document's quality aspect. When the measuring instrument has been formed, then the reliability test is performed using Cronbach's Alpha. If the reliability has been fulfilled, the next step is to measure the URD generated from the requirements elicitation process using interval analysis. Based on the goal, properties, and properties, a measuring instrument was made in the form of a questionnaire. Measurements were made on all aspects of the URD in Table III, measured based on each property. In the game application, 17 questions were represented by variable questions P1 through P17. Meanwhile, the non-game application questionnaire consisted of 20 items (P1-P20). The rating of each property was done by using Likert 1-5. Cronbach's Alpha was conducted for the reliability test of the questionnaire created. While the result of data from filling out the questionnaire was processed using interval analysis with the range of values listed in Table II.

#### C. Quality of user Requirements Document

The cause of an application failure is due to a defect in the collection and identification of user requirements [26]. The quality of the user requirements document needs to be considered so that the application developer clearly understands it. Measurements should be taken so that the document is reliable for use. Cronbach's Alpha will be used to measure reliability. Quantitative measurement is done by taking into account the quality aspects of the software requirements specification (SRS), namely, (i) Requirements Sentences Quality (RSQ) and (ii) Requirements Document Quality (RDQ) [15]. The RSQ aspect was measured to see the syntactic quality of a single sentence considered separately. RDQ is measured to determine the quality of sentences considered in the context of all the requirements documents.

Each goal property is measured using properties, as can be seen in Table I. RSQ's goal properties have non-ambiguity, completeness, and understandability properties. In contrast, RDQ has completeness and understandability properties. The following is an explanation of each RSQ properties:

- Non-Ambiguity: the ability of a Requirement to have a unique interpretation.
- Completeness: the ability of each requirement to make references to precisely identified entities.
- Understandable: the ability of each requirement to be fully understood when used to develop software.

That is an explanation for the properties of RDQ:

- Completeness: Requirements Specification document can avoid potential or actual differences.
- Understandable: Requirements Specification document can be fully understood when read by the user.
- The Quality Model Goal Properties and the Related Properties.

| Goal Properties                      | Properties        |
|--------------------------------------|-------------------|
| Requirements Sentences Quality (RSQ) | Non-Ambiguity     |
|                                      | Completeness      |
|                                      | Understandability |
| Requirements Document Quality (RDQ)  | Completeness      |

TABLE I. THE CATEGORY OF INTERVAL VALUE

| Interval Value | Description |
|----------------|-------------|
| 0% - 19.99%    | Very Bad    |
| 20% - 39.99%   | Bad         |
| 40% - 59.99%   | Neutral     |
| 60% - 79.99%   | Good        |
| 80% - 100%     | Very good   |

#### IV. RESULTS AND DISCUSSION

##### A. User Requirements Document of Children's Learning Application

User requirements were needed as a step to compile the system requirements so that they can describe in detail how the system must be run [1], [27]. Different user points of view need to be recorded to be modeled in the system correctly accurately. User requirements describe the user class and what users need [26], [13]. User class describes the user's profile or persona, such as age, gender, user experience. The user needs to explain what the user needs from the application to be made. In the case of learning applications, users can convey related forms of learning that need to be in the app—for example, the material presented and the learning style in the application. User requirements document presented in natural language will be challenging to show a different user perspective [7]. Understanding child learning application developers of application requirements that need to be explored also becomes an obstacle to the quality of user requirements document.

Based on interviews with learning application developers, children's learning applications can be made in games and simulations. The kind of application was determined based on the type of learning material to be delivered in the app. A simulation was used when the material to be presented. That is conveying the conditions of the situation in the real world [28]. The form was considered to provide considerable learning potential because it is more effective and interactive [29] for the kind of games made in serious games. Serious games are tools that are considered useful in the learning process [30], [31]. Serious games for a child can be used for several things, including increasing motivation to learn, stimulating physical activity, solving behavioral problems, and helping with therapy-related to health problems [18]. Both forms of application have a different structure of requirements aspect. The difference in aspects structure will undoubtedly have an impact on the user requirements document. The user requirements document will be adjusted according to the aspects structure of the application formed. Table III explains the structures of aspect requirements for game learning and non-game applications. Each aspect has attributes that can provide a detailed description of the learning application's user requirements to be built. Each aspect's attributes contain one or more values that were translated into user requirements. Each aspect's attributes can be seen in Table IV for game learning applications [18] and Table V non-game learning applications [32].

Requirements of user-profiles and application platform preferences are fundamental attributes that need to be explored from the user. The values of several attributes have been presented in the elicitation application. It is making it easier for the elicitation team and participants to define requirements. The value of some attributes has been determined based on established by the conditions—for example, the learning style attribute's value. The attributes of the learning style and thinking skills level need to be elaborated on pedagogical aspects. A learning style must

determine the children's preference for how the learning material was presented [33]. Knowing the learning style will make it easier for developers to design learning applications. The thinking skills level was created to limit the cognitive level, adjusted to the range criteria of the child development [23], [34]. Thinking skills level was also used to direct in achieving learning objectives. The application's thinking skills level was based on cognitive processes that refer to taxonomy blooms [35].

The attributes of the two aspects of the application are presented in table form and filled with several user points of view as participants in the requirements elicitation process. The structure was to accommodate a set of requirements from many users. That form also facilitates the readability of the information presented. Fig. 2 is an example display of the database structure of the elicitation application that was built. The elicitation application was created as a tool to assist the elicitation team.

The team can be used the apps when collecting elicitation requirements and automatically generating user requirements documents. The app also makes it easy for the elicitation team to change requirements quickly. The change also directly occurred in the requirements document. Thus, the agile concept can be applied in technical terms and in the user requirements document. The aspect parts of the results of filling each of these attributes were then made requirements prioritization.

The requirements prioritization generated are then written down on each aspect of the user requirements document's requirements. Simultaneously, the data collection results were stored in an attachment to the user requirements document. Requirements prioritization formed using a formula that was combined from several attributes. Then in each aspect, Requirements prioritization is done using ranking techniques tailored to their attributes [36]. For example, platform applications in the context of use ranking will be made to get requirements prioritization. The requirements prioritization displayed in the user requirements document were performed to display the required requirements based on the user's point of view without causing conflicts [37]. Fig. 3 is an example of a user requirements document for the context of use. In this aspect, the results of processing requirements prioritization were explained from the results of data collection. Fig. 4 is a flowchart for requirements prioritization for the context of use using ranking methods.

The proposed URD for children's learning applications has been made. When the developer uses the URD, that question can be answered by looking at Fig. 4. Fig. 4 explained that the URD could be suitably used when the application to be built is a learning application. Besides the type of learning application, it needs to be seen for whom the application was made. If the user is a child, the proposed URD can be used. Nevertheless, if the user is an adult, it is better to use URD for adult learning applications. The same thing also applies when the application to be made is not a learning application, so it better used the standard URD.

TABLE II. THE RELATIONSHIP BETWEEN EACH ASPECT

| User Requirements | Description   | The aspect of Game Application | The aspect of Non-Game Application |
|-------------------|---|--------------------------------|------------------------------------|
| User class        | It was related to the user group of the application to be built. This section will describe the user's profile or persona.  | User aspect                    | Generic environment issues         |
| User need         | Related to user needs for the type of application to be built. A set of requirements that user needs in the application need to be described and made based on user preferences | Context of use                 | Learning context                   |
|                   |   | Pedagogical aspect             | Learning experience                |
|                   |   | Game aspect                    | Learning objective                 |
|                   |   | Implementation aspect          | Design Issues                      |

TABLE III. STRUCTURE ASPECT OF A GAME APPLICATION

| Aspect                | Description   | Attribute  |
|-----------------------|---|--|
| User aspect           | The user aspect is the initial stage, aiming to determine the user's characteristics, especially those related to the learning process. | <ul style="list-style-type: none"> <li>- Name</li> <li>- Age</li> <li>- Sex</li> <li>- Education level</li> <li>- A course like and unlike</li> </ul>  |
| Context of use        | The context of use describes the specifications of the application platform.  | <ul style="list-style-type: none"> <li>- Platform application</li> </ul>   |
| Pedagogical aspect    | The pedagogical aspect contains information about learning patterns to be able to achieve the objectives of learning.                   | <ul style="list-style-type: none"> <li>- Learning outcome</li> <li>- Detail learning outcome</li> <li>- Thinking skills level</li> <li>- Difficulty</li> <li>- Learning Mechanics</li> </ul> |
| Game aspect           | The game aspect was explained related to issues related to develop game applications.   | <ul style="list-style-type: none"> <li>- Game genre</li> <li>- Game mechanics</li> <li>- Game format</li> <li>- Game form</li> </ul>   |
| Implementation aspect | Implementation aspects describe essential aspects needed in the implementation of application   | <ul style="list-style-type: none"> <li>- Implementation element</li> </ul>   |

TABLE IV. STRUCTURE ASPECT OF A NON-GAME APPLICATION

| Aspect                     | Description   | Attribute   |
|----------------------------|---|---|
| Generic environment issues | Generic environment issues explain matters related to user-profiles and the use of the digital platform.  | <ul style="list-style-type: none"> <li>- Name</li> <li>- Age</li> <li>- Sex</li> <li>- Education level</li> <li>- A course like and unlike</li> <li>- Platform application</li> </ul>   |
| Learning context           | Learning context explains learning activities, learning facilities, and collaboration.  | <ul style="list-style-type: none"> <li>- Learning outcome</li> <li>- Learning Objective</li> </ul>  |
| Learning experience        | Learning experience explains related to the learning experience that will be provided to users.   | <ul style="list-style-type: none"> <li>- Learning content</li> <li>- Result and feedback on learning</li> <li>- Aim and target of learning</li> <li>- Representation/storyline</li> <li>- Social interaction</li> </ul>         |
| Learning objective         | Learning objectives explain the purpose of learning, whether to improve abilities or add new abilities.   | <ul style="list-style-type: none"> <li>- Learning activity</li> <li>- Learning facility</li> <li>- Collaboration</li> </ul>   |
| Design Issues              | Design Issues aims to determine child preferences related to objects, colors, text, navigation, and multimedia presentation to help design the learning application interface that will be built. | <ul style="list-style-type: none"> <li>- Interface design object</li> <li>- Interface design color</li> <li>- Interface design type of text</li> <li>- Interface design navigation</li> <li>- Interface design voice</li> </ul> |

| Table Name          | Columns |
|---------------------|---------|
| pedagogical_aspect  | 1       |
| learning_outcome    | 1       |
| learning_objective  | 1       |
| learning_mechanic   | 0       |
| learning_experience | 1       |
| learning_context    | 1       |
| image_object        | 2       |
| genre_game          | 1       |
| game_mechanic       | 1       |

Fig. 2. Example of the Database Structure of the Elicitation App.

**2.Context of Use**  
This section explains the platform related to the application to be built based on recommendations by participants. Requirements prioritization of the platform based on the results of data collection from users can be seen in table 2. Smartphone have the highest rank of platform applications based on the results of prioritization requirements process.

Table 2. Requirements prioritization for platform applications

| No | Platform   |
|----|------------|
| 1  | Smartphone |
| 2  | Tablet     |

Fig. 3. Example of the user Requirements Document.

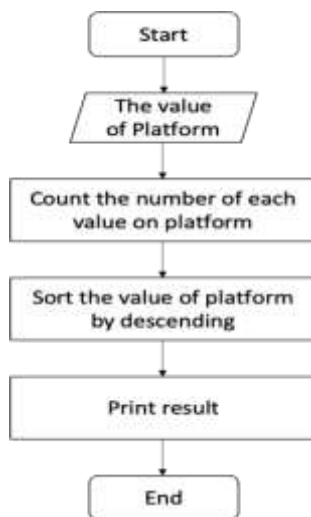


Fig. 4. Flowchart of Requirement Prioritization.

## B. Experiment and Results

1) *Participant*: As respondents involved in this study, participants were divided into two types, namely, child participants as users and application practitioner participants. The child participant is required to convey the requirements of the learning application to be built. There are 32 children aged 6-8 years who will express their wishes through the elicitation application that will be built. Practitioner application participants are participants who are involved in measuring document quality. Participants consisted of 37 respondents

who had experience in building children's learning applications.

2) *Materials*: In this study, there are materials used to assist in experiments. The materials used in the experiment are:

- Elicitation apps. This application was built to facilitate the elicitation team in communicating with children. The application was built in the form of mobile-based applications. The use of mobile technology has the effectiveness of interacting with children [36], and does not require a considerable cost [37]. Applications were built according to the characteristics of the child. Children also feel fun and joy when conveying their desires by using a mobile-based app [38]. Applications were also made to facilitate the elicitation team in documenting requirements.
- Questionnaires were used to measure the quality of the resulting URD.

3) *Case study*: Two cases were used in measuring URD. The aim is to produce URD game and non-game applications. The case for game applications is about introducing types of vegetables and fruits. While for the example of a non-game is about the introduction of rain. In the elicitation application, material choices were given for each instance with evaluation questions included. The material was presented with four learning styles: visual, audio, read/write, and kinesthetic. The thinking skills level gave evaluation questions.

Requirements elicitation process doing with 32 children with a span of about one month. Each child respondent expresses their needs through interaction through elicitation applications. After the elicitation process, the requirements were carried out, and then the URD processing is done through the app. The output of the use is URD by presenting information related to user requirements based on user preferences. Fig. 3 is an example of URD results of game applications. In the aspect of the context of use, the requirements prioritization for the application platform attribute are smartphones.

4) *Results*: The URD that has been generated from the application was then distributed to the app practitioner participants. A total of 37 respondents then studied URD from both types of applications and conducted an assessment through the quality questionnaire URD prepared. The questionnaire was filled in online. Respondents were asked to rate the URD according to the type of application. The assessment was done according to each goal properties for each aspect of the application type.

Questionnaire data processing was performed using Cronbach's Alpha. The results of data processing for game type applications obtained a questionnaire reliability test of 0.923. With this value, it can be concluded that the questionnaire has reliability. URD quality assessment for game applications found that most aspects have answer values in 80%-85% (very good). There are five aspects, namely, P3, P7, P14, P16, and P17, to answer value results in the range of

ethical values. Questions P3, P4, and P7 are questions that are in the RSQ goal properties. That means that the sentences of the three attributes have a good understanding of each sentence. While for P16 and P17 are questions in the RDQ goal properties. Both of these questions state in terms of the document as a whole is complete and well understood. The value of the answer from each aspect can be seen in Fig. 5. Based on that result, it can be concluded that the URD's general for game applications is perfect. That means that application developers can understand the user requirements for the children's learning application to be built.

Questionnaire data processing was performed using Cronbach's Alpha. The results of data processing for non-game type applications obtained a questionnaire reliability test of 0.946. With this value, it can be concluded that the questionnaire has reliability. As for the interval analysis results, it was found that 19 aspects had an answer value in the range of 80% -85% (perfect), and only one aspect, namely P10, had a value of 79% (excellent).

Based on these results, it can be concluded that the sentence in every aspect of the URD for the non-game application can be understood very well. That means that the presentation of data that combines sentences in the form of natural language and tables helps the developer understand user requirements. The availability of prioritization requirements also makes it easy for developers to decide on user requirements. However, the learning experience attribute assessment has an excellent rating (P10), the resulting general URD. The value of the answer from each aspect can be seen in Fig. 6.

## V. CONCLUSION

The conclusions of the research activities that have been carried out are as follows:

- The document's pedagogical aspects in detail help the development team when design learning applications from the user's point of view.
- The availability of requirements prioritization in documents also helps to reduce conflicts when developing the system.
- URD quality measuring instruments compiled have the reliability to be used in measuring document quality. That was evidenced by Cronbach's alpha measurements for games that are 0.923 and non-games is 0.946.
- URD measurement results for both types of children's learning applications are generally excellent. That was proof from each variable gives an average value range of 80-85%. In other words, URD can be understood by application developers.

The future work was to complete the URD by adding a form of notation for several attributes. The aim is that all attributes are understood very well by a child's learning application developers.

## ACKNOWLEDGMENT

This research was supported by The Indonesia Endowment Fund for Education (Lembaga Pengelola Dana Pendidikan/LPDP).

## REFERENCES

- [1] S. O. Mokhtar, R. Nordin, Z. A. Aziz, and R. M. Rawi, "Issues and Challenges of Requirements Review in the Industry," *Indian J. Sci. Technol.*, vol. 10, no. 3, 2017.
- [2] S. Smith, "Systematic development of requirements documentation for general purpose scientific computing software," *Proc. IEEE Int. Conf. Requir. Eng.*, pp. 205–214, 2006.
- [3] A. Zainol and S. Mansoor, "Requirements Management Tool Elements For The Malaysian Software Industry," *J. Inf. Commun. Technol.*, vol. 11, no. 1, pp. 179–192, 2012.
- [4] F. Paetsch, A. Eberlein, and F. Maurer, "Requirements engineering and agile software development," in *WET ICE 2003. Proceedings. Twelfth IEEE International Workshops on Enabling Technologies: Infrastructure for Collaborative Enterprises, 2003.*, 2003, pp. 308–313.
- [5] M. dos S. Soares and D. S. Cioquetta, "Analysis of Techniques for Documenting User Requirements," in *International Conference on Computational Science and Its Applications, 2012*, pp. 16–28.
- [6] J. J. Cho, "An Exploratory Study on Issues and Challenges of Agile Software Development with Scrum," *Issues Inf. Syst.*, vol. 9, no. 2, p. 599, 2010.
- [7] P. Thitisathienkul, "Quality Assessment Method for Software Requirements Specifications based on Document Characteristics and its Structure," in *2015 Second International Conference on Trustworthy Systems and Their Applications, 2015*, pp. 51–60.
- [8] C. Coulin and D. Zowghi, "Requirements Elicitation: A Survey of Techniques, Approaches," *Eng. Manag. Softw. Requir.*, pp. 19–46, 2005.
- [9] F. P. Brooks, "No Silver Bullet: Essence and Accidents of Software Engineering," *Computer (Long. Beach. Calif.)*, vol. 20, no. 4, pp. 10–19, 1987.

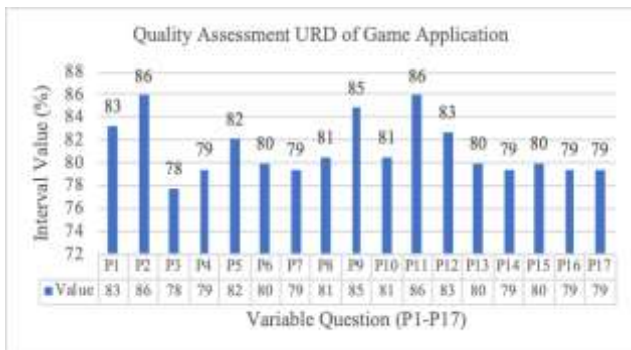


Fig. 5. Chart of Quality Assessment URD for a Game Application.

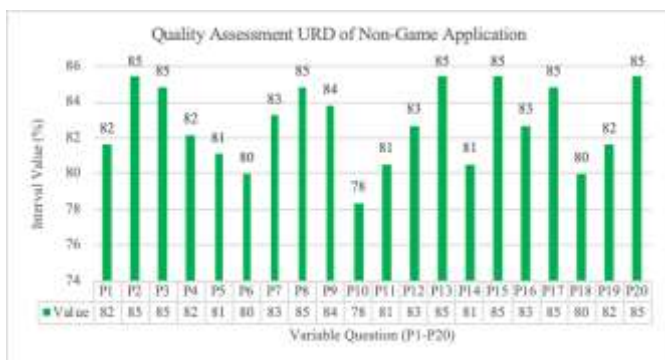


Fig. 6. Chart of Quality Assessment URD for a Non-Game Application.

- [10] M. Dos Santos Soares and J. Vrancken, "A framework for multi-layered requirements documentation and analysis," *Proc. - Int. Comput. Softw. Appl. Conf.*, pp. 308–313, 2011.
- [11] N. Power and T. Moynihan, "A theory of requirements documentation situated in practice," in *SIGDOC 2003: Finding Real-World Solutions for Doc.: How Theory Informs Pract. and Pract. Informs Theory. Proc. of the 21st Annu. Int. Conf. on Doc.*, 2003, pp. 86–92.
- [12] The Peak Performance Center, "How Children and Adults Learn," The Peak Performance Center, 2018. [Online]. Available: [www.thepeakperformancecenter.com/educational-learning/learning/principles-of-learning/adult-learning/children-adults-%0Alearn/](http://www.thepeakperformancecenter.com/educational-learning/learning/principles-of-learning/adult-learning/children-adults-%0Alearn/) The.
- [13] N. Maiden, "User requirements and system requirements," *IEEE Softw.*, vol. 25, no. 2, pp. 90–91, 2008.
- [14] J. Richardson, T. C. Ormerod, and A. Shepherd, "The role of task analysis in capturing requirements for interface design," *Interact. Comput.*, vol. 9, no. 4, pp. 367–384, 1998.
- [15] G. L. F. Fabbrini, M. Fusani, S. Gnesi, "Quality Evaluation of Software Requirement Specifications," 2001.
- [16] IEEE, *IEEE Standard for Software Requirement Specifications - IEEE Std 830-1998*, vol. 1998. 1998.
- [17] A. Hudaib, R. Masadeh, M. H. Qasem, and A. Alzaqebah, "Requirements Prioritization Techniques Comparison," *Mod. Appl. Sci.*, vol. 12, no. 2, p. 62, 2018.
- [18] O. De Troyer and E. Janssens, "Supporting the requirement analysis phase for the development of serious games for children," *Int. J. Child-Computer Interact.*, vol. 2, no. 2, pp. 76–84, 2014.
- [19] S. Cano, C. Collazos, H. M. Fardoun, and D. M. Alghazzawi, "Model Based on Learning Needs of Children," in *International Conference on Social Computing and Social Media*, 2016, vol. 3, pp. 324–334.
- [20] H. M. T. Tran and F. Anvari, "A Five-Dimensional Requirements Elicitation Framework for e-Learning Systems," *Int. J. Inf. Electron. Eng.*, vol. 6, no. 3, pp. 185–191, 2016.
- [21] A. G. Picciano, "Theories and frameworks for online education: Seeking an integrated model," *Online Learn. J.*, vol. 21, no. 3, pp. 166–190, 2017.
- [22] R. D. Vatavu, G. Cramariuc, and D. M. Schipor, "Touch interaction for children aged 3 to 6 years: Experimental findings and relationship to motor skills," *Int. J. Hum. Comput. Stud.*, vol. 74, pp. 54–76, 2015.
- [23] D. Charsky, "From edutainment to serious games: A change in the use of game characteristics," *Games Cult.*, vol. 5, no. 2, pp. 177–198, 2010.
- [24] M. P. Jacob Habgood and S. E. Ainsworth, "Motivating children to learn effectively: Exploring the value of intrinsic integration in educational games," *J. Learn. Sci.*, vol. 20, no. 2, pp. 169–206, 2011.
- [25] T. Rodgers, *Engineering Play: A Cultural History of Children's Software*, vol. 14, no. 7. Cambridge, Massachusetts: The MIT Press Cambridge, 2011.
- [26] K. Wiegers and J. Beatty, *Software Requirements*. Microsoft Press, 2013.
- [27] S. Al-Megren and A. Almutairi, "Analysis of User Requirements for A Mobile Augmented Reality Application to Support Literacy Development Amongst Hearing-Impaired Children," *J. Inf. Commun. Technol.*, vol. 18, no. 1, pp. 97–121, 2019.
- [28] M. E. Gredler, *Games and Simulations and Their Relationships to Learning*. 2004.
- [29] L. Sha, C. K. Looi, W. Chen, P. Seow, and L. H. Wong, "Recognizing and measuring self-regulated learning in a mobile learning environment," *Comput. Human Behav.*, vol. 28, no. 2, pp. 718–728, 2012.
- [30] A. Slimani, O. B. Yedri, F. Elouaai, and M. Bouhorma, "Towards a design approach for serious games," *Int. J. Knowl. Learn.*, vol. 11, no. 1, pp. 58–81, 2016.
- [31] P. Backlund and M. Hendrix, "Educational games - Are they worth the Effort?: A literature survey of the effectiveness of serious games," 2013 5th Int. Conf. Games Virtual Worlds Serious Appl. VS-GAMES 2013, 2013.
- [32] D. Parsons, H. Ryu, and M. Cranshaw, "A design requirements framework for mobile learning environments," *J. Comput.*, vol. 2, no. 4, pp. 1–8, 2007.
- [33] M. K. Sabariah, P. I. Santosa, and R. Ferdiana, "User experience analysis in software requirements specification (srs) of learning application for children," *Int. J. Pure Appl. Math.*, vol. 119, no. 15 Special Issue B, pp. 2983–2988, 2018.
- [34] A. N. Antle, "Exploring how children use their hands to think: An embodied interactional analysis," *Behav. Inf. Technol.*, vol. 32, no. 9, pp. 938–954, 2013.
- [35] P. Jamieson and L. Grace, "A framework to help analyze if creating a game to teach a learning objective is worth the work," *Proc. - Front. Educ. Conf. FIE*, vol. 2016-Novem, 2016.
- [36] S. Kujala, "Effective user involvement in product development by improving the analysis of user needs," *Behav. Inf. Technol.*, vol. 27, no. 6, pp. 457–473, 2008.
- [37] J. A. Khan, I. U. Rehman, Y. H. Khan, S. A. Shah, and W. Khan, "Enhancement in Agile Methodologies Using Requirement Engineering Practices.," *Sci. Int.*, vol. 28, no. 2, pp. 1525–1532, 2016.
- [38] H. Nang and A. Harfield, "A framework for evaluating tablet-based educational applications for primary school levels in Thailand," *Int. J. Interact. Mob. Technol.*, vol. 12, no. 5, pp. 126–139, 2018.

# Design and Implementation of Real Time Data Acquisition System using Reconfigurable SoC

(DAS using RSoC)

Dharmavaram Asha Devi<sup>1</sup>

Department of Electronics and  
Communication Engineering  
Sreenidhi Institute of Science and  
Technology, Hyderabad, India

Tirumala Satya Savithri<sup>2</sup>

Department of Electronics and  
Communication Engineering  
JNTUH College of Engineering,  
Kukatpalli, Hyderabad, India

Sai Sugun.L<sup>3</sup>

Project Assistant  
Sreenidhi Institute of Science and  
Technology  
Hyderabad, India

**Abstract**—System on chip (SoC) technology is widely used for high speed and efficient embedded systems in various computing applications. To perform this task, Application Specific IC (ASIC) based system on chips are generally used till now by spending maximum amount of research, development time and money. However, this is not a comfortable choice for low and medium-level capacity industries. The reason is, with ASIC or standard IC design implementation, it is very difficult task where quick time to market, upgradability and flexibility are required. Therefore, better solution to this problem is design with reconfigurable SoCs. Therefore, FPGAs can be replaced in the place of ASICs where we can have more flexible and reconfigurable platform than ASIC. In the embedded world, in many applications, accessing and controlling are the two important tasks. There are several ways of accessing the data and the corresponding data acquisition systems are available in the market. For defence, avionics, aerospace and automobile applications, high performance and accurate data acquisition systems are desirable. Therefore, an attempt is made in the proposed work, and it has been discussed that how a reconfigurable SoC based data acquisition system with high performance is designed and implemented. It is a semicustom design implemented with Zynq processing system IP, reconfigurable 7-series FPGA used as programmable logic, hygro, ambient light sensor and OLEDRgb peripheral module IPs. All these sensor and display peripheral modules are interfaced with processing unit via AXI-interconnect. The proposed system is a reconfigurable SoC meant for high-speed data acquisition system with an operating frequency of 100MHz. Such system is perfectly suitable for high speed and economic real time embedded systems.

**Keywords**—Application Specific Integrated Circuit (ASIC); Advanced eXtensible Interface (AXI); data acquisition system; Field Programmable Gate Array (FPGA); Peripheral Module (PMOD); System on Chip (SoC)

## I. INTRODUCTION

Many cases, the term ‘SoC’ was referred as an Application Specific Integrated Circuit (ASIC). The best example of an ASIC based application is mobile phone. Complex circuitry with multiple functions with high speed logic, interfacing of many peripherals including memory is implemented on single chip meant for specific application is known as application specific integrated circuit. The solution with SoC gives low

cost with bulk volume production and enables high speed data transfers between the various system blocks [1].

The data acquisition is an integral part of any measurement and control systems used in various applications. The principle of data acquisition is to acquire real world physical parameters such as temperature, light intensity, pressure, sound etc. through appropriate sensors that are connected through multiple channel data selectors using time division multiplexing with serial peripheral interface technique or parallel processing technique [2]. The real and physical data which is in analog nature has to be converted in to digital representation using digitization. The process of conversion of analog signal into digital signal is known as digitization. After digitization, the digital data is read by the processor and process the data in readable format and then it will be sent to the display devices if standalone measurement is required. Otherwise, the data can be accessed remotely by using web-based data acquisition techniques using WiFi network.

Development time, cost and lack of flexibility are the major drawbacks of ASIC based SoCs [1]. These are appropriate choice for bulk volume production and where there are no requirements for upcoming enhancements. Because of this reason, low or medium volume market industries depends on a convenient solution, system on a programmable IC, an exact essence of system on reconfigurable device [1], [2]. This can be done with FPGAs which are reconfigurable and flexible than ASIC based SoCs. It is a better solution that using an FPGA for applications where system enhancements are desired [3]. The technical details of peripheral modules, Hygro, ALS and OLEDRgb are referred from [4], [5] and [6], respectively.

The earlier researchers were implemented FPGA based data acquisition designs using ISE software and Spartan 3, Spartan 6 FPGAs respectively. Daniel Roggow et al explained a laboratory workstation configuration with ZedBoard. They have demonstrated workstation setups for MP-1: Quadcopter Interface, MP-2: Digital Camera, MP-3: Target Acquisition and MP-4: UAV Control and they have received a good feedback from the students [7].

It is to be noted in [8], described about the system with three modules named as signal processing, data acquisition



with FPGA and data storage. It was designed by VHDL and simulated with ISE.

FPGA based high speed ADC with a sampling rate of 80 Mega samples and used DMA without loss of data and DDR3 memory was used for manipulations of data was discussed in [9]. The design was implemented with VHDL.

DAQ with Network Control Module using FPGA was discussed in [11]. LabView software tool was used to design and development of their design and tested with National Instruments data acquisition and FPGA devices.

In [12], a 32 channel DAQ system for medical imaging and clinical application was presented. It was developed with FPGA, NI's PXI, ADC, signal generator, timing and synchronization modules. FFT-hardware was implemented in FPGA for the purpose of high frame rates, demodulation of the signal and higher order harmonics spectral characteristics.

In [13], Satellite tracking data with respect to the humidity, temperature and light data measurements with Zynq processor based reconfigurable SoC discussed efficiently.

Smart monitoring of automobile data logger design and prototype implementation with Zed Board and Xilinx platform was explained in [14].

The FPGA based processor designs, verifications of various applications are referred from [15], [16], [17], [18], [19] and [20].

As per the present survey of history in related work, the majority of the research work methodology is conventional FPGA front end flow designs using CAD tools. Therefore, the corresponding technology and performance of the designs are restricted within the technology scope of hardware and software tools used. Definitely, there is a requirement in the performance-based enhancements in the SoC based designs. Therefore, by adding today's technology towards SoC designs, especially in the design of high-speed real time embedded system-based applications, an attempt is made to design and develop the high-performance based data acquisition system in Zynq-7000 architecture platform.

## II. DESIGN OF SOC BASED DATA ACQUISITION SYSTEM

### A. Methodology

The proposed system methodology is a semicustom SoC block design using multiple IP integration using front end CAD tool, Xilinx Vivado System Design Suite and SDK software. The proposed system architecture consists of Zynq processing system interfaced with hygro, ALS and OLEDrgb via AXI interconnect. The hardware part of the design is loaded in to the Artix7 FPGA, which acts as a programmable logic device. The data processing and controlling part is implemented with application software with SDK and ARM9 processor, acts as processing system. The hardware used for verification are Artix 7 FPGA and the multiple sensors interfaced with FPGA. The results can be monitored by either standalone system or remote system through WiFi network. In the proposed system, peripheral modules, Hygro and ALS are used to access temperature, humidity and ambient light intensity. These parameters are displayed on OLEDrgb display

and hyper terminal. The complete procedural flow with respect to CAD tool is illustrated in Fig. 1.

### B. Design Procedure

The proposed research work is based on FPGA based SoC using Zynq processor. It provides a perfect stage for the implementation of flexible system on chips. Because, Zynq is a sandwich of a processing-system (PS) and programmable-logic (PL) [1], [3]. The PS is prepared with a dual-core ARM processor (ARM Cortex-A9). And the Artix 7 FPGA is used as programmable logic.

In the proposed work, Zed board is used that consists of Zynq architecture, integrated memory, a various number of peripherals, general purpose ports and high-speed interfacing ports for communication. The Programmable Logic is used for design and implementation of logic with high speed. And the processing system cares about software routines, operating systems, and provides the communication between software and hardware [2]. To meet this need, Xilinx high level synthesis tool, Artix 7 FPGA and ARM processor are used.

In the proposed architectural design of data acquisition system, AXI peripheral interconnect block is used to interface Zynq processing system with all peripheral interfaces as shown in Fig. 2.



Fig. 1. Design and Implementation flow of SoC based DAS.

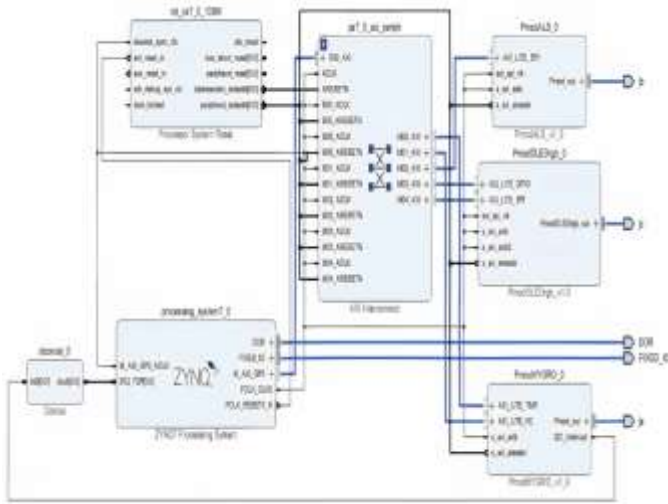


Fig. 2. Block Diagram of Data Acquisition System.

In this design, three peripheral modules are used, that are Hygro, ALS and OLEDrgb. The Hygro PMOD consists of temperature and humidity sensors. Hence, in order to select these parameters, one interrupt signal is used. And this interrupt is connected through Concat block. We can observe from Fig. 2, an interrupt-IRQ (Interrupt Request) is connected to output of Concat block and the input of Concat is connected to I2C Interrupt from PMOD Hygro. This Hygro PMOD is connected to port Ja of the Zed board. The other PMOD, ALS is connected to port Jb. To display the acquired and processed data, OLEDrgb is used which is connected to port Jc. In order to reset the system, one reset block, 'Processor System Reset' is used.

In the proposed system, we can view these parameters, in the hyper terminal. Hence, in its Peripheral I/O pins, one UART peripheral is enabled.

After completion of the block level design, validate it and make sure there will be no errors in the design. Then, generate HDL wrapper for the design, after that synthesize and implementation processes has to be done. Once if implementation is to be done successfully, then bit stream can be generated.

Next process is exporting the hardware including bitstream and launch software design kit (SDK) software to create and build the application project. Application project in the proposed system is developed using C. Here, the zed board consists of programmable logic (PL) and processing system (PS) are sandwiched in single IC. The created design will be loaded in PL, which is an ARTIX 7 FPGA. The control process, that is accepting the input data, processes it and displays the appropriate results at the display devices will be done by software logic instructed to processing system.

### III. HARDWARE INTERFACE

Once, application project is ready, then make sure we have to interface PMOD Hygro to output connector A, ALS to connector B and OLEDrgb is connected to the output connector C as we have created our design in Fig. 2. The

peripheral modules, hardware connectivity with zed board are shown in Fig. 3.

As there are three PMODs used in the proposed system, as illustrated in Fig. 4, their specifications and interface configuration details are explained as follows.

#### A. HYGRO Interface

The HYGRO PMOD consists of T1 HDC1080, which is an integrated digital temperature and humidity sensors. It provides accurate measurement with low power. It operates in a range of 2.7 V to 5.5 V supply. It is more economical and suitable for wide range of low power applications. The temperature and humidity sensors are calibrated with  $\pm 0.2^{\circ}\text{C}$  (typical) and  $\pm 2\%$  accuracy.

Fig. 5 illustrates the interfacing of HYGRO PMOD with Zynq Processor. It has dual modes of operations known as measurement mode and sleep mode. After power up, it will be in sleep mode and waits for I2C input and commands. The commands are used to trigger and read measurements, check the status condition of the battery and configure the timing conversions. Whenever it obtains a command to trigger a measurement, it switches to measurement mode from sleep mode. If it completes the measurement, it will return to sleep mode. The default mode of this device is, first it will measure temperature and then humidity. The bit pattern of the 16-bit temperature register is, first it will hold 14-bit acquisition value and the two least significant bits are always zero.

The result accuracy depends on the time conversion. The temperature and relative humidity can be calculated using equations 1 and 2, respectively [3].

$$\text{Temperature } (^{\circ}\text{C}) = \frac{\text{Temperature}[15:00]}{2^{16}} * 165^{\circ}\text{C} - 40^{\circ}\text{C} \quad (1)$$

$$\text{Relative Humidity } (\% \text{RH}) = \frac{\text{Humidity}[15:00]}{2^{16}} * 100\% \text{RH} \quad (2)$$

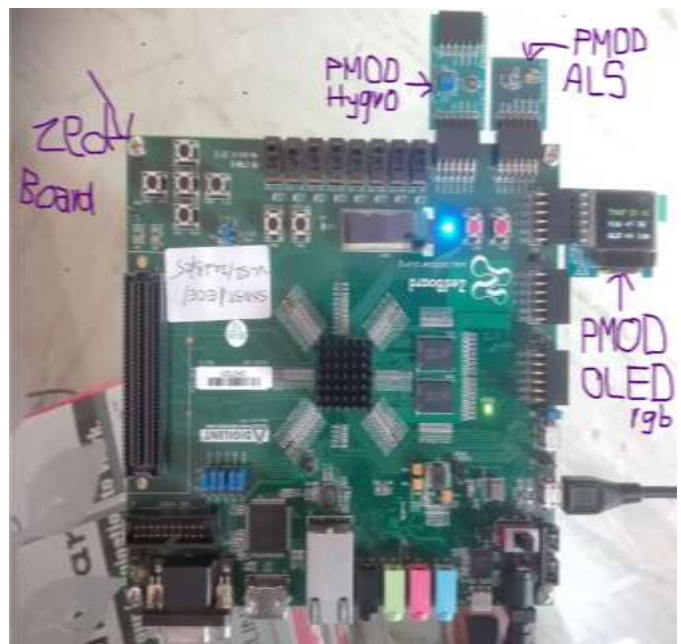


Fig. 3. Interfacing of HYGRO, ALS, OLEDrgb PMODs with PL via PS.

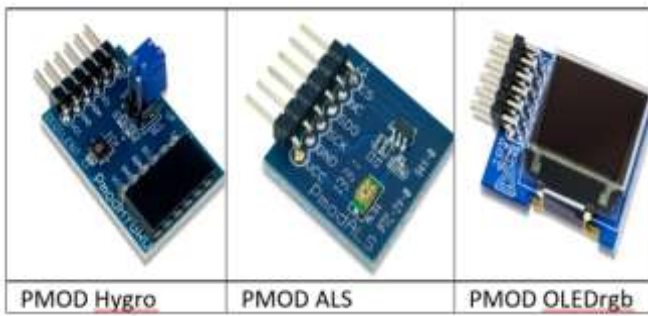


Fig. 4. PMODs used in Data Acquisition System.

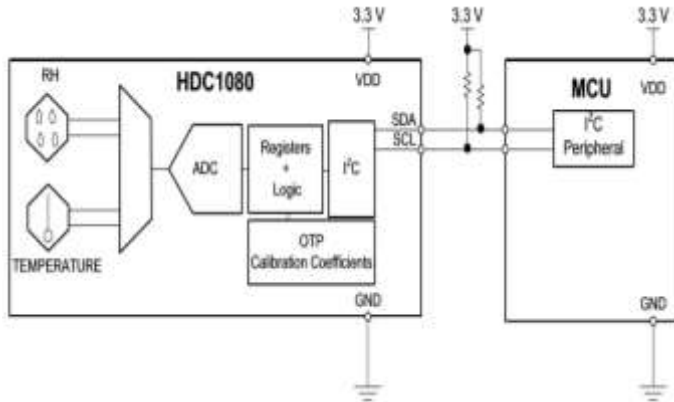


Fig. 5. Interfacing of HYGRO with Zynq Processor.

To complete the measurement of humidity and temperature, it is required to organize the register address to 0x02. We have to configure logic HIGH to Bit-12 to measure both humidity and temperature parameters. In order to set the resolution of a temperature, configure logic LOW to Bit-10 for 14-bit resolution or set it to logic HIGH for resolution of 11-bit.

In order to set the desired humidity measurement resolution, set the bits 9 and 8 of configuration register to 00 to achieve 14-bit resolution, or set to 01 to achieve 11-bit resolution or set to 10 to achieve 8-bit resolution. The measurements are triggered by setting an address pointer to 0x00 then, measurements waiting period will be completed, depending on the time conversion and read the output data [3]. Fig. 6 shows the timing signals under read operation when data is ready.

### B. ALS Interface

The peripheral module ALS is a single ambient light sensor. When the ALS is exposed with light, it will convert the light into voltage signal. This analog voltage signal is converted into 8-bits of digital data by the analog to digital converter. The range of values from 0 to 255 indicates low light level to a high light level [4].

The pin configuration of ALS [4], Pin.1 is Chip Select (CS), Pin.2 is no connection (NC), Pin.3 is Serial Data Out (SDO), Pin.4 is Serial Clock (SCK), Pin.5 is ground pin, and Pin.6 is Power Supply, VCC (3.3V/5V). The operating voltage range of this ALS is 2.7V to 5.25V. However, the proposed work requires 3.3V.

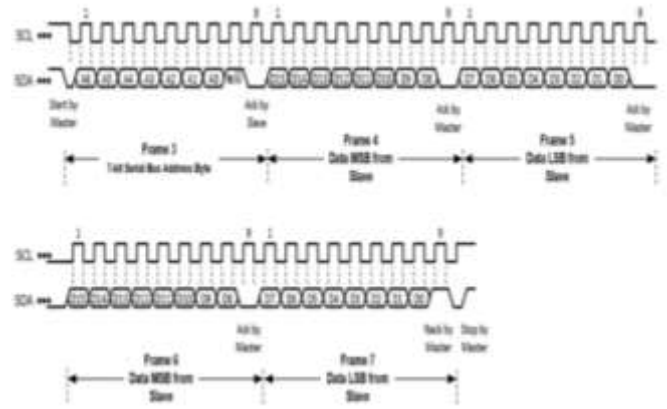


Fig. 6. Humidity and Temperature Measurement.

It communicates with the Zed board, via serial peripheral interface and requires a serial clock frequency (SCK) of 1 MHz and 4 MHz. When CS pin is made low, it operates in a regular mode of operation and brings a single reading in sixteen serial clock (SCLK) cycles [4]. The information bits are kept on the dropping edge of the serial clock. It is valid on the succeeding growing edge of SCLK. It consists of 3 zeroes at starting, the 8-bits of data with the most significant bit at first, and 4 zeroes at the end.

### C. OLEDrgb Interface

There are six ports are used between the peripheral module, OLEDrgb to Zynq processing system via AXI interconnect. AXI\_LITE\_GPIO is connected to M03\_AXI, AXI\_LITE\_SPI is connected to M04\_AXI. Because, it is connected through SPI protocol. The input clock to this module is, ext\_spi\_clk, driven by the processing system through FCLK\_clk0. s\_axi\_aclk of Hygro and s\_axi\_aclk of ALS are connected to s\_axi\_aclk and s\_axi\_aclk2 of OLEDrgb to have the communication appropriately.

The reset signal is generated from the processor reset block to all the peripheral modules including AXI interconnect as shown in the Fig. 2. The output port is connected to jC in order to interface the OLEDrgb to Zed Board as illustrated in Fig. 8.

## IV. DEVELOPMENT OF APPLICATION SOFTWARE

The process, export the created project into the software development kit using the Vivado software, will create “Hardware Base system” or “Hardware Platform”. Once, the hardware platform is created, software system is used or developed to complete the real time application. This software system is in the form of three layers over the hardware system base [10] as illustrated in Fig. 7. The second layer is board support package layer, which will set functions and drivers of low level that are needed by following layer over the board support package layer to communicate with hardware used. Application software is created with C or C++ and it will run over the operating system and it is the highest level of abstraction from the bottom hardware [1], [10].

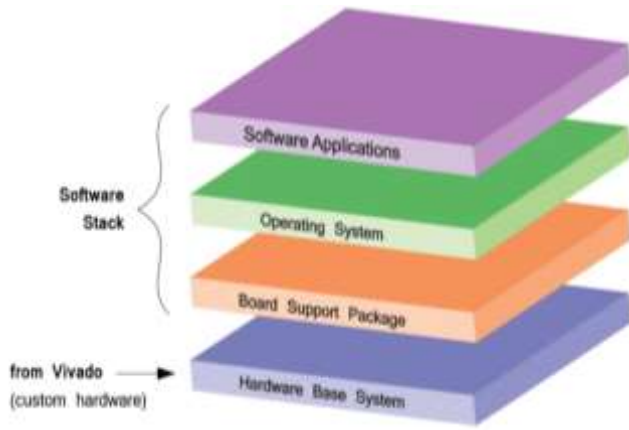


Fig. 7. Software Layers Over the Custom Hardware [10].

### V. RESULTS AND DISCUSSION

The Xilinx software development kit is used to create our data acquisition application. Compilation and debugging processes are also done within this tool. In the proposed application, we have connected three peripheral modules: Hygro, ALS and OLEDrgb as illustrated in Fig. 3. The application software is developed to read the temperature and humidity from Hygro Pmod and ambient light intensity from ALS Pmods respectively. These three data values are processed and displayed on OLEDrgb display Pmod as shown in Fig. 8.

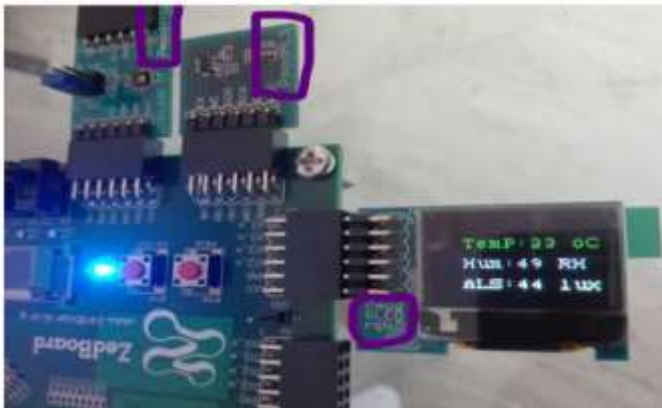


Fig. 8. Display of Temperature, Humidity and Ambient Light Intensity on OLEDrgb.

TABLE I. UTILIZATION REPORT

| Table I(A). LUT, Slice Registers, Muxes and Slices |                    |                          |                  |                  |               |
|--|--------------------|--------------------------|------------------|------------------|---------------|
| Name   | Slice LUTs (53200) | Slice Registers (106400) | F7 Muxes (26600) | F8 Muxes (13300) | Slice (13300) |
| DAC_i  | 2007               | 2530                     | 8                | 4                | 473           |
| PmodALS_0  | 387                | 590                      | 0                | 0                | 166           |
| PmodHYGRO_0  | 580                | 541                      | 8                | 4                | 211           |
| PmodOLEDrgb_0                                      | 431                | 666                      | 0                | 0                | 174           |
| PS7_Axi_peripheral                                 | 593                | 696                      | 0                | 0                | 242           |
| Rst_ps7_0_100M                                     | 18                 | 37                       | 0                | 0                | 13            |

TABLE I. (B). LUT AS LOGIC, MEMORY AND FLIPFLOP PAIRS

| Name               | LUT as Logic (53200) | LUT as Memory (17400) | LUT Flip Flop Pairs (53200) |
|--------------------|----------------------|-----------------------|-----------------------------|
| DAC_i              | 1901                 | 106                   | 1130                        |
| PmodALS_0          | 370                  | 17                    | 265                         |
| PmodHYGRO_0        | 570                  | 10                    | 273                         |
| PmodOLEDrgb_0      | 414                  | 17                    | 300                         |
| PS7_Axi_peripheral | 532                  | 61                    | 272                         |
| Rst_ps7_0_100M     | 17                   | 1                     | 15                          |

TABLE I. (C). IOBS AND LOGIC UTILIZATION

| Name               | Bonded IOB (200) | Bonded IOPADS (130) | ILOGIC (200) | OLOGIC (200) |
|--------------------|------------------|---------------------|--------------|--------------|
| DAC_i              | 1901             | 106                 | 1130         | 2            |
| PmodALS_0          | 370              | 17                  | 265          | 1            |
| PmodHYGRO_0        | 570              | 10                  | 273          | 0            |
| PmodOLEDrgb_0      | 414              | 17                  | 300          | 1            |
| PS7_Axi_peripheral | 532              | 61                  | 272          | 0            |
| Rst_ps7_0_100M     | 17               | 1                   | 15           | 0            |

The area of the proposed design occupied in the Artix 7 FPGA is represented in Table I. The LUT area for total DAC is 3.77%. Slice registers occupy 2.38%, total Slices occupy 5.59%. The detailed utilization report for LUT, Slice Registers, Muxes and Slices are illustrated in Table I(A). LUT as Logic, Memory and FlipFlop Pairs reports are specified in Table I(B). The number of input and output blocks and input logic and output logics are illustrated in Table I(C).

The proposed design, timing constraints are met satisfactorily as specified in the Table II. All the worst negative slacks for setup, hold and pulse widths are positive. There are zero negative slacks for setup, hold and pulse widths. Numbers of Failing end points are also zero. There were total 5636 end points for each setup and hold. And for pulse width, total end points are 2687.

The comparison of proposed system that is Reconfigurable data acquisition with respect to ASIC based data acquisition systems are illustrated in Table III and Fig. 9. Therefore, the summarized advantages of Reconfigurable data acquisition systems are listed as follows.

1) The existed data acquisition systems are implemented with ASIC (Application Specific Integrated Circuit) based, which does not include the enhancement facility in increasing the number of channels and the corresponding data acquisition and signal conditioning circuitry enhancement. The non-recurring engineering cost and time to market is very high. Whereas the proposed system is implemented with Reconfigurable device, and the corresponding cost and development time is very less. Hence, it is best suitable for low and medium industrial applications and even academic institutions can develop such systems.

2) The software and hardware used in the proposed system is suitable for low powered and high-speed applications which can meet with ASIC based designs.

3) Many embedded systems used LCD, LED or monitors for data acquisition purpose. Whereas in the proposed system, OLEDrgb is used for display of results which is a tiny, clarity colour display of 94 X 64 pixels [5].

TABLE II. DESIGN TIMING SUMMARY REPORT

| Table II. Timing Report             |                                 |  |
|-------------------------------------|---------------------------------|--|
| Setup                               | Hold                            | Pulse Width                                    |
| Worst Negative Slack (WNS): 2.417ns | Worst Hold Slack (WHS): 0.026ns | Worst Pulse Width Slack (WPWS): 3.750ns        |
| Total Negative Slack(TNS): 0.00ns   | Total Slack(TNS): 0.00ns        | Worst Pulse Width Negative Slack (WPWS):0.00ns |
| Number of Failing Endpoints: 0      | Number of Failing Endpoints: 0  | Number of Failing Endpoints: 0                 |
| Total Number of Endpoints: 5636     | Total Number of Endpoints: 5636 | Total Number of Endpoints: 2687                |

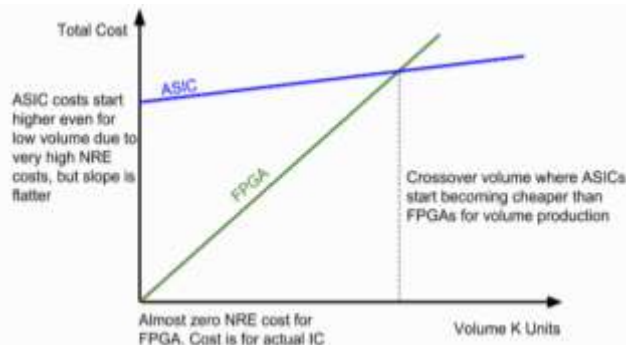


Fig. 9. ASIC Versus FPGA based Data Acquisition Systems.

TABLE III. COMPARISON ANALYSIS OF RECONFIGURABLE DATA ACQUISITION SYSTEM VS ASIC BASED DATA ACQUISITION SYSTEM

| S.No. | ASIC based Data Acquisition SoC  | Reconfigurable Data Acquisition SoC  |
|-------|--|--|
| 1     | Once the circuit design is taped out, alteration is not possible. Hence, it is less flexible for enhancement of the designs  | In Reconfigurable data acquisition, FPGAs are used. With these, either partial or full reconfiguration is possible.                    |
| 2     | Not suitable for applications where design up gradation is very much required such as Radars, mobile applications etc.   | These are extremely suitable for Radars, mobile base stations, high computing data applications etc.                                   |
| 3     | These are not suitable for prototype of a design. ASICs are also prototype by FPGAs.   | Best suitable for prototype and validation of a design.  |
| 4     | Difficult for low and medium level industries because of unbearable NRE cost. Hence, it is a hard barrier to design entry and lot of time is required to market the end product. | Many advantages such as less NRE cost, easy barrier to design entry, less time required to market.                                     |
| 5     | No doubt these are energy efficient, high performance and cost per unit will be less with bulk volume production.  | With latest 7-series FPGAs, these are also more or less equivalent to ASIC based designs in terms of performance and power consumption |

## VI. CONCLUSION

The Reconfigurable SoC based data acquisition system is designed, implemented and made successfully functional using Xilinx Vivado System Design Suite 2018.1, Xilinx Software Development Kit (SDK), Zed development board, and three peripheral modules. Advantages of this proposed system has been clearly discussed with the comparison analysis of Reconfigurable data acquisition with respect to ASIC based data acquisition system. It is a cost effective, low powered and high accurate Reconfigurable embedded application.

## ACKNOWLEDGMENT

This work has been carried out by using the equipment sanctioned by TEQIP Phase-III sponsored fund from the JNTUH. Therefore, the authors of this paper acknowledge the funding source, JNTUH, Kukatpalli and the supporting institute, Sreenidhi Institute of Science and Technology for giving the encouragement and providing facilities to do research and development activities.

## REFERENCES

- [1] Louise H. Crockett, Rose A. Elloit, Martin A. Enderwitz, David Northcote "The Zynq Book Tutorials for Zybo and Zed Board", Department of Electronic and Electrical Engineering, University of Strathclyde Glasgow, Scotland, UK, August 2015.
- [2] Survey on developing data acquisition system using ZYNQ architecture, ZedBoard\_HW\_UG\_v1\_9.pdf, ZedBoard (Zynq™ Evaluation and Development) Hardware's Users Guide, Version 1.9 January 29th, 2013.
- [3] Pmod-hygro-rm.pdf, PMOD HYGRO Reference Manual, Revised February 22, 2017.
- [4] pmodals\_rm.pdf., PmodALSTM, Reference Manual, Revised April 15, 2016.
- [5] Pmodoledrgb\_rm.pdf, PmodOLEDrgb TM, Reference Manual, Revised April 26, 2016.
- [6] Daniel Roggow, Paul Uhing, Phillip Jones and Joseph Zambreno, "A project-based embedded systems design course using a reconfigurable soc platform", 2015 IEEE International Conference on Microelectronics Systems Education (MSE) , 10.1109/MSE.2015.7160005, ISBN: 978-1-4799-9915-6, 2015.
- [7] AXI Reference Guide, UG761(v13.4) January 18, 2012.
- [8] Ye Fan, "FPGA based data acquisition system", 2011 IEEE International Conference on Signal Processing, Communications and Computing (ICSPCC), DOI:10.1109/ICSPCC.2011.6061756, Electronic ISBN: 978-1-4577-0894-7, October 2011.
- [9] Aboli Audumbar Khedkar, Khade. R.H. "High speed FPGA-based data acquisition system", Microprocessors and Microsystems, Volume 49, Pages 87-94, March 2017.
- [10] Louse H.Crockett, Ross A.Elliot, Martin A. Enderwitz, Robert W.Stewart "The Zynq Book Embedded Processing with the ARM cortex-A9 on the Xilinx, Zynq-7000 All Programmable SoC" 1st Edition, Published by Strathclyde Academic Media in July 2014.
- [11] Rajasekaran, C., Jeyabharath, R. and Veena, P. "FPGA SoC Based Multichannel Data Acquisition System with Network Control Module", Circuits and Systems, 8, 53-75, ISSN Online: 2153-1293, 2017.
- [12] S Khan et al 2013 FPGA Based High Speed Data Acquisition System for Electrical Impedance Tomography, J. Phys.: Conf. Ser. 434 012081.
- [13] Dharmavaram Asha Devi, Tirumala Satya Savithri, Sai Sugun.L, "Satellite Tracking Data Acquisition System using Reconfigurable SoC", Journal of Xi'an University of Architecture and Technology", Volume XII, Issue III, March 2020.
- [14] M. Bhavani and D. A. Devi, "Design of smart Monitor for automobiles using FPGA based Data Logger," 2019 International Conference on Communication and Electronics Systems (ICCES), Coimbatore, India, 2019, pp. 1940-1945, doi: 10.1109/ICCES45898.2019.9002034.

- [15] Dharmavaram Asha Devi, Muchukota Suresh Babu, "Design and Analysis of Power Efficient 64-Bit ALCCU", *International Journal of Recent Technology and Engineering (IJRTE)*, ISSN: 2277-3878, Volume-8 Issue-2, pp-162-166, July 2019.
- [16] Talal Bonny. "Chaotic or Hyper-chaotic Oscillator? Numerical solution, circuit design, Matlab HDL-coder implementation, VHDL code, security analysis, and FPGA realization", *Journal of Circuits, Systems, and Signal Processing*.
- [17] Talal Bonny and Qasim Nasir, "Clock Glitch Fault Injection Attack on an FPGA-based Non-autonomous Chaotic Oscillator", *Journal of Nonlinear Dynamics*.
- [18] Talal Bonny, Tamer Rabie, Mohammed Baziad, and Walid Balid. "SHORT: Segmented Histogram Technique for Robust Real-Time Object Recognition", *Journal of Multimedia Tools and Applications*.
- [19] D. A. Devi and N. S. Rani, "Design and Implementation of custom IP for Real Time Clock on Reconfigurable Device," 2019 Third International Conference on Inventive Systems and Control (ICISC), Coimbatore, India, 2019, pp. 414-418, doi: 10.1109/ICISC44355.2019.9036428.
- [20] D. A. Devi and L. S. Sugun, "Design, implementation and verification of 32-Bit ALU with VIO," 2018 2nd International Conference on Inventive Systems and Control (ICISC), Coimbatore, 2018, pp. 495-499, doi: 10.1109/ICISC.2018.8399122.

# GAIT based Behavioral Authentication using Hybrid Swarm based Feed Forward Neural Network

Gogineni Krishna Chaitanya<sup>1\*</sup>, Krovi Raja Sekhar<sup>2</sup>

Department of Computer Science and Engineering  
Koneru Lakshmaiah Education Foundation  
Vaddeswaram, 522502, Andhra Pradesh  
India

**Abstract**—Authentication of appropriate users for accessing the liable gadgets exists as one among the prime theme in security models. Illegal access of gadgets such as smart phones, laptops comes with an uninvited consequences, such as data theft, privacy breakage and a lot more. Straight forward approaches like pattern based security, password and pin based security are quite expensive in terms of memory where the user has to keep remembering the passwords and in case of any security issue risen then the password has to be changed and once again keep remembering the recent one. To avoid these issues, in this paper an effective GAIT based model is proposed with the hybridization of Artificial Neural Network model namely Feedforward Neural Network Model with Swarm based algorithm namely Krill Herd optimization algorithm (KH). The task of KH is to optimize the weight factor of FNN which leads to the convergence of optimal solution at the end of the run. The proposed model is examined with 6 different performance measures and compared with four different existing classification model. The performance analysis shows the significance of proposed model when compared with the existing algorithms.

**Keywords**—GIAT behavioral pattern recognition; feedforward neural network; krill herd algorithm

## I. INTRODUCTION

For the verification of identity one of the most trustworthy and effective approach is Biometric method [1]. Its advancements in the recent years are more significant than any other models due to its uniqueness in recognition [2]. One among the significance of biometric is, it has the tendency to connect with the authenticated user in a straightforward manner rather than linking the originality through a third-party mediator such as keywords or tokens [3]. From the other perspective of biometric, there exists GAIT recognition [4] in which the person can be identified by their gesture such as walk, running, etc. [5]. The uniqueness of a person gesture is often carried out in a more effective manner than when it compared with recognition through keywords or patterns. The verification of GAIT is carried out by any of the three forms through video mode, wearable and floor sensor [6]. The authentication of user is the first step for reducing the access of illegal persons [7]. The process is done through the verification of the given identity by the user. The other way of recognizing the authenticated user apart from biometric is through knowledge-based verification [8]. The knowledge-based verification is the process of information which is given by the authenticated user to verify the genuine of the person.

Example methods of knowledge-based verification are pattern, pin, passwords, etc. [9]. The disadvantage of this knowledge-based system is forgotten schema and stolen. To overcome this disadvantage an additional feature is added in mobile phones which is the recognition of users through fingerprint biometric [10]. This helps the equipment to verify the users which further gives another step of verification [11].

The recognition of human movements and their biometrics has reached a greater level of importance in the sectors such as military, airport, banks [12]. In some of the regions, authentication of users through a password are pin numbers are often leading to tedious process or it is sometimes recognized as less defense particularly is some of the assaults listed in [13]. Sometimes the catchy words of passwords are often easily predicted. On the other side if non-dictionary words are kept as passwords there is a high possibility of easy forgotten scenario [14]. One more way of authentication of authenticated users is proposed namely Speech recognition. However, the background noise is often making the system to be confused in recognizing the authenticate user. Also, it has the probability to be hacked [15]. Even in biometric recognition pattern also if there exists face-based recognition and if the face is unclear it also may lead to difficulty [16].

Sometimes after the usage of authenticated person there is a possibility of leaving the gadget without any lock and hence the probability of accessing it by an unauthenticated user is also high till the gadget gets locked [17]. Hence GAIT type of recognition is proposed to address their problems in smart phones. In this model, the embedded sensor with the device [18] observes the movement of user. If it matches with the authenticated user's movement then the phone will be unlocked. Through this model the user need not any other secondary verification activity to access the gadget. In this model, the unauthenticated user cannot access it since the sensor completely observes the user's motion and if it does not match it will be locked [19]. An optimal FNN network is used to classify the authenticated user and unauthenticated user [20]. The proposed model holds Krill Herd algorithm for optimizing the weights of FNN.

Thus, the rest of the paper is organized as follows: Section 2 deals with the literature study. Section 3 discusses the problem statement. Section 4 discusses the FNN and Krill herd algorithm. Section 5 holds the experimental evaluation of the proposed model and the last section concludes the paper.

\*Corresponding Author

## II. RELATED WORKS

The entry check point of a gadget often not secured enough, and it is open for data theft by any way means. In the year 2019, author Praveen kumar Rayani proposed a model [21] using Naïve Bayes classifier for effective verification of authenticated users using GAIT behavioral pattern recognition. The input for the proposed model is Boolean based banner model. Using this model, the authorization of intended data thief is being identified which improves the security of data in mobile devices. Sometimes the smart locks in the mobile device's issues certain kinds of problems such as holding on something for quite some time for being able to approve the authenticity. To address this issue Author Kazhuki, et al. [22] proposed a model that recognizes the walking gesture with the help of sensors in the mobile devices. Another model using key based arrangement model [23] which was proposed by Arne Bruschi, et al. to reduce the imperfection in recognition of GAIT based behavior. To prove the significance of the proposed model an effective power based attacking mechanism is used to test the integrity of the model. The output shows the consistency and the integrity of the proposed model by recognizing the GAIT based behavior in mobile device.

Answering based author authentication was proposed by the author Buriroa [24] which is used to record thousands of GAIT movement to propose an effective plan using behavioral pattern in GAIT. The other models include [25] waving of hand gesture etc. in this paper the contributions are listed as follows:

- An effective FNN hybridized with Swarm based Krill Herd Algorithm is developed for effective classification of authenticated and unauthenticated users [29].
- For the background subtraction from videos Kernel Compactness approximation is used to improve the quality of the frames.
- The proposed model is implemented and evaluated under suitable testbed.
- The performance of the proposed model is examined and proved its significance by comparing with the existing models.

## III. PROBLEM STATEMENT

One among the popular electric device is smartphone. Since it received a significant number of users worldwide, the problems such as hacking, phishing are also become common in nature for smart devices. One such mode of authentication is password or pattern lock or pin number. As it is discussed in the introduction section, these become void due to the restrictions it holds. To address this issue a mechanism called GAIT behavioral authentication system is proposed. Fig. 1 shows the model of normal authentication models such as key, face and pin-based authentication process.

In this plan, if the secret word space of the subsequent stage is distinguishable or arranged to parody, at that point it experiences in interior assault. This plan sets aside some additional effort for verification. To defeat interior assaults and data robbery, a potential arrangement is recognized through social biometrics of cell phone since the personal conduct standards of the individual client are indistinguishable.

### A. Information Gathering

The information for the recognition of authentication of users is collected using the sensor namely Gyroscope along with the accelerometer of the mobile devices and it collects the information such as the users style and gestures of running, sitting, walking and standing. The collected information are then sent to the model in the mobile device. The model is developed with the identification module of authenticated and unauthenticated users. In case if any illegal access is done on the gadget then the model identifies the irregularity and denote it to the user's knowledge.

### B. Pre-Processing

The preprocessing phase will remove the unwanted low-quality pixels for conducting the experiment in most effective way. The performance of the model can be well evaluated when the input is with less error values. It is removed using the pixilation model for reducing the unwanted pixels in the frames of the video.

### C. Feature Reduction using Kernel Compactness Approximation

The behavior recognition with the help of sensors lead to recognizing the actions of the users in 3D form. The input should be read in a 3-dimensional proforma so that it is effective to read the entire walking sitting or running model of the user [30]. The mode we used to generate all the key pairs that are useful for generating the Gram metric where the complexity to solve is in quadratic form.



Fig. 1. Problem Statement Features.



The kernel module of a frame can be shown in the form of  $k$  and it is represented in the following as:

$$k(m, n) = \int \beta_p \theta_p(m) \theta_p^*(n)$$

where  $\beta_p$  is the value of eigen and  $\theta_p$  is the eigen vector in normalized form.

$$T_r f = \int k(m) f(m) d\mu(m)$$

where  $f(m)$  denotes the features that are selected for user authentication purpose, the above equation is used for verification process.

#### IV. WORKING PRINCIPLE MULTI-LAYER PERCEPTRON IN FNN

The Feedforward Neural Network model is shown the Fig. 2. The number of input layer in any FNN will be 1. There may be more than one number of input variables in the input player. In Fig. 1, the input value ranges from 1 to  $n$ . The next layer is the hidden layer. Unlike input layer the hidden layer may range between any number of positive integer values and the number of neurons in the hidden layer is also not restricted. In general, the total number of hidden layers are between the range 3 and 6. The total number of output neurons should be at least one.

Each layer in FNN have certain input and output. In the input layer the input values depend on the problem. This will generate the input for the hidden layer and the output of the hidden layer will be calculated using Eq. (1).

$$s_j = \sum_{i=1}^n (W_{i,j} \times X_i) - \theta_j \quad (1)$$

And the output of the hidden layer will be gone through Equation (2) which is given below:

$$S_j = \frac{1}{(1 + \exp(-s_j))} \quad (2)$$

The variable  $j$  denotes the number of neurons in the hidden layer. The output computation is done using Eq. (3).

$$o_k = \sum_{j=1}^h (W_{j,k} \times S_j) - \theta_k \quad (3)$$

And finally, the output value is given to the activation function Eq. (4) and it is given as.

$$O_k = \frac{1}{(1 + \exp(-o_k))} \quad (4)$$

In the above model the weights ( $W$ ) and the bias values ( $\theta$ ) are random in general which ranges between 0 and 1. This often leads to non-optimal model construction. This can be eradicated by proposing the algorithm for handling the values such as weights and bias.

#### A. Krill Herd Algorithm

Krill herd algorithm has the potential to search for an optimal solution in the given stamp of time. It is inspired from the Krills concept. There are three processes which holds the search process in krills.

- 1) Movements imposed by other krills
- 2) Foraging
- 3) Random diffusion

And the combined solution can be represented as

$$\frac{dX_i}{dt} = N_i + F_i + D_i \quad (5)$$

The calculation of  $N, F$  and  $D$  are computed as follows:

The induction of movement by other solutions (Krills) can be computed as.

$$N_i^{new} = N_i^{max} \alpha_i + \omega_n N_i^{old} \quad (6)$$

And the value of  $\alpha$  is computed as

$$\alpha_i = \alpha_i^{local} + \alpha_i^{target} \quad (7)$$

Every Krill has its own foraging behavior which can be computed as follows:

$$F_i = V_f \beta_i + \omega_f F_i^{old} \quad (8)$$

The computation of  $\beta_i$  is given as follows:

$$\beta_i = \beta_i^{food} + \beta_i^{best} \quad (9)$$

The generation of new solutions can be done using Diffusion factor as follows:

$$D_i = D_i^{max} \left(1 - \frac{t}{t_{max}}\right) \delta \quad (10)$$

The Pseudo code of the Krill Herd algorithm is given in Algorithm 1 and the flowchart is shown in Fig. 3.

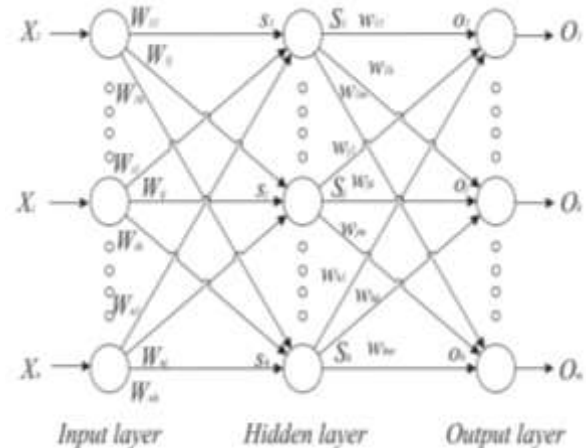


Fig. 2. Feedforward Neural Network.

### Krill Herd Algorithm for Tuning weight parameters

**Begin**

**Step 1:** Initialize  $Iter \leftarrow 1, i \leftarrow 0, k \leftarrow 0$

**Step 2:** for each  $KHIndiv \in KHSizedo$

$$Pop_{KHIndiv} \leftarrow LB + (UB - LB) * rand()$$

**end for**

**Step 3:** for each  $KHIndiv \in KHSizedo$

$$Fit(Pop_{KHIndiv}) \leftarrow f(Pop_{KHIndiv})$$

**end for**

**Step 4:** Repeat through Step 8 **Until**  $Max_{IT} \leq Iter$ , then go to Step 9

**Step 5:** Movement Induced by other Krill's

**Step 5.1:** Repeat through Step 5.3 **Until**  $i < KHSize$  |  $i \in KHIndiv$  else goto Step 6

$$Step 5.2: \alpha_i = \alpha_i^{local} + \alpha_i^{target}$$

$$Step 5.3: N_i^{new} = N^{max} \alpha_i + \omega_n N_i^{old}$$

**Step 6:** Foraging motion of individual Krill's without Swarm colonial behavior

**Step 6.1:** Repeat through Step 5.3 **Until**  $i < KHSize$  |  $i \in KHIndiv$  else goto Step 7

$$Step 6.2: \beta_i = \beta_i^{food} + \beta_i^{best}$$

$$Step 6.3: F_i = V_f \beta_i + \omega_f F_i^{old}$$

**Step 7:** Individual movement of Krill Herd in a random manner

**Step 7.1:** Repeat through Step 7.2 **Until**  $i < KHSize$  |  $i \in KHIndiv$  else goto Step 8

$$Step 7.2: D_i = D^{max} \left(1 - \frac{i}{i_{max}}\right) \delta$$

**Step 8:** Repeat through Step 8.2 **Until**  $i < KHSize$  |  $i \in KHIndiv$  else goto Step 9

$$Step 8.1: \frac{dX_i}{dt} = N_i + F_i + D_i$$

$$Step 8.2: \Delta t = C_t \sum_{j=1}^{NV} (UB_j - LB_j)$$

$$Step 8.3: X_i(t + \Delta t) = X_i(t) + \Delta t \frac{dX_i}{dt}$$

**Step 9:** Return  $\min(Fit(Pop_{KHIndiv}))$

**End**

**Output:**  $\min(Fit(Pop_{KHIndiv}))$

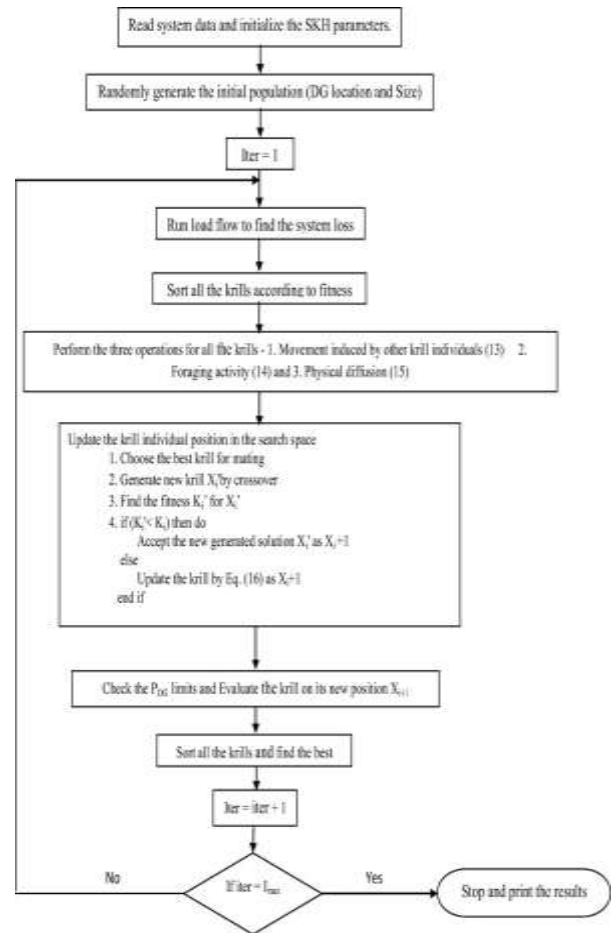


Fig. 3. Flowchart of Krill Herd Algorithm.

### V. EXPERIMENTAL EVALUATION

The proposed model is implemented in Python 3.7 with the system configuration Windows 10 Pro, with Intel core i7 processor speed 3.2GHz, with 16 GB primary storage and 1TB secondary storage. The proposed model identifies the user authentication with the help of GAIT based behavioral pattern recognition. To prove the significance of the proposed model, effective classification algorithms are used for evaluation on the same information which is gathered through sensors.

#### A. Case Study

The verification model of users is given in Fig. 4. The flow of identification between authenticated and unauthenticated is clearly given for better understanding of the proposed model.



Fig. 4. GAIT Verification Process.

Initially the authenticated behavior will be recorded based on the sensors to provide the training using legible users for the smart device. Later the unauthenticated behaviors are also recorded and given as input for training to find the illegal access. Once the training phase is over, the model will be ready for deployment of identifying the legal and illegal access and the users.

**B. Performance Measures**

The prove the performance of the proposed model 4 different classification algorithms are chosen namely Naïve Bayes classifier [21], Decision Forest-Decision jungle [26], Random forest [27] and SVM [28]. The different performance measures include Accuracy, Precision, Recall, F-Measure, False Accept Rate and False Reject Rate.

The confusion matrix useful for the interpretation of results is shown in Fig. 5.

1) *Accuracy*: The accuracy denotes the ratio between sum of true positive and true negative to the sum of all true positive, true negative, false positive and false negative value. Table I shows the overall percentage acquired by all algorithms along with the proposed algorithm.

$$Accuracy = \frac{True\ Positive + True\ Negative}{True\ Positive + True\ Negative + False\ Positive + False\ Negative}$$

From the Fig. 6, on comparing the results of accuracy with existing algorithms, our proposed model proves its significance in terms of percentage over NB with 1%, DF-DJ with 7%, RF with 2% and SVM with 1%.

2) *Precision*: Precision denotes the ratio between total number of identified correct answers with the total number of actual correct classification. The formula for calculation of precision from the confusion matrix is given in Table II and depicted as graph in Fig. 7.

$$Precision = \frac{True\ Positive}{True\ Positive + False\ Positive}$$

On comparing the results of accuracy with existing algorithms, our proposed model proves its significance in terms of percentage over NB with 4%, DF-DJ with 15%, RF with 9% and SVM with 10%.

3) *Recall*: The performance measure recall denotes the ratio between actual positive values to the overall identification of positive values from the dataset. The formula for calculation of recall from the confusion matrix is given in Table III and depicted as graph in Fig. 8.

$$Recall = \frac{True\ Positive}{True\ Positive + False\ Negative}$$

|                     |                               |                    |                    |
|---------------------|-------------------------------|--------------------|--------------------|
|                     |                               | True Condition     |                    |
|                     | Total Population              | Condition Positive | Condition Negative |
| Predicted Condition | Prediction Condition Positive | True Positive      | False Positive     |
|                     | Prediction Condition Negative | False Negative     | True Negative      |

Fig. 5. Confusion Matrix.

TABLE I. COMPARISON OF ACCURACY OF FNN-KL WITH EXISTING ALGORITHMS

| Algorithms                      | Accuracy (%) |
|---------------------------------|--------------|
| Naïve Bayes classifier          | 98.08        |
| Decision Forest-Decision jungle | 91.92        |
| Random forest                   | 97.16        |
| SVM                             | 97.11        |
| FNN-KL                          | 99.12        |

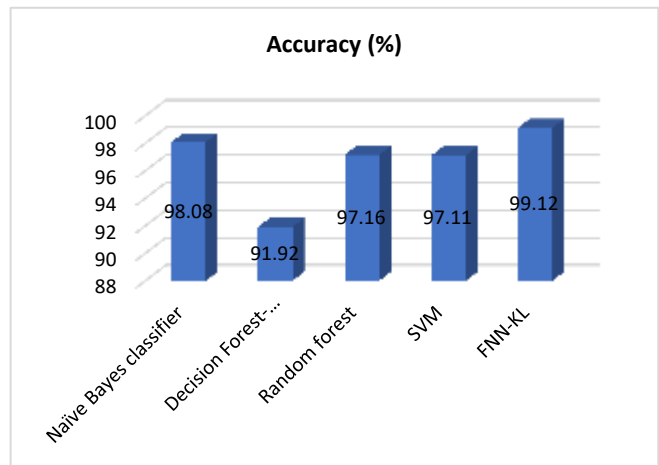


Fig. 6. Comparison Chart on Accuracy (%).

TABLE II. COMPARISON OF PRECISION OF FNN-KL WITH EXISTING ALGORITHMS

| Algorithms                      | Precision (%) |
|---------------------------------|---------------|
| Naïve Bayes classifier          | 95.53         |
| Decision Forest-Decision jungle | 84.44         |
| Random forest                   | 89.65         |
| SVM                             | 88.65         |
| FNN-KL                          | 99.58         |

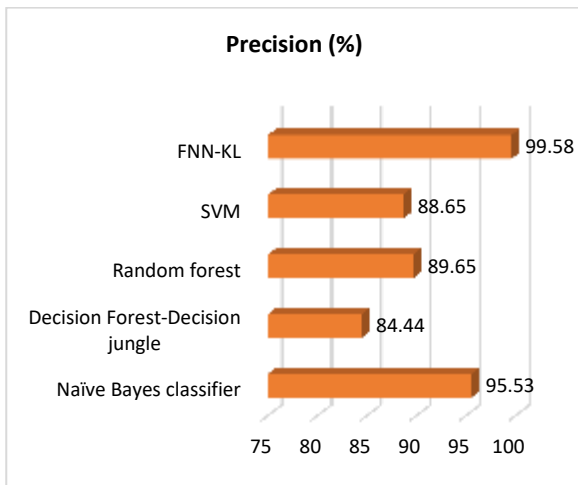


Fig. 7. Comparison Chart on Precision (%).

TABLE III. COMPARISON OF RECALL OF FNN-KL WITH EXISTING ALGORITHMS

| Algorithms                      | Recall (%) |
|---------------------------------|------------|
| Naïve Bayes classifier          | 86.11      |
| Decision Forest-Decision jungle | 90.34      |
| Random forest                   | 80.21      |
| SVM                             | 80.26      |
| FNN-KL                          | 95.47      |

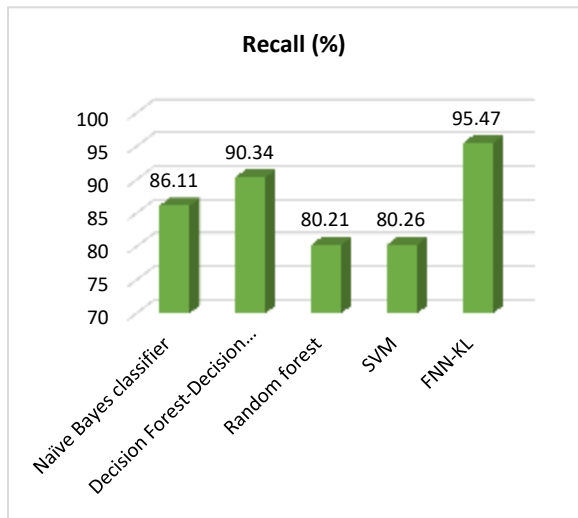


Fig. 8. Comparison Chart on Recall (%).

On comparing the results of accuracy with existing algorithms, our proposed model proves its significance in terms of percentage over NB with 9%, DF-DJ with 5%, RF with 15% and SVM with 16%.

4) *F-Measure*: The F-score is a meta measurement taken from precision and recall. The mathematical model of calculating F-Score is given in Table IV and depicted as graph in Fig. 9.

$$F - Score = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

TABLE IV. COMPARISON OF F-MEASURE OF FNN-KL WITH EXISTING ALGORITHMS

| Algorithms                      | F-measure (%) |
|---------------------------------|---------------|
| Naïve Bayes classifier          | 93.27         |
| Decision Forest-Decision jungle | 81.73         |
| Random forest                   | 92.50         |
| SVM                             | 94.48         |
| FNN-KL                          | 95.12         |

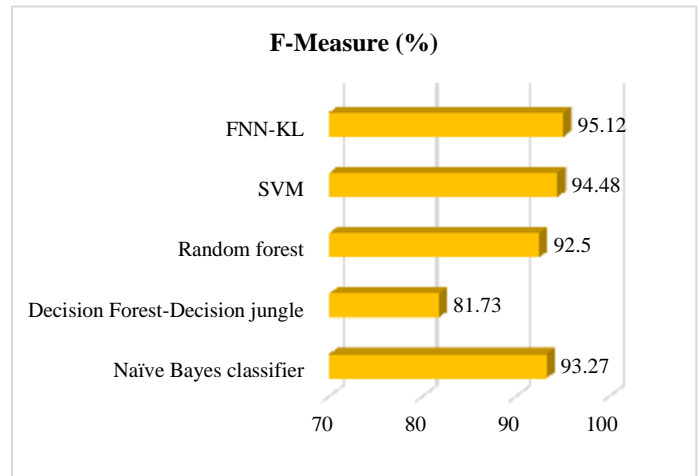


Fig. 9. Comparison Chart on F-Measure (%).

On comparing the results of accuracy with existing algorithms shown in Fig. 9, our proposed model proves its significance in terms of percentage over NB with 1%, DF-DJ with 1%, RF with 1% and SVM with 2%.

5) *False Accept Rate*: False accept rate for the given model is computed based on the number of unauthenticated users used the gadget. The values for the model are given in Table V. It shows that the proposed model significantly shows less error rate when compared with existing systems.

On comparing the results of accuracy with existing algorithms shown in Fig. 10, our proposed model proves its significance in terms of percentage over NB with 81%, DF-DJ with 76%, RF with 87% and SVM with 86%.

6) *False Reject Rate*: False reject rate is the identification of unauthenticated users to access the gadget and it lies as the ratio between the two models is given in Table VI.

TABLE V. COMPARISON OF FALSE ACCEPT RATE OF FNN-KL WITH EXISTING ALGORITHMS

| Algorithms                      | FAR (%) |
|---------------------------------|---------|
| Naïve Bayes classifier          | 3.93    |
| Decision Forest-Decision jungle | 1.33    |
| Random forest                   | 2.31    |
| SVM                             | 1.38    |
| FNN-KL                          | 1.35    |

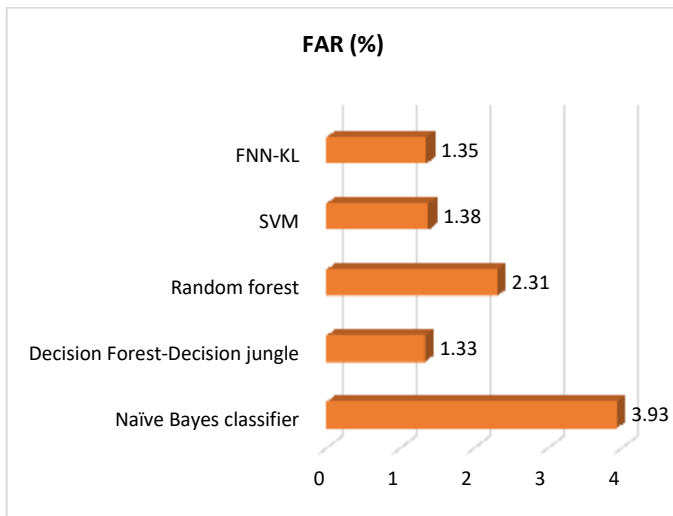


Fig. 10. Comparison Chart on FAR (%).

TABLE VI. COMPARISON OF FALSE REJECT RATE OF FNN-KL WITH EXISTING ALGORITHMS

| Algorithms                      | FRR (%) |
|---------------------------------|---------|
| Naïve Bayes classifier          | 3.26    |
| Decision Forest-Decision jungle | 2.22    |
| Random forest                   | 2.20    |
| SVM                             | 2.41    |
| FNN-KL                          | 6.85    |

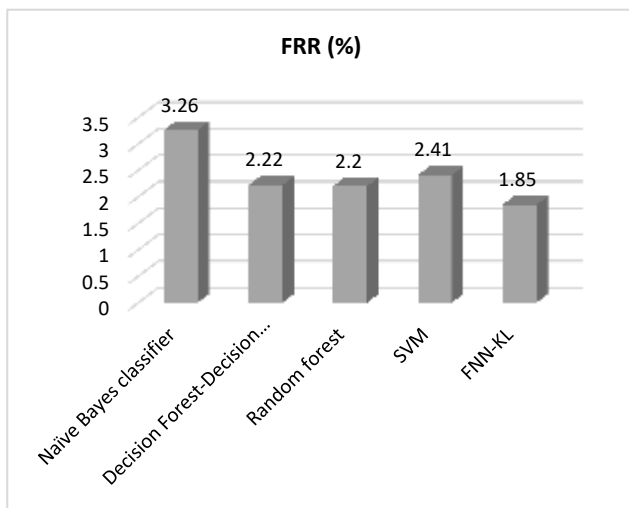


Fig. 11. Comparison Chart on FRR (%).

On comparing the results of accuracy with existing algorithms shown in Fig. 11, our proposed model proves its significance in terms of percentage over NB with 86%, DF-DJ with 79%, RF with 79% and SVM with 81%.

## VI. CONCLUSION

In this paper, an GIA based behavioral pattern recognition namely Feedforward Neural Network-Krill Herd (FNN-HL) is proposed for solving recognition of user for accessing the gadget. The proposed model comprises of FNN for

classification of authenticated and unauthenticated users and KH algorithm for tuning the weight and bias values in FNN. The performance of the proposed model is compared with four existing classification algorithms. The performance measures of the proposed model indicate the significance of FNN-KL when compared with other existing algorithms. The future work of this model can be extended with reducing the time complexity of processing the input frames.

## REFERENCES

- [1] Belkhouja, Taha, et al. "Biometric-based authentication scheme for Implantable Medical Devices during emergency situations." *Future Generation Computer Systems* 98 (2019): 109-119.
- [2] Mohsin, A. H., et al. "Based Blockchain-PSO-AES techniques in finger vein biometrics: A novel verification secure framework for patient authentication." *Computer Standards & Interfaces* 66 (2019): 103343.
- [3] Nazarkevych, Mariya, et al. "Biometric Identification System with Ateb-Gabor Filtering." 2019 XIth International Scientific and Practical Conference on Electronics and Information Technologies (ELIT). IEEE, 2019.
- [4] Khan, Muhammad Hassan, Muhammad Shahid Farid, and Marcin Grzegorzek. "A non-linear view transformations model for cross-view gait recognition." *Neurocomputing* (2020).
- [5] Huitzil, Ignacio, et al. "Gait recognition using fuzzy ontologies and Kinect sensor data." *International Journal of Approximate Reasoning* 113 (2019): 354-371.
- [6] Bai, Guifeng, and Yunqiang Sun. "Application and research of MEMS sensor in gait recognition algorithm." *Cluster Computing* 22.4 (2019): 9059-9067.
- [7] Spagnoletti, Paolo, et al. "Securing national e-ID infrastructures: Tor networks as a source of threats." *Organizing for the Digital World*. Springer, Cham, 2019. 105-119.
- [8] Zhou, Yuchen, et al. "Cyber-Physical-Social Systems: A State-of-the-Art Survey, Challenges and Opportunities." *IEEE Communications Surveys & Tutorials* (2019).
- [9] Joudaki, Zeinab, Julie Thorpe, and Miguel Vargas Martin. "Enhanced Tacit Secrets: System-assigned passwords you can't write down, but don't need to." *International Journal of Information Security* 18.2 (2019): 239-255.
- [10] Choudhary, Swati K., and Ameya K. Naik. "Multimodal Biometric Authentication with Secured Templates—A Review." 2019 3rd International Conference on Trends in Electronics and Informatics (ICOEI). IEEE, 2019.
- [11] Sharma, Rohit, Rajendra Prasad Mahapatra, and Naresh Sharma. "The Internet of Things and Its Applications in Cyber Security." *A Handbook of Internet of Things in Biomedical and Cyber Physical System*. Springer, Cham, 2020. 87-108.
- [12] Kakkad, Vishruti, Meshwa Patel, and Manan Shah. "Biometric authentication and image encryption for image security in cloud framework." *Multiscale and Multidisciplinary Modeling, Experiments and Design* 2.4 (2019): 233-248.
- [13] Liu, Yu, et al. "Account Lockouts: Characterizing and Preventing Account Denial-of-Service Attacks." *International Conference on Security and Privacy in Communication Systems*. Springer, Cham, 2019.
- [14] Karim, Nader Abdel, Zarina Shukur, and AbedElkarim M. AL-banna. "UIPA: User authentication method based on user interface preferences for account recovery process." *Journal of Information Security and Applications* 52 (2020): 102466.
- [15] Melnik, S. V., and N. I. Smirnov. "Voice Authentication System for Cloud Network." 2019 Systems of Signals Generating and Processing in the Field of on Board Communications. IEEE, 2019.
- [16] Mehraj, Tehseen, et al. "Critical Challenges in Access Management Schemes for Smartphones: An Appraisal." *Smart Network Inspired Paradigm and Approaches in IoT Applications*. Springer, Singapore, 2019. 87-113.
- [17] Ku, Yeeun, et al. "Draw It As Shown: Behavioral Pattern Lock for Mobile User Authentication." *IEEE Access* 7 (2019): 69363-69378.

- [18] Ahmadi, S. Sareh, Sherif Rashad, and Heba Elgazzar. "Machine Learning Models for Activity Recognition and Authentication of Smartphone Users." 2019 IEEE 10th Annual Ubiquitous Computing, Electronics & Mobile Communication Conference (UEMCON). IEEE, 2019.
- [19] Do, Quang, Ben Martini, and Kim-Kwang Raymond Choo. "The role of the adversary model in applied security research." *Computers & Security* 81 (2019): 156-181.
- [20] Fairley, Michael, David Scheinker, and Margaret L. Brandeau. "Improving the efficiency of the operating room environment with an optimization and machine learning model." *Health care management science* 22.4 (2019): 756-767.
- [21] Rayani, Praveen Kumar, and Suvamoy Changder. "Continuous Gait Authentication against Unauthorized Smartphone Access through Naïve Bayes Classifier." International Conference on Intelligent Computing and Communication. Springer, Singapore, 2019.
- [22] Watanabe, Kazuki, et al. "Gait-Based Authentication Using Anomaly Detection with Acceleration of Two Devices in Smart Lock." International Conference on Broadband and Wireless Computing, Communication and Applications. Springer, Cham, 2019.
- [23] Bruesch, Arne, et al. "Security Properties of Gait for Mobile Device Pairing." *IEEE Transactions on Mobile Computing* (2019).
- [24] Buriro, Attaullah, Bruno Crispo, and Mauro Conti. "AnswerAuth: A bimodal behavioral biometric-based user authentication scheme for smartphones." *Journal of information security and applications* 44 (2019): 89-103.
- [25] Shen, Chao, et al. "Waving Gesture Analysis for User Authentication in the Mobile Environment." *IEEE Network* 34.2 (2020): 57-63.
- [26] Kumar, Vivek, Chirag Gupta, and Vatsal Agarwal. "Gait-Based Authentication System." *Emerging Technologies in Data Mining and Information Security*. Springer, Singapore, 2019. 685-691.
- [27] Kececi, Aybuke, et al. "Implementation of machine learning algorithms for gait recognition." *Engineering Science and Technology, an International Journal* (2020).
- [28] Lamiche, Imane, et al. "A continuous smartphone authentication method based on gait patterns and keystroke dynamics." *Journal of Ambient Intelligence and Humanized Computing* 10.11 (2019): 4417-4430.
- [29] Odili, Julius Beneoluchi, Mohd Nizam Mohmad Kahar, and Shahid Anwar. "African buffalo optimization: A swarm-intelligence technique." *Procedia Computer Science* 76 (2015): 443-448.
- [30] Mukuta, Yusuke, and Tatsuya Harada. "Kernel approximation via empirical orthogonal decomposition for unsupervised feature learning." *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2016.

#### AUTHORS



Gogineni Krishna Chaitanya research scholar received his Bachelors Degree in Computer Science from Acharya Nagarjuna University and Masters Degree from JNTUK. He is currently pursuing Ph.D degree with Department of Computer Science and Engineering Koneru Lakshmaiah Education Foundation, Vaddeswaram, 522502 Andhra Pradesh, India. His research interests include digital forensics, Biometrics, Authentication and Machine Learning.



Dr. K Raja Sekhar received his Ph.D Degree in Computer Science and Engineering from Acharya Nagarjuna University. He is currently a Professor with the of Computer Science and Engineering Koneru Lakshmaiah Education Foundation, Vaddeswaram, 522502 Andhra Pradesh, India. He has published more than 50 articles in journals and Conference proceedings. His research interests include Digital forensics, Biometrics, Network Security and Usable security. He received several Excellence awards and several best paper awards. He has been on the editorial boards of several journals.

# Electricity Cost Prediction using Autoregressive Integrated Moving Average (ARIMA) in Korea

Safdar Ali<sup>1</sup>

Department of Software Engineering  
University of Lahore, Lahore  
Lahore, Pakistan

Do-Hyeun Kim<sup>2</sup>

Department of Computer Engineering  
Jeju National University  
Jeju, Republic of Korea

**Abstract**—Electricity cost plays a vital role due to the immense increase in power utilization, rise in energy rates and alarms about the variations and impact on the environment which ultimately affects electricity cost. We claim that electrical power utilization data became more beneficial if it is presented to the customers along with the prediction of power consumption, prediction of energy prices and prediction of its expected electricity cost. It will assist the residents to alter their power utilization behavior, and thus will have an optimistic influence on the electricity production companies, dissemination network and electricity grid. In this study, we present a residential area power cost prediction by applying the Autoregressive Integrated Moving Average (ARIMA) technique in Korean apartments. We have investigated the energy utilization data on the foundation of daily, weekly and monthly power utilization. The accumulated data constructed on daily, weekly and monthly utilization are selected. Then we predict the maximum and average power consumption cost for each of the predicted daily, weekly and monthly power consumption. The power consumption and general price (General Electricity Price in Korea) data of Korea are used to analyze the efficiency of the prediction algorithm. The accuracy of the power cost prediction using the ARIMA model is verified using the absolute error.

**Keywords**—Electricity price; electricity cost; Autoregressive Integrated Moving Average (ARIMA); prediction; energy consumption

## I. INTRODUCTION

Power cost prediction and forecasting have become one of the main areas of interest to the researchers and experts in energy markets due to the fluctuation of electricity cost. This leads to the requirement of accurate and efficient electricity cost prediction methodology. All the stakeholders of energy including, energy market players and electricity price regulators now think to give attention to the cost prediction and its evolution process. Similar to price prediction, market electricity cost estimation is vital information to organize bidding policies and increased their profits in electricity the market. Furthermore, it can also support customers to reduce their electricity expenses by appropriate consumption of power and the smart grid will be operated smoothly and efficiently with the satisfactory level of consumers' needs and power generation companies.

In the literature there are many approaches to price prediction that can be extended to the cost prediction. These techniques can be generally categorized into two classes. The

first one is Artificial Neural Networks (ANNs) techniques. The second kind of price forecasting method is the time series approaches. ANN's techniques that possess outstanding strength and error clemency are the finest way to address the multi-layered nonlinear problems. ANN has acknowledged the responsiveness of investigators due to its rich modeling process, ease of implementation and decent performance in resolving linear and nonlinear problems. Until now it is effective to formulate and predicts variations in complicated power system using ANN techniques. Various ANN-based techniques have been presented to estimate energy prices in almost all energy markets [1–6].

To intensify the forecasting correctness, it has been employed based on supervised neural learning methodologies [7, 8]. These studies repeatedly applied neural network approaches, which comprises several attributes. These attributes are uninterruptedly examined by knowledge, and the approaches become hard to be established properly [9]. Moreover, it has been apparent that while the ANN models give minor inaccuracy throughout training the input data patterns, the inaccuracy in testing these patterns is typical of a greater order [10], in other ways we can say that this technique when applied to systems in the real world, the prediction accuracy is compromised to an unacceptable level. Moreover, this method is required to convert the characters of all the glitches into quantities and alter all the implications into the arithmetic calculation. Yet, it will certainly result in the loss of important data which will degrade the prediction accuracy.

The non-stationary time series approaches like Autoregressive Integrated Moving Average (ARIMA) [11], stationary time series models like autoregressive (AR) [12], Dynamic Regression (DR) and Transfer Function (TF) [13] are formulated previously for prediction of energy price. Stationary time series methods can be extended for forecasting of power costs. In the modest electricity power marketplaces, the series of electricity costs describes features like: great occurrence, non-constant average and adjustment, every day, weekly, monthly and seasonable schedule outcome on holidays and community vacations; great instability and a great proportion of infrequent power consumption cost [14]. It is not simple to forecast electricity costs precisely, consequently, it requires exceptional dealing in case of predicting electricity cost changes. A hybrid technique to estimate day-ahead electricity price is presented in [15]. The technique is constructed on wavelet transform, ARIMA method and Radial Basis Function Neural Networks (RBFN).

---

Funded by National Research Foundation (NRF) of Korea. Any correspondence associated to this paper should be directed to Dohyeun Kim

The price estimating model for electric energy market stakeholders to minimize the danger of price volatility is proposed in [16]. The method integrates the Enhanced Probability Neural Network (EPNN), Probability Neural Network (PNN) and Orthogonal Experimental Design (OED). Another work described the fusion of feature selection method which is constructed on common information technique and wavelet transformation [17].

Some evolutionary algorithms are used for estimation in different fields of study which can also be used for electricity price and cost prediction. For example, a technique called sunflower optimization algorithm is used for prediction of circuit-based model known as proton exchange membrane fuel cell (PEMFC) [18]. The model is used to reduce the error of sum of squared of predicted and real output voltage. Another method proposed for the reduction of sum of the squared error for PEMFC is presented in [19]. Both models achieved acceptable results and minimized the gap between actual parameters and predicted parameters. A cost optimization model for hybrid energy system is presented in [20]. The model achieved better results as compared to existing methodologies. A multi-objective optimization algorithm for heat pump problem is proposed in [21]. The model produces better results.

In this paper we proposed electricity cost prediction methodology using ARIMA model in republic of Korea. The power consumption in Korean apartments and general electricity price are considered to evaluate the cost and its prediction accuracy. The prediction results of ARIMA (0, 1, 1) model are evaluated using the absolute error (AE).

The rest of the paper is structured as follows. Section II describes the power consumption and cost prediction model. Section III explains the ARIMA-base cost prediction in detail. Section IV describes the case study of price prediction in Korea. Section V puts light on the analysis and result discussion and Section VI concludes the paper.

## II. POWER CONSUMPTION SCENARIO AND COST PREDICTION MODEL

The power consumption scenario which is adopted for cost prediction is based on the different electricity power consumption of the day, week and month for the apartments of Hwasung building number 417. We considered hourly power consumption data for the daily electricity cost prediction. For weekly basis cost prediction we considered power consumption on each day of the week and divided each day in three slots each of 8 hours. Then we calculated the maximum and average power consumption for each slot of the day. For monthly basis cost prediction we considered daily basis power consumption and calculated the maximum and average power consumption for each day of the month. The flow chart of the power consumption scenarios and power prediction model is shown in Fig. 1. The flow chart consists of several components. Smart meters are used to record the power consumption data of the building. The collected consumed power and price statistics are stored in database for further processing. The energy consumption and price data are loaded to the program for power cost prediction. After reading the data in the program, first the data is analyzed with respect to

the duplication and if there is any duplicated data found we simply remove it by averaging technique. For each of the daily, weekly and monthly basis cost prediction we did not merely considered power consumption for one specific day, or week but we considered all the data of 30 days power consumption for each of the daily and weekly basis cost prediction.

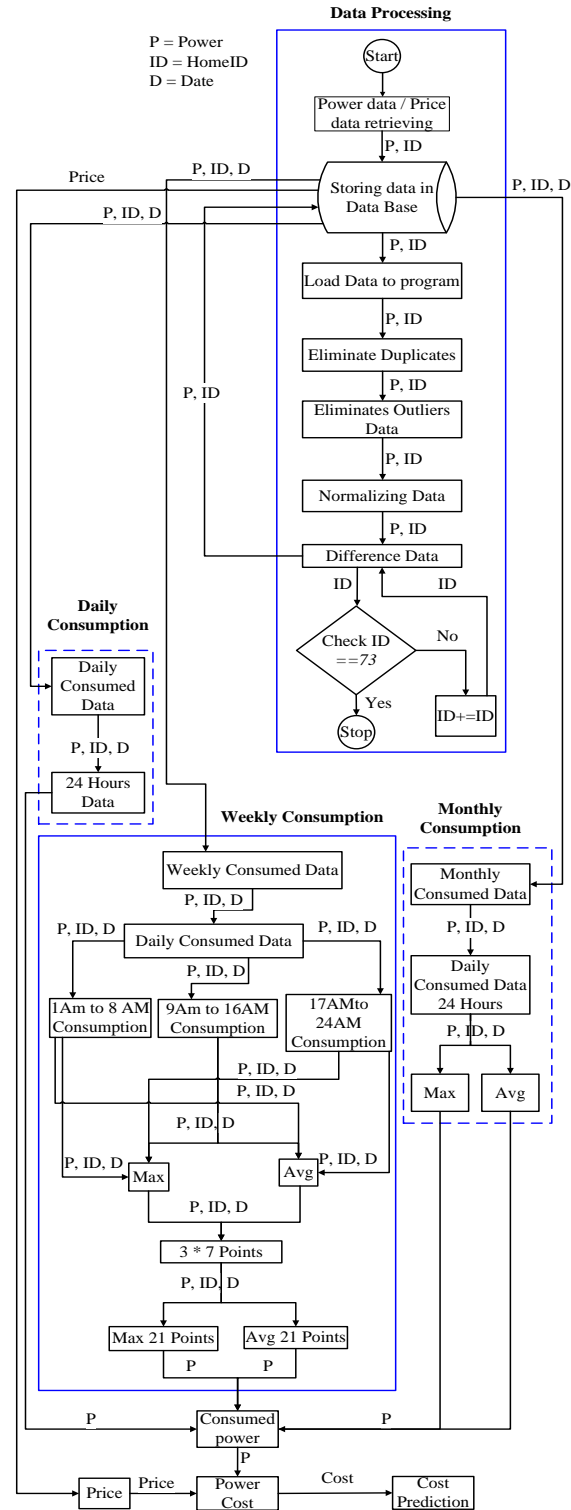


Fig. 1. Power Consumption Scenarios and Cost Prediction Model.



For example, for daily basis power cost prediction we choose a representative daily basis power consumption patterns for the month of January 2010. The representative day is the average power consumption for that particular month against each hour of the whole month for all of the Hwaung building number 417. Each corresponding hour of the day is considered and then we took the average of that hour for whole month. At the end we get a representative daily basis power consumption pattern for the month of January 2010.

Similarly, for weekly basis power consumption we considered daily basis power consumptions for all of the 30 days. The month is divided in to 4 weeks appropriately. Then we got one representative week for January 2010 by taking the average of each corresponding day of the week. At the end each day of the representative week is divided in to three slots as described in the start of this section. For monthly basis power cost prediction, we considered January 2010 data and then find the maximum and average power consumption during each day of the month.

### III. COST PREDICTION USING ARIMA

ARIMA (0, 1, 1) forecasting method is one of the modified version of the ARIMA (p, d, q) model. The ARIMA (p, d, q) prediction method is the mutual category of modeling for predicting time series data. The model can be made static by using some kind of alterations like differencing and logging. In reality, one of the coolest method to talk regarding ARIMA technique is as fine-tuned kinds of random-walk and random-trend technique. The fine-tuning contains adding lags of the differenced series and lags of the predicted errors to the forecasting equation as mandatory to remove any preceding touches of autocorrelation from the predicted errors. In ARIMA (p, d, q) model,  $p$  is the number of autoregressive terms,  $d$  is the number of non-seasonal variances, and  $q$  is the number of lagged predicted error in the forecasting equation.

$$Z_{(k)} = \mu + Z_{(t-1)} - \beta \times e_{(t-1)} \quad (1)$$

$$\mu = Z_t - Z_{(t-1)} \quad (2)$$

$$\beta = 1 - \epsilon \quad (3)$$

Where  $Z_{(t)}$  is the estimation,  $\beta$  is the coefficient of the lagged estimation error,  $e_{(t-1)}$  represents the error at time epoch  $t-1$ ,  $\epsilon$  rate varies within limit [0, 1]. To get optimal predicted values, we ran ARIMA (0, 1, 1) technique twenty times for each 24-hours of the day and then took the average of the twenty points.

### IV. A CASE STUDY IN REPUBLIC OF KOREA

The proposed ARIMA (0, 1, 1) based prediction is verified using a case study of predicting electricity cost on one day, one week and on monthly basis. The electricity price and power consumption data are collected on an hourly basis for a one month of January 2010 and then took an average of each corresponding hour to get one sample day of 24 hours. The cost for each hour is calculated. The unit of price is in Korean WON. Fig. 2 and 3 shows the 24 points of each hour for one day price and energy consumption, respectively. In this work, as an example we have shown only electricity price and power consumption data values of one day to calculate and then

predict hourly basis one day power cost. In order to perform sufficiently and to show the flexibility and reliability of the ARIMA (0, 1, 1) model for electricity cost prediction, section V shows that the ARIMA algorithm perform well when we consider maximum and average power consumption during each slot of the day on different week days. The hourly electricity cost values have been estimated based on the ARIMA (0, 1, 1) model.

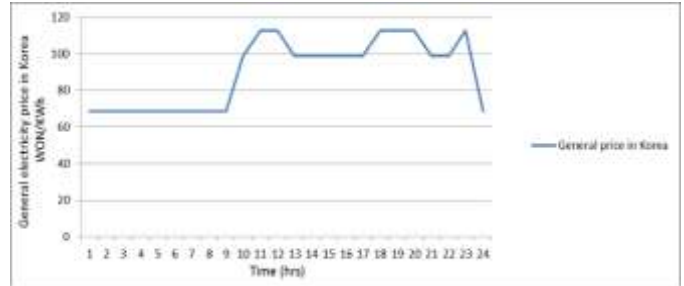


Fig. 2. General Electricity (Actual/Original) Prices, on (01/01/2010).

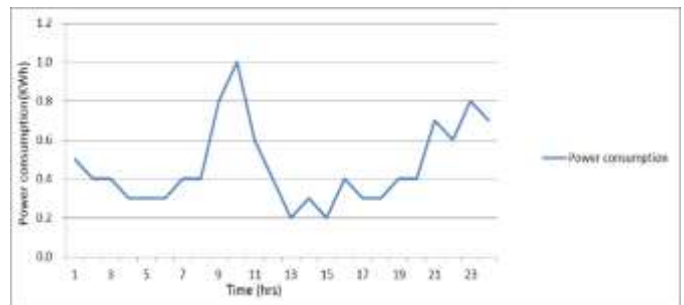


Fig. 3. Power/Electricity Consumption on (01/01/2010).

### V. ANALYSIS AND RESULT DISCUSSION

In this section we are evaluating ARIMA (0, 1, 1) model performance with respect to electricity cost prediction. In the literature many objective functions exist to assess the prediction algorithm. We used absolute error (AE) method as the evaluation criteria. Like other objective functions, estimation performance of AE is much improved when the objective function value is minor. Equation (4) below describes the (AE) objective function.

$$AE = |p_p - a_p| \quad (4)$$

In equation (4)  $a_p$  means real time original/actual electricity cost,  $p_p$  means predicted electricity cost.

Fig. 4 shows the daily basis actual cost and predicted cost based on the ARIMA (0, 1, 1) technique for Hwasung building number 417. Each hour of the day represents the average power cost of all homes of building number 417. We have considered the experimental investigation of the electricity cost data by applying ARIMA (0, 1, 1) method. In Fig. 4 there are three lines. Sky color line described electricity cost based on the general price in Korea during one representative day of January 2010, Indian red color shows the cost prediction line using ARIMA (0, 1, 1) technique and yellow green color represents the absolute error for ARIMA (0, 1, 1) technique. From the Fig. 4 we can see that the prediction based on the ARIMA (0, 1, 1) almost follows the actual cost line. If we see

the absolute error line, then we can say that the absolute error line some time approaching to zero during some hours of the day, which means that the prediction of cost using proposed ARIMA (0, 1, 1) method in this case up-to somehow follows the actual electricity cost line. As the prediction accuracy between actual and predicted cost increase, the absolute error decreases.

Fig. 5 shows the efficiency of the proposed ARIMA (0, 1, 1) model based on the weekly and monthly power consumption scenario as shown in Fig. 1.

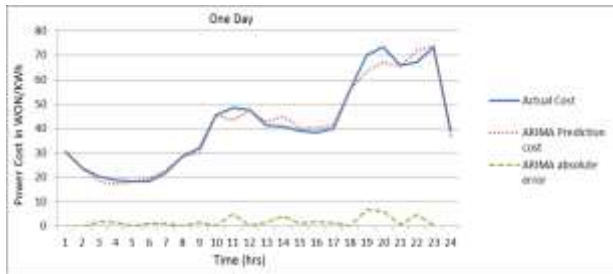


Fig. 4. Actual Cost, Predicted Cost using ARIMA (0, 1, 1) Model and Absolute Error for One Day.

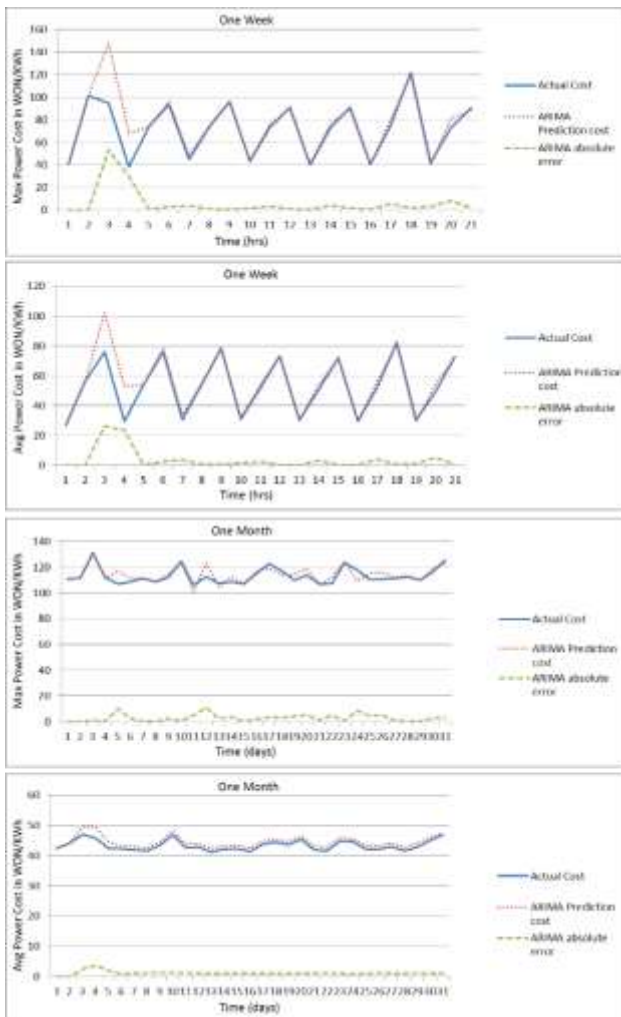


Fig. 5. Actual Cost, Predicted Cost using ARIMA (0, 1, 1) Model and Absolute Error based on the Power Consumption Scenarios.

The aim of testing ARIMA model on different power consumption scenarios is to authenticate the efficiency of the prediction algorithm. In order to estimate daily electricity cost for each day of the week and month, the general electricity price and power consumption data are collected for one month of January 2010. From the Fig. 5 we can see that (AE) for weekly basis electricity cost prediction is very less except for the hour 4 & 5. The performance of the ARIMA model for the maximum and average weekly electricity cost prediction is good.

Fig. 5 also shows the efficiency of the proposed ARIMA (0, 1, 1) model based on the monthly power consumption scenario. The (AE) values in Fig. 5 also describe the efficiency of the proposed algorithm for prediction of electricity cost on monthly basis. The maximum and average monthly electricity predicted cost is very good by looking into (AE) values. The AE values are much near to zero which means the error between actual cost and predicted cost is less and almost zero.

So the bottom line is that proposed ARIMA (0, 1, 1) model effectively predicted the electricity cost using different power consumption scenarios.

## VI. CONCLUSIONS

The management of energy system depends on the good knowledge of the electricity cost due to the deregulation of the energy market prices and cost. In the literature a number of different simulation and prediction methodologies have been proposed and applied for a number of times to forecast the electricity prices. These methodologies can also be used for power cost prediction. The Autoregressive Integrated Moving Average (ARIMA) method is an optimal prediction technique and has been established and used successfully previously in other fields for prediction. In this paper we proposed it for electricity cost prediction. The prediction results of the ARIMA (0, 1, 1) model is satisfactory. The efficiency and reliability of the ARIMA (0, 1, 1) method is verified using a case study in Republic of Korea. The proposed methodology has many benefits with respect to electricity cost. The proposed technique validates the effectiveness of ARIMA (0, 1, 1) approach for electricity cost prediction.

Electricity cost of Republic of Korea reveal periodicity and time-varying uncertainty credited to the non-storability of electricity, insignificant costs and potential for electricity power market. Our aim is to propose and evaluate a technique to separate periodicity that could be effortlessly reasonable in application. The proposed ARIMA (0, 1, 1) technique also produced good results which are reasonably acceptable for current research. The proposed ARIMA (0, 1, 1) technique takes into account the periodicity of the electricity cost data and up-to somehow can successfully solve the prediction issue along with periodicity. The proposed model is essentially simple, typical and do not need to make tough variations and judgments concerning the clear form of the model for different trends of the power cost.

The proposed prediction technique using ARIMA (0, 1, 1) model results in minimum AE. Finally based on the results and above arguments we can say that the ARIMA (0, 1, 1)

prediction technique performs satisfactory up-to some level. From the above discussions, it is vibrant that the proposed ARIMA (0, 1, 1) model is effective for electricity cost prediction.

There is limitation of ARIMA model to predict electricity cost and price. If there is various kind of non-linearity in the data and there is no specific patterns in the data, then AE values getting increases which affects the accuracy rate.

#### ACKNOWLEDGMENT

This research is supported by Energy Cloud R&D Program via the National Research Foundation of Korea (NRF) funded by the Ministry of Science, ICT (2019M3F2A1073387), and this research is also supported by Institute for Information & communications Technology Planning & Evaluation (IITP) grant funded by the Korea government (MSIT)(No.2018-0-01456, AutoMaTa: Autonomous Management framework based on artificial intelligent Technology for adaptive and disposable IoT). Any correspondence related to this paper should be addressed to Dohyeun Kim.

#### REFERENCES

- [1] H.Y. Yamin; S.M. Shahidehpour, and Z. Li, "Adaptive short-term electricity price forecasting using artificial neural networks in the restructured power markets," *Electr Power Energy Syst*, vol. 26, pp. 571–81, 2004.
- [2] P. Mandal, T. Senjyu, N. Urasaki, T. Funabshi, and A. K. Srivastava, "A novel approach to forecast electricity price for PJM using neural network and similar days method," *IEEE Trans Power Syst*, vol. 22, pp. 2058–2065, 2007.
- [3] H. T. Pao, "Forecasting electricity market pricing using artificial neural networks," *Energy Convers Manage*, vol. 48, pp. 907–912, 2007.
- [4] R. Pino, J. Parreno, A. Gomez, and P. Priore, "Forecasting next-day price of electricity in the Spanish energy market using artificial neural networks," *Eng Appl Artif Intell*, vol. 21, pp. 53–62, 2008.
- [5] N. M. Pindoriya, S. N. Singh, and S. K. Singh, "An adaptive wavelet neural network-based energy price forecasting in electricity markets," *IEEE Trans Power Syst*, vol. 23, pp. 1423–1432, 2008.
- [6] N. Amjady, A. Daraeepour, and F. Keynia, "Day-ahead electricity price forecasting by modified relief algorithm and hybrid neural network," *IET Generat Trans Distribut*, vol. 4, pp. 432–444, 2010.
- [7] G. Raquel, M. Romeo-Luis, and A. Gil, "Forecasting of electricity prices with neural networks," *Int J Energy Convers Manage*, vol. 47, pp. 1770–1778, 2006.
- [8] J. P. S. Catalao, S. J. P. S. Mariano, V. M. F. Mendes, and L. A. F. M. Ferreira, "Short-term electricity prices forecasting in a competitive market: a neural network approach," *Electr Power Syst Res*, vol. 77, pp. 1297–1304, 2007.
- [9] R. Lapedes, R. Farber, "Nonlinear signal processing using neural networks: prediction and system modeling," Technical report LA-VR87-2662. Los Alamos, New Mexico: Los Alamos National Laboratory, 1987.
- [10] E. R. J. Van, "The application of neural network in the forecasting of share prices," *Finance and Technology Publishing*, 1996.
- [11] Y. Y. Hong, and C. F. Lee, "A neuro-fuzzy price forecasting approach in deregulated electricity markets," *Elect Power Syst Res*, vol. 73, pp. 151–157, 2005.
- [12] C. P. Rodriguez, and G. J. Anders, "Energy price forecasting in the Ontario competitive power system market," *IEEE Trans Power Syst*, vol. 19, pp. 366–374, 2004.
- [13] G. Li, C. C. Liu, C. Mattson, and J. Lawarree, "Day-ahead electricity price forecasting in a grid environment," *IEEE Trans Power Syst*, vol. 22, pp. 266–274, 2007.
- [14] F. J. Nogales, J. Contreras, A. J. Conejo, and R. Espinola, "Forecasting next-day electricity prices by time series models," *IEEE Trans Power Syst*, vol. 17, pp. 342–348, 2002.
- [15] M. Shafie-Khah, M. M. Parsa, and Sheikh-El-Eslami, "Price forecasting of day-ahead electricity markets using a hybrid forecast method," *Energy Conversion and Management* vol. 5, pp. 2165-2169, 2011.
- [16] Lin. Whei-Min, Gow. Hong-Jey, Tsai. Ming-Tang, "Electricity price forecasting using Enhanced Probability Neural Network," *Energy and Conversion Management* vol. 51, pp. 2707-2714, 2010.
- [17] H. Shayeghi, and A. Ghasemi, "Day-ahead electricity prices forecasting by a modified CGSA technique and hybrid WT in LSSVM based scheme," *Energy Conversion and Management* vol. 74, pp. 482-491, 2013.
- [18] Z. Yuan, W. Wang, H. Wang, N. Razmjoooy, "A new technique for optimal estimation of the circuit-based PEMFCs using developed Sunflower Optimization Algorithm," *Energy Rep*, vol. 6, pp. 662-671, 2020.
- [19] Z. Yang, Q. Liu, L. Zhang, J. Dai, N. Razmjoooy, "Model parameter estimation of the PEMFCs using improved Barnacles Mating Optimization algorithm," *Energy Rep*, vol. 212, pp. 0360-5442, 2020.
- [20] Y. Guo, X. Dai, K. Jermsittiparsert, N. Razmjoooy, "An optimal configuration for a battery and PEM fuel cell-based hybrid energy system using developed Krill herd optimization algorithm for locomotive application," *Energy Rep*, vol. 6, pp. 885-894, 2020.
- [21] X. Fan, H. Sun, Z. Yuan, Z. Li, R. Shi, N. Razmjoooy, "Multi-objective optimization for the proper selection of the best heat pump technology in a fuel cell-heat pump micro-CHP system," *Energy Rep*, vol. 6, pp. 325-335, 2020.

# Deep Learning Bidirectional LSTM based Detection of Prolongation and Repetition in Stuttered Speech using Weighted MFCC

Sakshi Gupta<sup>1</sup>, Ravi S. Shukla<sup>2</sup>, Rajesh K. Shukla<sup>3</sup>, Rajesh Verma<sup>4</sup>

Department of Computer Science and Engineering, Invertis University, Bareilly, Uttar Pradesh, India<sup>1,3</sup>  
Department of Computer Science, Saudi Electronic University, Kingdom of Saudi Arabia<sup>2</sup>  
Department of Electrical Engineering, King Khalid University, Kingdom of Saudi Arabia<sup>4</sup>

**Abstract**—Stuttering is a neuro-development disorder during which normal speech flow is not fluent. Traditionally Speech-Language Pathologists used to assess the extent of stuttering by counting the speech disfluencies manually. Such sorts of stuttering assessments are arbitrary, incoherent, lengthy, and error-prone. The present study focused on objective assessment to speech disfluencies such as prolongation and syllable, word, and phrase repetition. The proposed method is based on the Weighted Mel Frequency Cepstral Coefficient feature extraction algorithm and deep-learning Bidirectional Long-Short term Memory neural network for classification of stuttered events. The work has utilized the UCLASS stuttering dataset for analysis. The speech samples of the database are initially pre-processed, manually segmented, and labeled as a type of disfluency. The labeled speech samples are parameterized to Weighted MFCC feature vectors. Then extracted features are inputted to the Bidirectional-LSTM network for training and testing of the model. The effect of different hyper-parameters on classification results is examined. The test results show that the proposed method reaches the best accuracy of 96.67%, as compared to the LSTM model. The promising recognition accuracy of 97.33%, 98.67%, 97.5%, 97.19%, and 97.67% was achieved for the detection of fluent, prolongation, syllable, word, and phrase repetition, respectively.

**Keyword**—Speech; stuttering; deep learning; WMFCC; Bi-LSTM

## I. INTRODUCTION

For communication between human beings, speech proves to be the most habitually and widely used verbal means to precise feelings, ideas, and thought. Not all human beings are blessed with normal means of speech. The potency of speech in delivering data during communication depends on fluency. Fluency is defined by normal speech flow, which connects different phonemes to make a message [1]. Speech is fluent if continuity among semantic units, rhythm, speed, and energy applied for flow is normal. Any kind of disruption in fluency is known as dysfluency. Stuttering is a complex type of dysfluency. In stuttering, there is a disturbance in continuity and rhythm due to pauses and blocks, the rate is much slower, and efforts are higher than normal. Researchers have categorized the factors that lead to stuttering as of three types, namely, development, neurogenic, and psychogenic.

People who stutter (PWS) may have three sorts of disfluencies: repetition of a sound, syllable, word or phrase, sound prolongation during which a sound is sustained for a markedly more extended period that may be traditional and silent blocks at starting of vocalization or word or within the middle of a word. Johnson [2] introduced this classification for the first time. It has been used by clinicians and researchers ever since.

Even though stuttering may not be considered as a disability by many people, it incites a speech constraint. People who stutter loses not only their confidence but also generate a negative attitude towards their communication skills. Furthermore, it ruins their self-confidence, relationship with others, employment opportunities, and opinions of others about them [3]. Stuttering influence individuals of all ages, culture, and races irrespective of their intelligence and financial status. Many pieces of research have stated that stuttering affects approximately 1% of the world population and is more common in males as compared to females [4]. Therefore, this area is mainly a knowledge base field of analysis for different domains like speech pathology, psychology, speech physiology, acoustics, and signal analysis.

Stuttering is one of the intense issues found in speech pathology. Speech-Language Pathologists (SLP) diagnoses the individual who stutters and measures the fluency to gauge the response of the stutterer throughout the treatment process. Traditionally SLPs used to assess the extent of stuttering manually. They counted and divided the frequency of stuttered events with total spoken words. Such sorts of stuttering assessments are arbitrary, incoherent, lengthy, and error-prone. Over the past two decades, SLPs gave great attention to objective assessment techniques for assessing the stuttered events, as discussed in our previous work [5].

Automatic evaluation of stuttered speech is therefore necessary, to automate the count and classification of stuttered events. The proposed work has employed Weighted Mel Frequency Cepstral Coefficients (WMFCC) feature extraction method and deep-learning-based classification method Bi-directional Long-Short Term Memory (Bi-LSTM) for the automatic assessment of four forms of disfluency prolongation and syllable, word, and phrase repetition. The efficacy of the Bi-LSTM model is assessed as compared to other

classification models, based on the accuracy of the classification of stuttered events.

In this paper, the University College London Archive of Stuttered Speech (UCLASS) database is utilized for analysis. The experimental analysis in this study reveals that WMFCC and Bi-LSTM based proposed method performs more efficiently as compared to other models.

The results elucidate that the model proposed has improved performance and advantages compared with other models. This study makes two significant contributions.

- Firstly, it uses WMFCC instead of traditional MFCC for feature extraction. WMFCC includes the dynamic information of the speech samples, which increases the detection accuracy of stuttered events; and also reduces the computational overhead to the classification stage.
- Secondly, it employs Bi-LSTM rather than traditional RNN and LSTM. Bi-LSTM provides the solution for gradient disappearance in RNN, as well as overcomes the unidirectional flow of information of LSTM.

The paper is structured according to the following. Section 2 reviews the work related to automatic detection of stuttering speech disorders. Section 3 elaborates on the framework for the system proposed. It also includes brief descriptions of the database used, feature extraction, and classification techniques applied. Section 4 consists of experimental results and a comparative analysis of the classification model. Section 5 provides a conclusion.

## II. RELATED WORKS

This section reviews work relating to recognition systems designed to detect or classify stuttering speech disorders; previous research has presented various methods and algorithms that have been applied to recognizing stuttering events from speech signals. Table I displays a comprehensive comparative analysis of various feature extraction and classification methods based on the dataset used, type of disfluency, and accuracy. The previous works conducted signifies the importance of feature extraction and classification methods in the stuttered events detection.

Traditional machine learning techniques are being gradually replaced by Deep learning technology. Deep learning provides a more accurate representation of objects

and can automatically obtain objects features from a vast amount of data [26]. These are progressively used to further refine computers' capacities in order to understand what humans can do, including speech recognition. Deep structured learning models based on these functional attributes include convolutional neural network (CNN) [27], recurrent neural network (RNN) [28][24], and long-short term memory (LSTM) [25]. The conventional machine learning techniques for recognition employed shallow structured architectures such as hidden Markov model (HMM), Support Vector Machines (SVM), Artificial Neural Network (ANN), and linear and non-linear dynamical system [29]. These architectures are ideally suited for simple or constrained problems, since their limited capabilities can cause problems in complicated large-scale real-world problems [30]. Such real-world problems involve human speech, language recognition, and visual scenes, requiring a more profound and layered architecture to extract the complex information.

Tian Swee et al. [6] and Thiang and Wanto [9] trained Hidden Markov Model (HMM) model to classify speech samples as fluent and non-fluent. The HMM model determines the likelihood of being in a state depends on its prior state at (t-1) while disregarding all other dependence. It also requires a large number of parameters and data for building and training the model [31]. In [8] and [14], Ravikumar et al. and Hariharan et al. discussed the classification of extracted features through Support Vector Machines (SVM). However, SVM deals with only fixed-size input are not efficient for large databases as well as its computational cost is directly proportional to the number of classes to be classified. Savin et al. [19] employed an ANN for classification. ANN does not have structured methodology as well as time-consuming for large networks [32].

The deep learning technique CNN performs very well on non-sequential data while fails in interpreting temporal information. However, the RNN is good at modeling the temporal data but suffers from the problem of short-term memory caused by vanishing gradient [33][34][35]. Thus, LSTM was created as a solution to short-term memory [36]. They are capable of learning long term dependencies [37]. Based on the above considerations, this paper applies Bi-LSTM for the classification of a vast amount of speech data [38]. Bi-LSTM model processes the information in two directions and links them to obtain the output class of stuttering.

TABLE I. COMPREHENSIVE ANALYSIS OF VARIOUS ON RESEARCH ACTIVITIES ON STUTTERING DETECTION, DESCRIBING THE FEATURES USED, CLASSIFIER EMPLOYED, NUMBER OF SUBJECTS, TYPE OF CLASSIFICATION AND EXPERIMENTAL RESULTS

| Year      | Feature Used                                   | Classifier Used  | Dataset Used   | Type of Classification  | Result  |
|-----------|--|--|--|---|---|
| 2007 [6]  | MFCC   | HMM  | Malay language-based 20 normal and 15 artificial speech samples                | Repetition, Prolongation, and Blocks                          | Normal data- 96%, Artificial Stutter Speech data- 90%   |
| 2009 [7]  | Kohonen Network                                | Multi-layer Perceptron and RBF   | 59 800ms samples of 8 stuttering Polish speakers                               | Repetition and Prolongation                                   | MLP- 92%<br>RBF- 91%  |
| 2009 [8]  | MFCC   | SVM  | 12 training and 3 testing samples of 15 adults who stutter                     | Syllable Repetition   | 94.35%  |
| 2010 [9]  | LPC  | HMM  | 5, 10, 15, 20 samples per command and 40-50 observation symbols of HMM         | -   | 5 samples-93.75%,<br>10- 98.75%,<br>15- 100% and<br>20-97.5%                                      |
| 2010 [10] | MFCC   | KNN and LDA  | 10 samples of 8 males and 2 females (11 to 20 years) from UCLASS               | Repetition and Prolongation                                   | 90%   |
| 2010 [11] | LPCC   | KNN and LDA  | 10 samples of 8 males and 2 females (11 to 20 years) from UCLASS               | Repetition and Prolongation                                   | 88.05%  |
| 2011 [12] | 12, 13, 26 and 39 Dimensional MFCC             | DTW  | 8 training and 2 testing samples   | Repetition  | 12 D- 80.69%,<br>13 D- 68.4%,<br>26 D- 84.01%,<br>39 D- 84.58%,                                   |
| 2012 [13] | MFCC and LPCC                                  | KNN and LDA  | UCLASS database  | Repetition and Prolongation                                   | MFCC- 92.55%, LPCC- 94.51%  |
| 2012 [14] | Spectral Entropy using Bark, Mel and Erb Scale | SVM  | UCLASS database  | Repetition and Prolongation                                   | Average accuracy- 96%. Beat result of 96.84% in Erb scale   |
| 2013 [15] | MFCC, PLP, and LPC                             | KNN, LDA, and SVM  | UCLASS database  | Repetition and Prolongation                                   | Best average classification accuracy is given by SVM using the WLPCC, PLP, and MFCC features- 95% |
| 2013 [16] | SOM  | Hierarchal ANN, MLP  | 153 recordings of 19 PWS   | Blocks, syllable repetition and syllable initial prolongation | Blocks- 96%<br>Syllable Repetition- 84% and Prolongation-99%                                      |
| 2014 [17] | MFCC   | SVM  | UCLASS database  | Repetition and Prolongation                                   | 97.6%   |
| 2015 [18] | MFCC   | KNN  | 80 speech samples for training and 20 for testing                              | Repetition with 0db to 10db babble noise                      | 60-95% depending on the sound used  |
| 2016 [19] | MFCC, Formant, Pitch, ZCR, and Energy          | ANN  | 78 recordings of 4 PWS (25-40 years)   | Repetition and Prolongation                                   | 88.29%  |
| 2016 [20] | MFCC, Formant and Shimmer                      | DTW  | 50 repetition events   | Repetition  | 94%   |
| 2016 [21] | MACV   | Thresholding   | 5 Stuttering person speech samples from UCLASS database                        | Repetition and Prolongation                                   | 73.29%  |
| 2016 [22] | MFCC and PLP                                   | Cross-correlation, Euclidean distance using Morphological Image Processing | UCLASS database  | Prolongation, word repetition, and phrase repetition          | Prolongation- 99.84%, Word repetition- 98.07% and Phrase repetition- 99.87%                       |
| 2017 [23] | MFCC   | I-Vector   | 1380 segments of 18 PWS from UCLASS. 80% used for training and 20% for testing | Repetition, Prolongation, and Repetition-Prolongation         | Normal- 52.43%, Repetition- 69.56%, Prolongation- 40%, Rep-Pro- 50%                               |
| 2020 [24] | MFCC   | Gated Recurrent CNN  | UCLASS database  | Prolongation and Repetition                                   | Prolongation- 95%<br>Repetition- 92%  |
| 2020 [25] | MFCC   | LSTM   | UCLASS database  | Prolongation, Blocks, and Repetition                          | 4% and 6% higher than ANN and SVM   |

### III. CONSTRUCTION OF MODEL

The proposed work has employed the WMFCC feature extraction method and deep-learning-based classification method Bi-directional Long-Short Term Memory (Bi-LSTM) for the automatic assessment of four forms of disfluency prolongation and syllable, word, and phrase repetition. The process for detection of repetition and prolongation in stuttered speech is split into five stages: signal pre-processing, disfluent speech sample segmentation and labeling, labeled sample splitting into training, validation and test sets, feature extraction and classification using network training and model (Fig. 1). The University College London Archive of Stuttered Speech (UCLASS) database is utilized for analysis [39]. The study evaluates the efficacy of Bi-LSTM model, based on the accuracy of the classification of stuttered events.

#### A. Signal Pre-Processing

A signal is pre-processed by removing the silence regions [40][41]. There is no excitation in the vocal tract during the silence region, hence no speech production. Thus, pre-processing reduces not only the amount of processing but also enhances the overall efficiency and accuracy of the system proposed. The combination of two widely known approaches, namely Short Time Energy (STE) and Zeros Crossing Rate (ZCR) (Fig. 2), has been used in this work [42][43]. It is a fast and straightforward approach and gives a better result of classifying the speech into voiced/unvoiced.

The short-term energy is the energy-related to short term region of speech [41]. The total energy of a speech frame is determined by the following (1).

$$E(n) = \sum_{m=-\infty}^{\infty} (s(m) \cdot w(n - m))^2 \quad (1)$$

Where  $w(n)$  represents the windowing function, and  $n$  is the shift in the number of samples. The voiced region energy is high in comparison with the unvoiced region. The silent region displays marginal energy content.

Zero-Crossing Rate specifies the number of zero crossings in a given signal [41]. The zero-crossing rate of a stationary signal is calculated by (2):

$$ZCR = \sum_{n=-\infty}^{\infty} |sgn(s(n)) - sgn(s(n - 1))| \quad (2)$$

Where  $sgn(s(n))$  is a signum function and is described as by the (3).

$$sgn(s(n)) = \begin{cases} 1 & \text{if } s(n) \geq 0 \\ -1 & \text{if } s(n) < 0 \end{cases} \quad (3)$$

The zero-crossing rates in unvoiced sounds are comparatively high as compared to the voiced sounds. The combination of these two features overcome the issue of categorizing the speech into a voiced/unvoiced signal (Fig. 3).

#### B. Disfluent Speech Sample Segmentation and Labeling

The disfluent speech signals are obtained from the University College London Archive of Stuttered Speech (UCLASS) [39]. It is released in version 1 and version 2, consisting of three types of recording: monologues, reading, and spontaneous conversation. Version 1 has 138 “monologue” recordings contributed by 81 speakers. The

database used in this work refers to 20 samples of speech for experimentation [44]. It comprises two female speakers and 18 male speakers aged 7years 8 months to 17 years 9 months. The selection of speech signals aims at covering a wide variety of stuttering rate and age. The samples provided with text script are only included in the database.

This paper investigates only four forms of disfluencies, prolongation, and syllable, word, and phrase repetition. They are easily detectable in monosyllabic words. After pre-processing the selected speech samples, disfluent speech samples were marked and segmented manually by listening to the pre-processed signals. The segmented samples were labeled as five classes, namely, Fluent, Prolongation, Syllable Repetition, Word Repetition, and Phrase Repetition (Fig. 4).

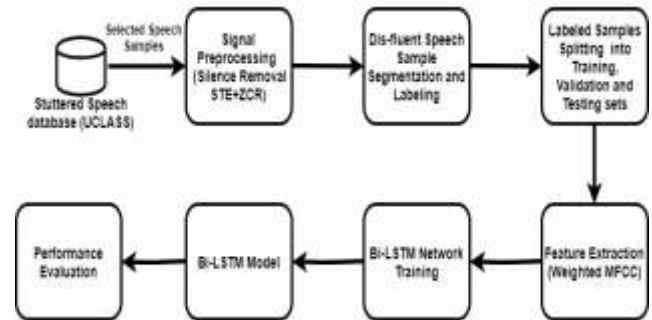


Fig. 1. Block Diagram of the Proposed Model.

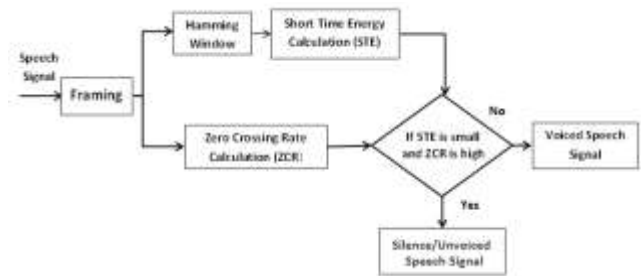


Fig. 2. Speech Pre-Processing by Silence Removal.

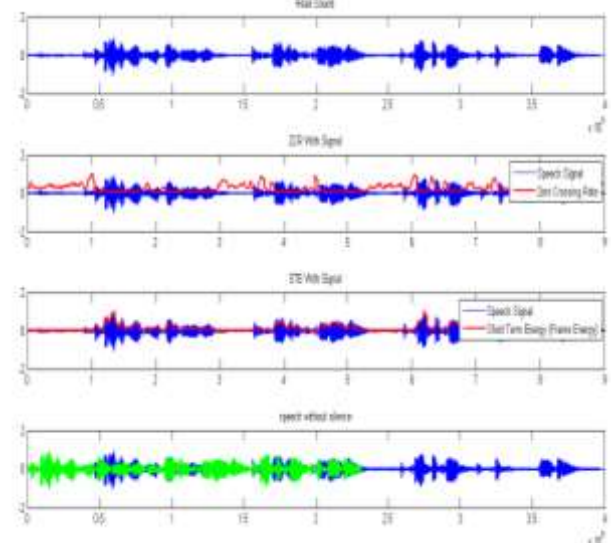


Fig. 3. Silence Removal using STE-ZCR Method.

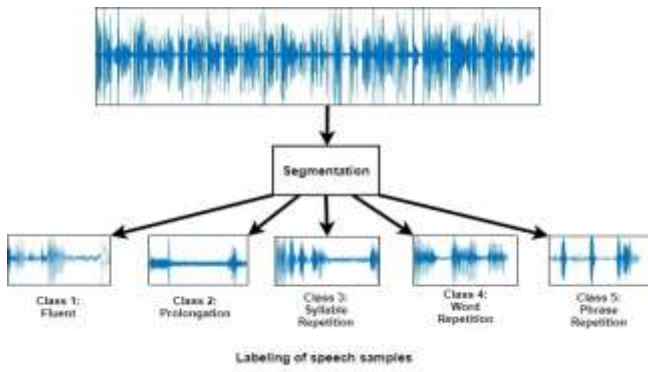


Fig. 4. Disfluent Speech Sample Segmentation and Labeling.

### C. Labeled Samples Splitting

The segmented disfluent speech samples were divided into three sets for training, validation, and testing. The training set is a subset of labeled stuttered speech samples used to train the model. The validation set evaluates the performance of the model with different hyperparameter values. It is smaller than the training set. The test set determines the final accuracy of the model and analyses the performance of different models. In this study, the datastore of disfluent speech samples is split into training, validation, and test set in the ratio of 60%, 20%, and 20%, respectively.

The process of pre-processing, segmentation, labeling, and sample splitting is described through an algorithm in Table II.

TABLE II. ALGORITHM OF SPEECH SAMPLE PRE-PROCESSING AND SEGMENTATION

|   |
|---|
| <p><b>Input:</b> Selected speech samples of UCLASS database <math>T_{unpreprocess}</math><br/> <b>Output:</b> Pre-processed and labeled speech samples dataset <math>T_{train}</math>, <math>T_{valid}</math>, and <math>T_{test}</math></p>  |
| <ol style="list-style-type: none"> <li>1. Loading the selected speech samples of UCLASS database <math>T_{unpreprocess}</math>.</li> <li>2. For each sample <math>S \in T_{unpreprocess}</math></li> <li>3. Divide the samples into frames of 30msec.</li> <li>4. For each frame <math>f \in</math> sample <math>S</math></li> <li>5. Calculate STE and ZCR of frame <math>f</math>.</li> <li>6. If <math>(STE \geq 0.01)</math> and <math>(ZCR \leq 0.2)</math><br/>Append the frame <math>f</math> to new pre-processed sample <math>\in S_{preprocess}</math><br/>Else discard the frame.</li> <li>7. End for.</li> <li>8. End for.</li> <li>9. The set of pre-processed samples <math>T_{preprocess}</math> is derived.</li> <li>10. For each sample <math>S \in T_{preprocess}</math></li> <li>11. Manually divide the speech sample into a set of segments <math>T_{unlabelled}</math>.</li> <li>12. For each segment <math>St \in T_{unlabelled}</math></li> <li>13. Identify and label the segment as fluent, prolongation, syllable repetition, word repetition, or phrase repetition.</li> <li>14. End for.</li> <li>15. End for.</li> <li>16. The set of labeled samples <math>T_{labelled}</math> is derived.</li> <li>17. Split the set <math>T_{labelled}</math> into training <math>T_{train}</math>, validation <math>T_{valid}</math>, and testing <math>T_{test}</math> datasets in the ratio 60%, 20% and 20% respectively.</li> </ol> |

### D. WMFCC Feature Extraction

The extraction of speech features is a sort of dimension reduction technique that is employed to minimize the data that is giant to be processed by an algorithm. The key objective of feature extraction is to upbraid the speech signal into the various acoustically recognizable elements and to get the feature vectors with a nominal amendment to keep the processing efficient. In our previous work [45], a comparative analysis of extensions of MFCC feature extraction techniques [46], namely Delta MFCC, Delta-delta MFCC, and Weighted MFCC [47] was conducted. Its experimental results displayed, WMFCC slightly outperforms Delta-delta MFCC and significantly outperforms Delta MFCC and MFCC in all situations of frame length, alpha values, and frame overlap percentage [45]. The proposed work has applied frequency-domain based Weighted Mel Frequency Cepstral Coefficients. WMFCC is a fusion of MFCC and its derivatives delta and delta-delta. The resultant vector contains both static as well as dynamic information of the signal. Moreover, the feature vector is of size 14; thus, incur less computational overhead to the classification stage. Table III describes the WMFCC feature extraction algorithm, and the results of the algorithm are displayed in Fig. 5.

TABLE III. ALGORITHM OF WMFCC FEATURE EXTRACTION

|  |
|--|
| <p><b>Input:</b> Pre-processed and labeled speech samples dataset <math>T_{train}</math>, <math>T_{valid}</math>, and <math>T_{test}</math>.<br/> <b>Output:</b> WMFCC feature vector of <math>T_{train}</math>, <math>T_{valid}</math>, and <math>T_{test}</math> datasets.</p>   |
| <ol style="list-style-type: none"> <li>1. Load the datasets <math>T_{train}</math>, <math>T_{valid}</math>, and <math>T_{test}</math>.</li> <li>2. Initialize the parameters of the WMFCC feature extraction method.</li> <li>3. For each labeled sample <math>S \in (T_{train}, T_{valid}, \text{ and } T_{test})</math></li> <li>4. Pre-emphasize <math>S</math> using <math>\alpha</math> filter as 0.98.</li> <li>5. Divide the sample into 30msec frames with an overlapping percentage of 75%.</li> <li>6. For each frame <math>f \in S</math></li> <li>7. Apply the Hamming window function to frame <math>f</math>.</li> <li>8. Calculate the power spectrum of the windowed signal using FFT.</li> <li>9. Calculate the Mel spectrum by passing the power spectrum through 20 Mel filters.</li> <li>10. Calculate the log-energy of each filter bank part.</li> <li>11. Calculate MFCC by applying energies to DCT.</li> <li>12. Compute Delta MFCC as:<br/> <math display="block">\Delta c_t = \frac{\sum_{k=1}^M (c_{t+k} - c_{t-k})}{2 \sum_{k=1}^M k^2}</math> </li> <li>13. Compute Delta-Delta MFCC as:<br/> <math display="block">\Delta \Delta c_t = \frac{\sum_{k=1}^M (\Delta c_{t+k} - \Delta c_{t-k})}{2 \sum_{k=1}^M k^2}</math> </li> <li>14. Compute 14-dimensional WMFCC as:<br/> <math display="block">wc_t = c_t + p \cdot \Delta c_t + q \cdot \Delta \Delta c_t, \quad q &lt; p &lt; 1</math> </li> </ol> |



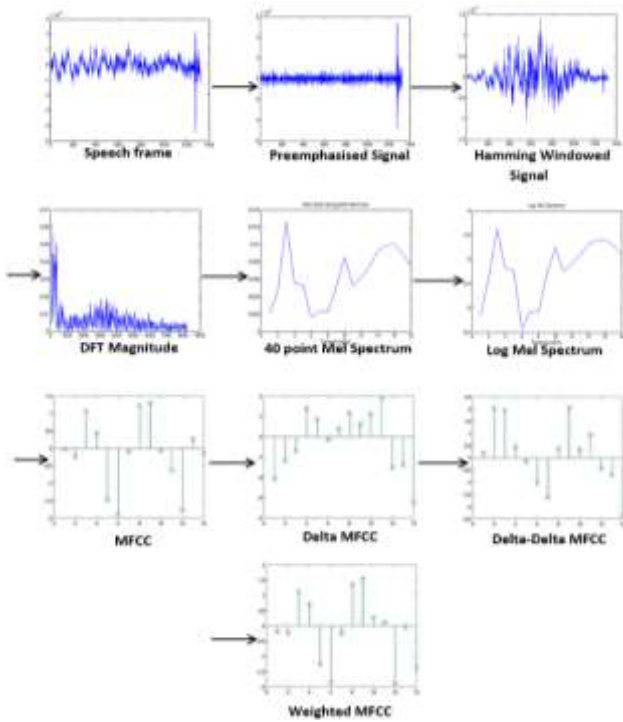


Fig. 5. WMFCC Feature Extraction Process.

### E. Bi-Directional Long-Short Term Memory

Deep learning Bi-LSTM is applied for the classification of stuttered speech samples. It is composed of LSTM cells (Fig. 6). The set of features vectors discussed in the above section are set as input to the classifier. The model is trained and validated with 60% and 20% of the speech samples of the datastore, respectively. The remaining of the samples are used for testing the model.

1) *Long-Short Term Memory*: LSTM is a specialized Recurrent Neural Network (RNN) architecture, competent in learning long term dependencies [48]. RNN suffers from short-term memory, caused by vanishing gradient problem. To mitigate this problem, LSTM has a hidden layer known as the LSTM cell. LSTM cells are built with various gates and cell state that can regulate the flow of information. Like RNNs, at each time iteration,  $t$ , the LSTM cell has the layer input,  $x_t$ , and the layer output,  $h_t$ . The cell also takes the cell input state,  $\tilde{C}_t$ , the cell output state,  $C_t$ , and the previous cell output state,  $C_{t-1}$ . LSTM architecture has three gates, namely, forget, input, and output gate denoted as  $f_t$ ,  $i_t$ , and  $o_t$ , respectively.

The cell state act as the network memory, conveying valuable information across the entire sequence. The gates are specific neural networks that determine which information is permitted on the cell state. Throughout the training, the gates will learn which information is essential to retain or forget. The value of gates and cell state can be determined by using the following (4) to (7):

$$f_t = \sigma(W_f x_t + U_f h_{t-1} + b_f) \quad (4)$$

$$i_t = \sigma(W_i x_t + U_i h_{t-1} + b_i) \quad (5)$$

$$o_t = \sigma(W_o x_t + U_o h_{t-1} + b_o) \quad (6)$$

$$\tilde{C}_t = \tanh(W_c x_t + U_c h_{t-1} + b_c) \quad (7)$$

where  $W_f$ ,  $W_i$ ,  $W_o$ , and  $W_c$  are the weights connecting the hidden layer input to all the gates and input cell state. The  $U_f$ ,  $U_i$ ,  $U_o$ , and  $U_c$  are the weight matrices mapping previous cell output state to all the gates and input cell state. The  $b_f$ ,  $b_i$ ,  $b_o$ , and  $b_c$  are bias vectors. The  $\sigma$  and  $\tanh$  are the sigmoid and tanh activation function, respectively. The cell output state,  $C_t$ , and the layer output,  $h_t$ , at each time iteration  $t$ , can be calculated as in (8)-(9):

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t \quad (8)$$

$$h_t = o_t * \tanh(C_t) \quad (9)$$

The result of the LSTM layer should be a vector of all the outputs, represented as  $Y_T = [h_{T-n}, \dots, h_{T-1}]$ .

2) *Bidirectional LSTM*: The Bi-LSTM are originated from bidirectional RNN [50]. It processes sequential data with two different hidden layers, in both forward and backward directions, and links them to the same output layer. Across certain areas, bidirectional networks are considerably stronger than unidirectional ones, such as speech recognition [51].

Fig. 7 represents an unfolded Bi-LSTM layer structure containing a forward and a backward LSTM layer [52]. The output sequence of the forward layer,  $\vec{h}$ , is determined iteratively using inputs in a definite sequence, while the output sequence of backward layer,  $\overleftarrow{h}$ , is determined using the reversed input. The forward and backward layer outputs are computed using standard LSTM by (4) - (9). The Bi-LSTM layer produces an output vector,  $Y_T$ , which defines each element by the following Equation (10).

$$y_t = \sigma(\vec{h}_t, \overleftarrow{h}_t) \quad (10)$$

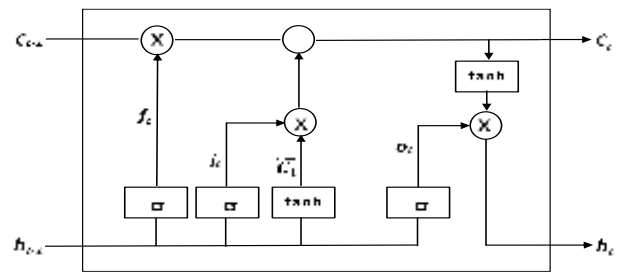


Fig. 6. LSTM Cell [49].

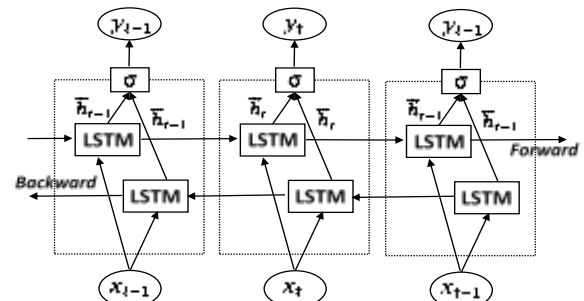


Fig. 7. Structure of an unfolded Bi-LSTM Layer [52].

where  $\sigma$  function combines the two output sequences. It can be a summation function, a multiplication function, a concatenating function, or an average function. Similar to the LSTM layer, a vector,  $Y_T = [h_{T-n}, \dots, h_{T-1}]$ , represents the final output of a Bi-LSTM layer.

F. Bi-LSTM Model Training and Testing

Although LSTM can acquire long speech sequence information but only takes one direction into consideration. It assumes that only previous frame affects the current frame. But not considers that the next frame is also related to current state. This signifies that there is a two-way relationship and the next speech frame should also be considered. Bi-LSTM provides the solution for this problem (Fig. 8).

Bi-LSTM is capable of solving the relationship between two speech frames. It also strengthens the two-way relationship between the current and next speech frame. Due to the bi-directional time structure of Bi-LSTM, it captures more structural information. Hence gives better classification accuracy as compared to one-way LSTM [53].

From Fig. 8, it can be seen that speech features vectors are obtained through the WMFCC feature extraction technique, and then the feature sequences are passed through Bi-LSTM for training and testing. The Bi-LSTM links the output of the feature extraction module to the further layers. Table IV describes the complete training and testing algorithm.

1) Sort data for padding: During training, the training feature vectors are split into mini-batches. The training data is padded so that they all have the same length. However, a large amount of padding degrades network performance. In order to prevent too much padding in the training process, the training data is sorted by sequence length.

2) Define Bi-LSTM network: Bi-LSTM network is a layered architecture shown in Fig. 8. The first layer embedding layer is also called as the sequence input layer. It takes the sorted 14-dimensional WMFCC feature vector as input. The second and third layers are the hidden forward and backward LSTM, forming the Bi-LSTM layer with 100 hidden units. Due to these two layers, the current input is related to the previous and next sequence. The input sequence

reaches the model in both directions through the hidden layer. After the processing of the hidden layers, the outputs are combined to obtain the final output of the Bi-LSTM layer. The output from both the LSTM layers can be computed by the following (11):

$$h_t = \alpha h_t^f + \beta h_t^b \tag{11}$$

where  $h_t^f$  and  $h_t^b$  represents the output of forward and backward LSTM layer, when it takes sequence from  $x_1$  and  $x_T$  as input.  $\alpha$  and  $\beta$  are to control the factors of Bi-LSTM.  $h_t$  is the sum of two unidirectional LSTM elements at time  $t$ .

The output of the Bi-LSTM layer is the input to the fully connected layer of size equal to the number of classes, i.e., five. This layer links each piece of input feature information with a piece of output information for classification by the next layers.

Finally, the softmax and classification layers categorize speech frames into various disfluencies classes such as prolongation, syllable repetition, word repetition, and phrase repetition. The softmax layer applies the softmax function as an activation function that converts the real vector values into a vector with values between 0 and 1, so it can be interpreted as probabilities. The probability of classifying  $x$  into class  $k$  in the softmax regression [54] can be defined by (12).

$$P(y^{(i)} = k|x; \theta) = \frac{\exp(\theta^{(k)T}x)}{\sum_{j=1}^K \exp(\theta^{(j)T}x)} \tag{12}$$

where  $K$  represents the number of classes and  $\theta$  are the model parameters.

In the classification layer, the model receives the values from the softmax function and assigns each input to one of the classes using the cross-entropy function (13).

$$loss = -\sum_{i=1}^N \sum_{j=1}^K t_{ij} \ln y_{ij} \tag{13}$$

where  $N$  represents the number of samples,  $K$  is the number of classes,  $t_{ij}$  indicates that  $i$ th sample belongs to  $j$ th class and  $y_{ij}$  represents the value obtained from the softmax function.

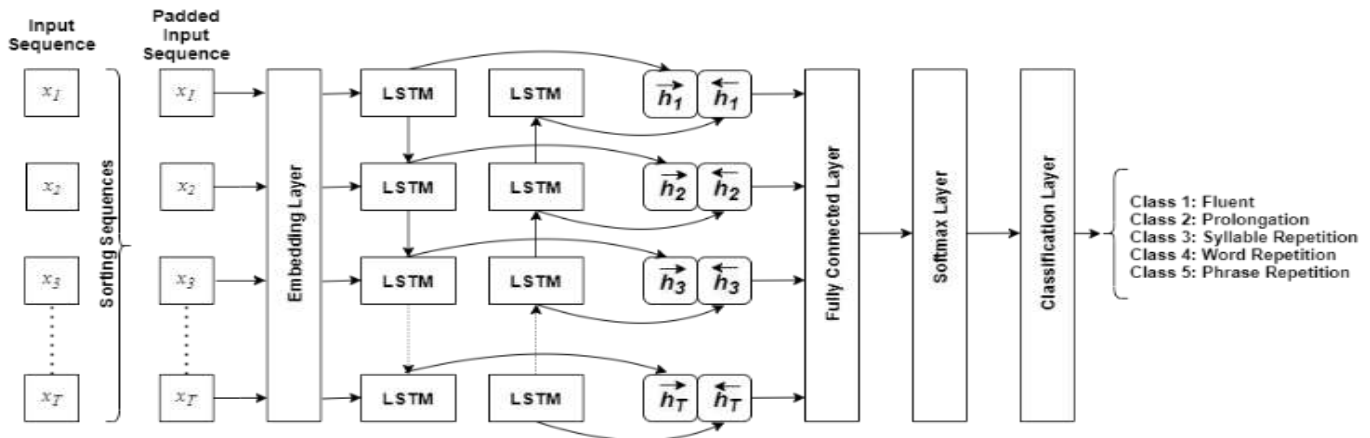


Fig. 8. Bi-LSTM Network.

TABLE IV. ALGORITHM OF BI-LSTM CLASSIFICATION

|   |
|---|
| <p><b>Input:</b> WMFCC feature vector of <math>T_{train}</math>, <math>T_{valid}</math>, and <math>T_{test}</math> datasets.<br/> <b>Output:</b> Stuttered events classification accuracy</p> <ol style="list-style-type: none"> <li>1. Load the training <math>T_{train}</math> and validation <math>T_{valid}</math> dataset.</li> <li>2. Sort the datasets by sequence length.</li> <li>3. Build the Bi-LSTM network.</li> <li>4. Initialize the Bi-LSTM training hyper-parameters.</li> <li>5. Specify the training options.</li> <li>6. Train the Bi-LSTM network with <math>T_{train}</math> dataset.</li> <li>7. Validate the Bi-LSTM network with <math>T_{valid}</math> dataset.</li> <li>8. If Bi-LSTM network is not optimized<br/>then reinitialize the hyperparameters from step 4.</li> <li>9. Load the testing dataset <math>T_{test}</math>.</li> <li>10. Classify the <math>T_{test}</math> samples using a trained Bi-LSTM model.</li> <li>11. Match the similarity between the test labels and predicted labels.</li> <li>12. Evaluate the stuttered events classification accuracy of the model.</li> <li>13. If classification accuracy is optimal<br/>then output classification accuracy<br/>else rebuild the model from step 3</li> </ol> |
|---|

3) *Initializing the hyper-parameters of the network:* Once the network is defined, the hyper-parameters of the network are initialized. Model hyper-parameters are properties on which the entire training process depends [55]. They are divided into two categories: Optimizer and model-specific hyper-parameters. The optimization parameters determine how the network is trained and is more related to optimization, such as the number of epochs, batch size, and learning rate. In contrast, the model-specific parameters are variables that determine the model structure, such as the number of hidden units and hidden layers. These parameters should be defined before training.

Hyper-parameter directly controls the training algorithm's behavior and thus have a significant difference in improving model performance [55]. Therefore, choosing appropriate parameters is an integral part of the optimization of the learned model. The process of selecting good hyper-parameters involves a large number of experiments, which is a time-consuming and tedious task. Most researchers rely on their experience of selecting appropriate parameters for a deep neural network.

In order to determine appropriate hyper-parameters, the classification accuracy of the validation set is used for evaluation. This work applies a diagnostic approach, in which various hyperparameters performance is investigated on both training and validation datasets. The analysis determines how a given configuration performs and how to be adjusted to obtain better performance. The hyper-parameters such as learning rate, batch size, number of epochs, and number of hidden units are taken into consideration for analysis.

4) *Training and testing of the datasets:* Once the Bi-LSTM model and its hyper-parameters are defined, the model is trained by using the training dataset. After the training process is over, the model is validated through the validation dataset. If the classification accuracy of the model is optimized, then its performance is tested; otherwise, the model

hyper-parameters are reconfigured. The parameters such as learning rate, batch size, number of epochs, and number of hidden units are considered for reconfiguration. These parameters are tested for various ranges of values. The process of reconfiguration of hyper-parameters is repeated until the model is optimized, as represented in Fig. 9.

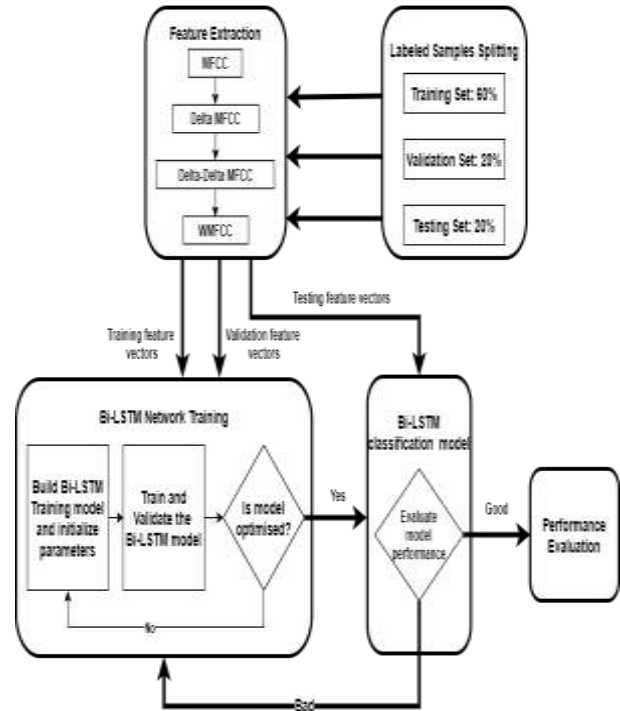


Fig. 9. Bi-LSTM Training and Testing Process.

The classification performance of the optimized model is compared with the traditional LSTM model using the testing dataset. After the testing process is over, the process of performance evaluation is carried out. If the results of the evaluation are optimal, then the process is stopped; otherwise, the complete model is redefined, and the complete training process is repeated until the model is optimized.

#### IV. EXPERIMENTS AND RESULTS

This section discusses the efficacy and performance of the proposed algorithm based on WMFCC feature extraction and Bi-LSTM classification for four forms of disfluencies. This study evaluates the stuttered events recognition model using the stuttered samples obtained from the UCLASS database. The dataset used in this work refers to 20 samples of speech from UCLASS for experimentation. It comprises two female speakers and 18 male speakers aged 7 years 8 months to 17 years 9 months. The stuttered speech samples are manually identified and segmented from the selected speech samples. The segmented samples were labeled as five classes, namely, Fluent, Prolongation, Syllable Repetition, Word Repetition, and Phrase Repetition. The speech samples were split into training testing and validation datasets. Firstly, the signals are pre-processed by removing the silent regions from the samples using the combination of STE and ZCR techniques. Then 14-dimensional acoustic features were extracted from the

segmented samples using the WMFCC feature extraction algorithm. Finally, the extracted feature vectors are inputted to the deep learning Bi-LSTM model. The Bi-LSTM model is trained and optimized through training and validation sets by reconfiguring the hyperparameters. The performance of the proposed model is compared with the traditional LSTM model by using the test set.

**A. Adjustments of Parameters**

In the training model, various hyperparameters of deep learning classification such as learning rate, batch size, number of epochs, and number of hidden units, also play a vital role in the performance of the learned model.

When training the Bi-LSTM network, these parameters are tuned, and their accuracy on the validation set is observed. The experiments were performed based on the hyperparameters' configuration tabled in Table V.

For the first experiment, the best value of the initial learning rate was determined while fixing the typical values for mini batch-size as 16, the number of epochs as 100, and the number of hidden units as 100. The learning rate was varied from  $10^{-2}$  to  $10^{-4}$  for analysis, and the result is presented in Fig. 10. It can be seen that  $10^{-2}$  as the initial learning rate, generated better classification accuracy of 86.67% for available stuttered data.

In the second experiment, the effect of batch-size values 4,8,16 and 32 was determined by fixing the initial learning rate to the best value obtained in the last experiment while the other two with their typical values. The average classification accuracy versus batch size is represented in Fig. 11. The experiment showed that the model produced the highest classification accuracy of 96.67% for the value of mini batch-size as 8.

The effect of the number of epochs was analyzed in the third observational study by fixing the learning rate and batch size as their best values while the typical value for the number of hidden units. The study discussed the effect of different values of epochs, such as 5, 10, 30, 50, and 100. The results are presented in Fig. 12. It can be figured out that number of epochs as 50 outputs best recognition accuracy of speech disfluencies with a value of 96.67%.

Finally, the last experiment was carried out to determine the effect of the various number of hidden units by using the best parameters obtained from the last three experiments. The number of hidden units was varied from 50 to 200 for analysis, and the result is presented in Fig. 13. It can be seen that hidden units as 100 generated better classification accuracy of 96.67%.

TABLE V. EXPERIMENTS OF HYPER-PARAMETERS CONFIGURATION

| Experiments   | Learning Rate          | Batch-Size | Epochs   | Hidden Units |
|---------------|------------------------|------------|----------|--------------|
| Learning Rate | $10^{-2}$ to $10^{-4}$ | 16         | 100      | 100          |
| Batch-Size    | $10^{-2}$              | 4 to 16    | 100      | 100          |
| Epochs        | $10^{-2}$              | 8          | 5 to 100 | 100          |
| Hidden Units  | $10^{-2}$              | 8          | 50       | 100          |

From the experiments, it was determined that the optimal value for learning rate, batch size, number of epochs, and number of hidden units was  $10^{-2}$ , 8, 50, and 100, respectively.

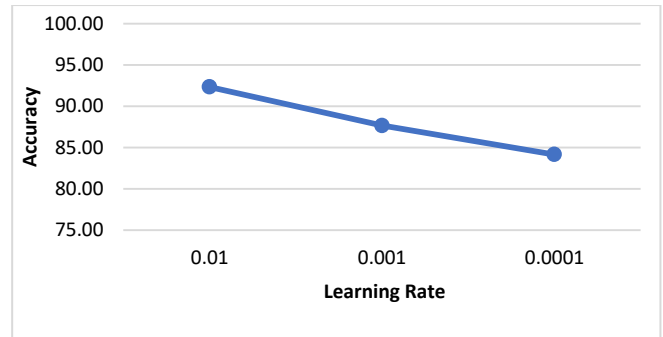


Fig. 10. Changes between Learning Rates and Accuracy.

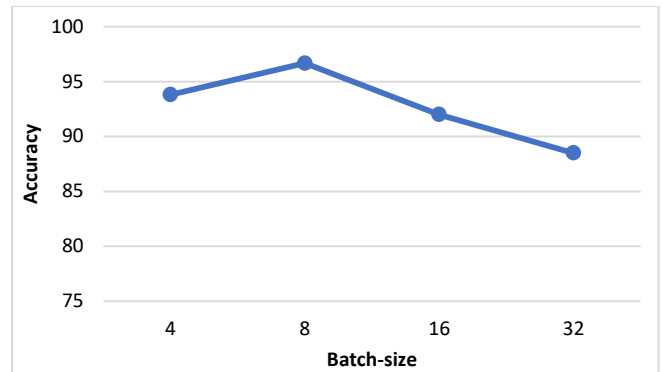


Fig. 11. Changes between Batch Size and Accuracy.

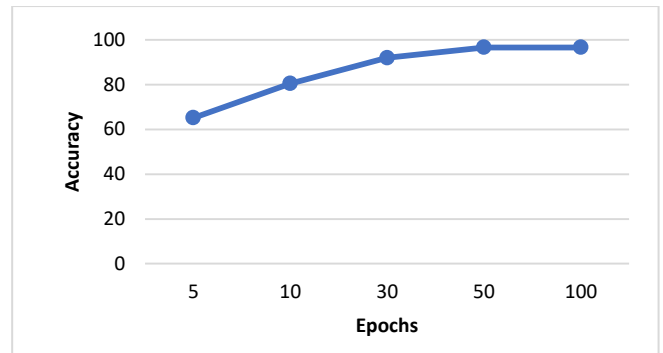


Fig. 12. Changes between Epochs and Accuracy.

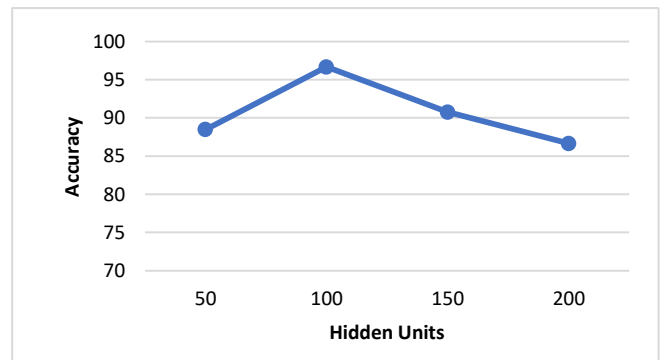


Fig. 13. Changes between the Hidden units and Accuracy.

B. Analysis of Experimental Results

The classification efficiency of the proposed WMFCC and Bi-LSTM based model is verified by carrying out the comparison experiments of the proposed model and unidirectional LSTM. During the experiment, the dimension of the WMFCC feature vector was 14, the frame length was 30ms with overlapping of 75%, the pre-emphasis factor alpha was 0.98, the single Bi-LSTM layer with 100 hidden units, the activation function was Adam, the epochs was 50, the batch-size was 8, and the learning rate was set to  $10^{-2}$ .

The accuracy and loss function of Bi-LSTM and LSTM is represented in Fig. 14 and 15. From Fig. 14, it can be observed that the Bi-LSTM model has slow convergence speed and high accuracy as compared to the LSTM model. From Fig.15, it can be seen that the Bi-LSTM model decreases the loss value to a shallow stable value as compared to LSTM. Thus, it is concluded that the proposed model accomplished a stronger convergence effect.

The complete illustration of the validity of the proposed model can be performed by using the evaluation indicators of relevant experiments such as precision, recall, specificity, and F measure according to the confusion matrix, on test datasets.

The comparison of the LSTM model and the proposed Bi-LSTM model is displayed in Table VI. The results elucidated that WMFCC and Bi-LSTM based model proposed in this work provides the best and efficient performance and the average overall classification accuracy as 96.67%.

Table VII displays the accuracy, sensitivity and specificity of various disfluency classes. In terms of detecting stuttered events, prolongation detection, and phrase detection displayed the highest sensitivity of 97.5%. Classification of word repetition samples gave the best specificity of 99.37%. The prolongation detection achieved the highest accuracy of 98.67%.

From the analysis of the above results, it is concluded that the proposed model performs better than other models, thus determining the effectiveness of long term and bidirectional dependence on information for stuttered speech analysis. Further, the feature extraction of WMFCC includes the dynamic information of the speech samples, which increases the detection accuracy of stuttered events; and also reduces the computational overhead to the classification stage.

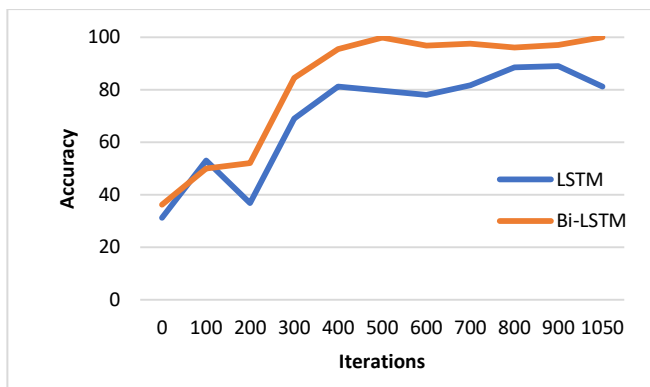


Fig. 14. Accuracy Comparison of LSTM and Bi-LSTM Models.

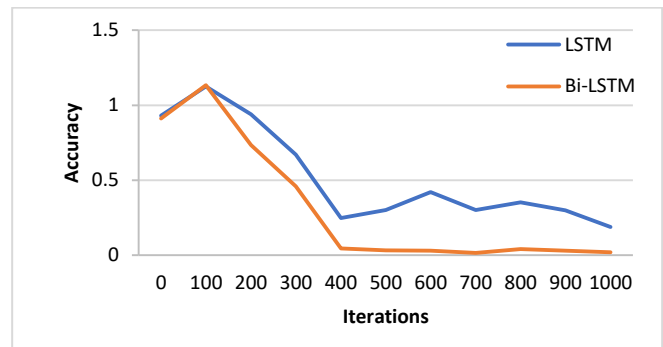


Fig. 15. Loss Comparison of LSTM and Bi-LSTM Models.

TABLE VI. CLASSIFICATION RESULTS OF DISFLUENCY CLASSES

| Disfluency type     | Accuracy (%) | Sensitivity (%) | Specificity (%) |
|---------------------|--------------|-----------------|-----------------|
| Fluent              | 97.33        | 90              | 98.75           |
| Prolongation        | 98.67        | 97.5            | 98.12           |
| Syllable Repetition | 97.5         | 92.5            | 98.12           |
| Word Repetition     | 97.19        | 87.5            | 99.37           |
| Phrase repetition   | 97.67        | 97.5            | 96.87           |

TABLE VII. COMPARISON OF LSTM MODEL AND PROPOSED MODEL

| Model   | Precision (%) | Recall (%) | F-Score (%) | Accuracy (%) |
|---------|---------------|------------|-------------|--------------|
| LSTM    | 83.67         | 84.11      | 83.88       | 83.33        |
| Bi-LSTM | 96.18         | 96.31      | 96.01       | 96.67        |

The result summary of this study (Table VII) and previous works results in Table I give comparable results. However, a direct comparison cannot be made due to different languages, different classifiers, and different types, size, and categorical distribution of stuttered speech database, as well as ways of segmentation of database for gathering, stuttered speech samples.

V. CONCLUSION

The present research proposed an automated and efficient method based on the WMFCC feature extraction algorithm and deep-learning Bi-LSTM network for automatic assessment of the stuttered speech. The disfluencies such as prolongation and syllable, word, and phrase repetition are accurately detectable using this method. The speech samples are parameterized into 14-dimensional WMFCC feature vectors. This model can extract static as well as dynamic acoustic features by using WMFCC, which enhances the detection accuracy of stuttered events; and also reduces the computational overhead to the classification stage. The feature vectors are modeled by Bi-LSTM in both forward and backward directions and capable of learning the long dependencies, taking full account of disfluency patterns in speech frames. Experiments show that when the hyper-parameters are reconfigured during the training of the model, results in an optimal configuration of parameters and leads to a highly accurate model. The optimally configured model proposed in this study is compared with the unidirectional

LSTM model. The disfluency classification accuracy of the proposed model has a better classification accuracy of 96.67% than the LSTM model. It can be concluded that the WMFCC and Bi-LSTM based proposed model effectively improves the recognition accuracy of stuttered events.

In the future study, other feature extraction and classification techniques may be applied for improving the process of detection of speech disfluencies.

#### REFERENCES

- [1] Guitar, *Stuttering: an integrated approach to its nature and treatment*, Fifth edition. Philadelphia; Baltimore: Wolters Kluwer; Lippincott Williams & Wilkins, 2019.
- [2] W. JOHNSON, "Measurements of oral reading and speaking rate and disfluency of adult male and female stutterers and nonstutterers.," *J. Speech Hear. Disord.*, vol. (Suppl 7), pp. 1–20, Jun. 1961.
- [3] S. Erickson and S. Block, "The social and communication impact of stuttering on adolescents and their families.," *J. Fluency Disord.*, vol. 38, no. 4, pp. 311–324, Dec. 2013.
- [4] O. BLOODSTEIN and N. B. RATNER, *A handbook on stuttering*, 6th ed. NY: Thomson Delmar Learning, 2008.
- [5] S. Gupta, R. S. Shukla, and R. K. Shukla, "Literature survey and review of techniques used for automatic assessment of Stuttered Speech.," *Int. J. Manag. Technol. Eng.*, vol. IX, no. X, pp. 229–240, 2019.
- [6] T. S. Tan, Helbin-Liboh, A. K. Ariff, C. M. Ting, and S. H. Salleh, "Application of Malay speech technology in Malay speech therapy assistance tools.," 2007 Int. Conf. Intell. Adv. Syst. ICIAS 2007, pp. 330–334, 2007.
- [7] I. Świetlicka, W. Kuniszyk-Józkowiak, and E. Smółka, "Artificial neural networks in the disabled speech analysis.," *Adv. Intell. Soft Comput.*, vol. 57, pp. 347–354, 2009.
- [8] K. M. Ravikumar, R. Rajagopal, and H. C. Nagaraj, "An Approach for Objective Assessment of Stuttered Speech Using MFCC Features.," vol. 9, no. 1, pp. 19–24, 2009.
- [9] Thiang and Wanto, "Speech Recognition Using LPC and HMM Applied for Controlling Movement of Mobile Robot.," *Semin. Nas. Teknol. Inf.* 2010, no. Seminar Nasional Teknologi Informasi 2010 SPEECH, 2010.
- [10] L. S. Chee, O. C. Ai, M. Hariharan, and S. Yaacob, "MFCC based recognition of repetitions and prolongations in stuttered speech using k-NN and LDA.," SCORed2009 - Proc. 2009 IEEE Student Conf. Res. Dev., pp. 146–149, 2009.
- [11] L. S. Chee, O. C. Ai, M. Hariharan, and S. Yaacob, "Automatic detection of prolongations and repetitions using LPCC.," *Int. Conf. Tech. Postgraduates 2009, TECHPOS 2009*, 2009.
- [12] Ravi Kumar K M and S. Ganesan, "Comparison of Multidimensional MFCC Feature Vectors for Objective Assessment of Stuttered Disfluencies.," *Int. J. Adv. Netw. Appl.*, vol. 860, pp. 854–860, 2011.
- [13] O. Chia Ai, M. Hariharan, S. Yaacob, and L. Sin Chee, "Classification of speech dysfluencies with MFCC and LPCC features.," *Expert Syst. Appl.*, vol. 39, no. 2, pp. 2157–2165, 2012.
- [14] M. Hariharan, V. Vijejan, C. Y. Fook, and S. Yaacob, "Speech stuttering assessment using sample entropy and Least Square Support Vector Machine.," *Proc. - 2012 IEEE 8th Int. Colloq. Signal Process. Its Appl. CSPA 2012*, pp. 240–245, 2012.
- [15] C. Y. Fook, H. Muthusamy, L. S. Chee, S. Bin Yaacob, and A. H. B. Adom, "Comparison of speech parameterization techniques for the classification of speech disfluencies.," *Turkish J. Electr. Eng. Comput. Sci.*, vol. 21, no. SUPPL. 1, pp. 1983–1994, 2013.
- [16] I. Świetlicka, W. Kuniszyk-Józkowiak, and E. Smółka, "Hierarchical ANN system for stuttering identification.," *Comput. Speech Lang.*, vol. 27, no. 1, pp. 228–242, 2013.
- [17] J. Pálffy, "Analysis of Dysfluencies by Computational Intelligence.," *Information Sci. Technol. Bull. ACM Slovakia*, vol. 6, no. 2, pp. 45–58, 2014.
- [18] S. Jabeen and K. M. Ravikumar, "Analysis of 0dB and 10dB babble noise on stuttered speech.," *Proc. IEEE Int. Conf. Soft-Computing Netw. Secur. ICSNS 2015*, pp. 0–4, 2015.
- [19] P. S. Savin, P. B. Ramteke, and S. G. Koolagudi, "Recognition of repetition and prolongation in stuttered speech using ANN.," in *Smart Innovation, Systems and Technologies*, 2016, vol. 43, pp. 65–71.
- [20] P. B. Ramteke, S. G. Koolagudi, and F. Afroz, "Repetition detection in stuttered speech.," in *Smart Innovation, Systems and Technologies*, 2016, vol. 43, pp. 611–617.
- [21] P. Mahesha and D. S. Vinod, "Automatic segmentation and classification of dysfluencies in stuttering speech.," *ACM Int. Conf. Proceeding Ser.*, vol. 04-05-Marc, 2016.
- [22] I. Esmaili, N. J. Dabanloo, and M. Vali, "Automatic classification of speech dysfluencies in continuous speech based on similarity measures and morphological image processing tools.," *Biomed. Signal Process. Control*, vol. 23, pp. 104–114, 2016.
- [23] S. Ghonem, S. Abdou, M. Esmael, and N. Ghamry, "Classification of Stuttering Events Using I-Vector.," *Egypt. J. Lang. Eng.*, vol. 4, no. 1, pp. 11–19, Apr. 2017.
- [24] G. Bhatia, B. Saha, M. Khamkar, A. Chandwani, and R. Khot, "Stutter Diagnosis and Therapy System Based on Deep Learning.," *Jul. 2020*, Accessed: Aug. 16, 2020.
- [25] Girirajan S, Sangeetha R, Preethi T, and Chinnappa A, "Automatic Speech Recognition with Stuttering Speech Removal using Long Short-Term Memory (LSTM).," *Int. J. Recent Technol. Eng.*, vol. 8, no. 5, pp. 1677–1681, Jan. 2020.
- [26] A. B. Nassif, I. Shahin, I. Attili, M. Azzeh, and K. Shaalan, "Speech Recognition Using Deep Neural Networks: A Systematic Review.," *IEEE Access*, vol. 7, pp. 19143–19165, 2019.
- [27] P. J. Lou, P. Anderson, and M. Johnson, "Disfluency detection using auto-correlational neural networks.," in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, EMNLP 2018*, 2020, pp. 4610–4619.
- [28] M. Reisser, "Recurrent Neural Networks in Speech Disfluency Detection and Punctuation Prediction.," 2015.
- [29] Y. Cho and L. K. Saul, "Kernel Methods for Deep Learning.," *Adv. neural Inf. Process. Syst.*, pp. 342–350, 2009.
- [30] G. Zhong, X. Ling, and L. N. Wang, "From shallow feature learning to deep learning: Benefits from the width and depth of deep architectures.," *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, vol. 9, no. 1. Wiley-Blackwell, Jan. 01, 2019.
- [31] C. Chakraborty and P. H. Talukdar, "Issues and Limitations of HMM in Speech Processing: A Survey.," *Int. J. Comput. Appl.*, vol. 141, no. 7, pp. 13–17, May 2016.
- [32] C. P. Lim, S. C. Woo, A. S. Loh, and R. Osman, "Speech recognition using artificial neural networks.," in *Proceedings of the 1st International Conference on Web Information Systems Engineering, WISE 2000*, 2000, vol. 1, pp. 419–423.
- [33] S. Hochreiter, Y. Bengio, P. Frasconi, and J. Schmidhuber, "Gradient Flow in Recurrent Nets: the Difficulty of Learning Long-Term Dependencies.," 2001.
- [34] Y. Bengio, P. Simard, and P. Frasconi, "Learning Long-Term Dependencies with Gradient Descent is Difficult.," *IEEE Trans. Neural Networks*, vol. 5, no. 2, pp. 157–166, Mar. 1994.
- [35] R. Pascanu, T. Mikolov, and Y. Bengio, "On the difficulty of training recurrent neural networks.," 2013.
- [36] C. Olah, "Understanding LSTM Networks.," 2015, Accessed: Aug. 16, 2020.
- [37] S. Hochreiter, "The vanishing gradient problem during learning recurrent neural nets and problem solutions.," *Int. J. Uncertainty, Fuzziness Knowledge-Based Syst.*, vol. 6, no. 2, pp. 107–116, Apr. 1998.
- [38] A. Graves and J. Schmidhuber, "Framewise phoneme classification with bidirectional LSTM networks.," in *Proceedings of the International Joint Conference on Neural Networks*, 2005, vol. 4, pp. 2047–2052.
- [39] P. Howell, S. Davis, and J. Bartrip, "Europe PMC Funders Group The UCLASS archive of stuttered speech.," vol. 52, no. 2, pp. 556–569, 2010.

- [40] Rabiner, L. R. and B. H. Juang, Fundamentals of speech recognition. Prentice-Hall, Inc., USA., 1993.
- [41] L. Rabiner, Digital processing of speech signals. Englewood Cliffs N.J.: Prentice-Hall, 1978.
- [42] D. Enqing, L. Guizhong, Z. Yatong, and C. Yu, "Voice activity detection based on short-time energy and noise spectrum adaptation," Int. Conf. Signal Process. Proceedings, ICSP, vol. 1, no. 1, pp. 464–467, 2002.
- [43] R. G. Bachu, S. Kopparthi, B. Adapa, and B. D. Barkana, "Voiced/unvoiced decision for speech signals based on zero-crossing rate and energy," Adv. Tech. Comput. Sci. Softw. Eng., pp. 279–282, 2010.
- [44] P. Howell and M. Huckvale, "Facilities to assist people to research into stammered speech.," Stammering Res., vol. 1, no. 2, pp. 130–242, Jul. 2004, Accessed: Jul. 19, 2020.
- [45] S. Gupta, R. S. Shukla, R. K. Singh, and R. K. Shukla, "Weighted Mel Frequency Cepstral Coefficient based feature extraction for automatic assessment of stuttered speech using Bi-directional Long-Short Term Memory", unpublished.
- [46] X. Huang, Spoken language processing : a guide to theory, algorithm and system development. Upper Saddle River N.J.: Prentice Hall PTR, 2001.
- [47] S. V.Chapaneri, "Spoken Digits Recognition using Weighted MFCC and Improved Features for Dynamic Time Warping," Int. J. Comput. Appl., vol. 40, no. 3, pp. 6–12, 2012.
- [48] S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," Neural Comput., vol. 9, no. 8, pp. 1735–1780, 1997.
- [49] M. Phi, "Illustrated Guide to LSTM's and GRU's: A step by step explanation | by Michael Phi | Towards Data Science," 2018. <https://towardsdatascience.com/illustrated-guide-to-lstms-and-gru-s-a-step-by-step-explanation-44e9eb85bf21>.
- [50] M. Schuster and K. K. Paliwal, "Bidirectional recurrent neural networks," IEEE Trans. Signal Process., vol. 45, no. 11, pp. 2673–2681, 1997.
- [51] A. Graves, N. Jaitly, and A. Mohamed, "Hybrid speech recognition with deep bidirectional LSTM Alex Graves , Navdeep Jaitly and Abdelrahman Mohamed University of Toronto Department of Computer Science 6 King ' s College Rd . Toronto , M5S 3G4 , Canada," pp. 273–278, 2013.
- [52] Z. Cui, R. Ke, Z. Pu, and Y. Wang, "Stacked bidirectional and unidirectional LSTM recurrent neural network for forecasting network-wide traffic state with missing values," Transp. Res. Part C Emerg. Technol., vol. 118, Sep. 2020.
- [53] S. Siami-Namini, N. Tavakoli, and A. S. Namin, "The Performance of LSTM and BiLSTM in Forecasting Time Series," in Proceedings - 2019 IEEE International Conference on Big Data, Big Data 2019, Dec. 2019, pp. 3285–3292.
- [54] Y. Ren, P. Zhao, Y. Sheng, D. Yao, and Z. Xu, "Robust Softmax Regression for Multi-class Classification with Self-Paced Learning," 2017.
- [55] J. Leonel, "Hyperparameters in Machine /Deep Learning | by Jorge Leonel | Medium," Apr. 2019. <https://medium.com/@jorgesleonel/hyperparameters-in-machine-deep-learning-ca69ad10b981>.

# Using of Redundant Signed-Digit Numeral System for Accelerating and Improving the Accuracy of Computer Floating-Point Calculations

Otsokov Sh.A<sup>1</sup>

Dept. of Computing Machines, Systems and Networks  
National Research University "Moscow Power Engineering  
Institute", Moscow, Russian Federation

Magomedov Sh.G<sup>2</sup>

Dept. of Intelligent Information Security Systems  
MIREA Russian Technological University  
Moscow, Russian Federation

**Abstract**—The article proposes a method for software implementation of floating-point computations on a graphics processing unit (GPU) with an increased accuracy, which eliminates sharp increase in rounding errors when performing arithmetic operations of addition, subtraction or multiplication with numbers that are significantly different from each other in magnitude. The method is based on the representation of floating-point numbers in the form of decimal fractions that have uniform distribution within a range and the use of redundant signed-digit numeral system to speed up calculations. The results of computational experiments for evaluating the effectiveness of the proposed approach are presented. The effect of accelerating computations is obtained for the problems of calculating the sum of an array of numbers and determining the dot product of vectors. The proposed approach is also applicable to the discrete Fourier transform.

**Keywords**—High-precision computation; redundant signed-digit numeral system; signed-digit floating-point format; redundant signed-digit arithmetic; decimal fractions

## I. INTRODUCTION

Most computer calculations are carried out in floating-point format and double precision computer calculations are sufficient for solving many computational problems.

However, there are a number of problems, for example, in computational geometry and other areas where double precision floating-point arithmetic is not sufficient [1]. To solve such problems, the well-known libraries of high-precision computations are used, such as ZREAL (Russia), MPARITH (Germany), GMP (USA), which implement floating-point calculations at the software level with a mantissa length set by the user [2, 3, 6, 7, 8, 9, 10].

But these libraries have the property of sharply increasing the calculation time with the increasing in the length of the mantissa and the number of arithmetic operations. In addition, they have the inherent disadvantages of the floating-point format itself, which does not always guarantee an accurate result of computer calculations.

One of such disadvantages is the uneven distribution of floating-point numbers. Fig. 1 below shows the uneven distribution of normalized floating-point numbers with the

mantissa length of 3 binary digits and the order from 0 to 4 [4].

As an example of the loss of accuracy in computer calculations consider the problem of determining the dot product of two vectors with following coordinates:

$$\vec{x} = (10^\alpha, 1223, 10^{\alpha-1}, 10^{\alpha-2}, 3, -10^{\alpha-5}),$$
$$\vec{y} = (10^\beta, 2, -10^{\beta+1}, 10^\beta, 2111, 10^{\beta+3}),$$

where  $\alpha, \beta$  are given parameters.

The true result is 8779. The dot product was calculated in the single precision format, the relative error was calculated with the constant value of  $\alpha = 1$  and  $\beta$  ranging from 1 to 21.

The relative error of the dot product was calculated using the formula:

$$\delta = \frac{|(x, y) - 8779|}{8779} \cdot 100 \quad (1)$$

The dependence of the relative error of the dot product of vectors on the parameter  $\beta$  in single precision floating-point format is presented in the graph shown in Fig. 2.

Fig. 2 shows that for significantly different values of  $\alpha$  and  $\beta$ , starting from the value 18 for  $\beta$ , there is a sharp loss in the accuracy of the dot product results, which is due to the fact that calculations are performed with numbers that differ greatly from each other in magnitude.

Using double precision floating-point format, the increase in the relative error occurs at larger values of  $\beta$  compared to the single precision format.





Fig. 1. Distribution of Floating-Point Numbers.

### Relative error, %

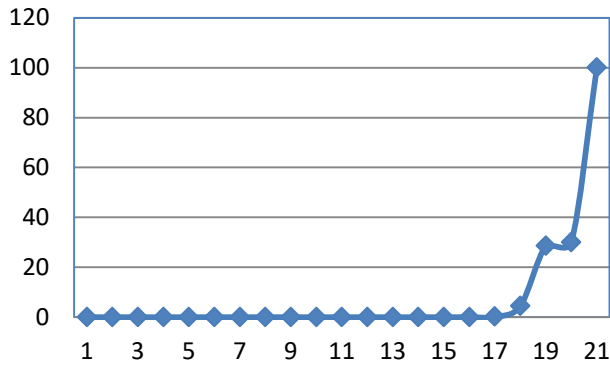


Fig. 2. Dependence of the Relative Error on the Parameter  $\beta$ .

Consider the following example that demonstrates the loss of precision in floating-point calculations associated with the discrete Fourier transform.

Let a vector  $\vec{\mathbf{x}} = (x_1, x_2, x_3, \dots, x_n)$

be given with coordinates defined as follows:

$$x_i = \begin{cases} \frac{i}{10}, & \text{if } i \text{ is odd} \\ \frac{i}{10} + \alpha, & \text{if } i \text{ is even} \end{cases} \quad (2)$$

$$i = 1, \dots, n, \alpha > 0$$

The discrete Fourier transform of a vector  $\vec{\mathbf{x}}$  into a vector  $\vec{\mathbf{y}}$  is performed using the following formula:

$$y_k = \frac{1}{n} \cdot \sum_{j=0}^{n-1} x_j \cdot \left[ \cos \frac{2 \cdot \pi \cdot j \cdot k}{n} - i \cdot \sin \frac{2 \cdot \pi \cdot j \cdot k}{n} \right],$$

$$0 \leq k \leq n-1 \quad (3)$$

The inverse Fourier transform is performed using the following formula:

$$x_k = \sum_{j=0}^{n-1} y_j \cdot \left[ \cos \frac{2 \cdot \pi \cdot j \cdot k}{n} + i \cdot \sin \frac{2 \cdot \pi \cdot j \cdot k}{n} \right],$$

$$0 \leq k \leq n-1 \quad (4)$$

Obviously, if we carry out the direct discrete Fourier transform, then the inverse Fourier transform of the vector.

$$\vec{\mathbf{x}} = (x_1, x_2, x_3, \dots, x_n)$$

we get the vector

$$\vec{\mathbf{x}}^* = (x_1^*, x_2^*, x_3^*, \dots, x_n^*)$$

which should approximately coincide with  $\vec{\mathbf{x}}$ , and the maximum relative error of the transformations can be estimated using the formula:

$$\delta = \max_i \left( \frac{|x_i^* - x_i|}{x_i} \cdot 100 \right) \quad (5)$$

Fig. 3 also demonstrates the loss of accuracy of the floating-point calculations with increasing  $\alpha$ .

The first goal of this work is to eliminate the sharp loss of accuracy in calculations with numbers that differ greatly from each other in magnitude. The second goal is to speed up computations by parallelizing them.

The first goal is achieved by moving from floating-point representation to decimal representation that is evenly distributed within the range, as shown in Fig. 4, for example for decimal fractions of the first degree.

The second goal is achieved through the use of a redundant signed-digit numeral system, in which redundant negative digits are introduced into the system of bases in such a way that the propagation of the carry when adding is not allowed further than one digit [4,5]. Due to this, the arithmetic operations of addition, subtraction and multiplication are parallelized, which leads to their acceleration, especially when the number of digits increases. The time required for addition or subtraction of numbers does not depend on the digit capacity of the numbers.

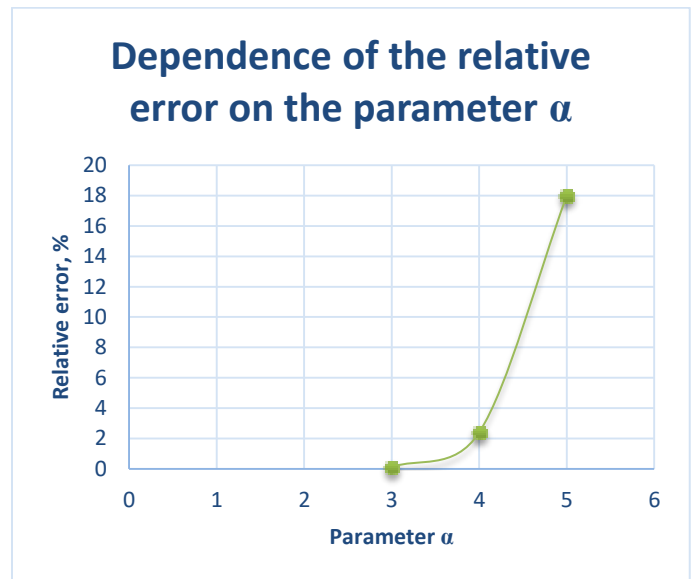


Fig. 3. Dependence of the Relative Error on the Parameter  $\alpha$  using Single Precision Format.

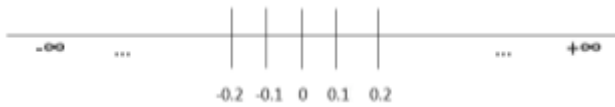


Fig. 4. Even Distribution of Decimal Fractions.

In [11,12,14,15,16] methods of representation and algorithms for performing arithmetic operations in a modular numeral system and a method for their acceleration due to parallelization in several modules are presented.

This article proposes a different approach based on the transition to a redundant signed-digit numeral system.

Such numeral system has an advantage over modular systems in that it is simpler to convert directly and inversely to the traditional numeral system and there is no need for overflow control.

A method for summing a group of numbers and calculating the dot product, oriented towards parallel implementation on a GPU, is considered. The results of experimental studies of the effectiveness of this method of high-precision calculations are obtained.

The next section considers a possible way to represent numbers in a redundant signed-digit numeral system.

## II. REPRESENTING NUMBERS IN A REDUNDANT SIGNED-DIGIT NUMERAL SYSTEM

Consider the representation of floating-point numbers of the following form [4]:

$$A = K \cdot q^t \tag{6}$$

where

$A$  is a floating-point number,

$K$  is the mantissa of the number  $A$ , an integer such that satisfies the inequality

$$|K| \leq q^{n_f} - 1,$$

$q$  is the base (radix) of the numeral system,

$t$  is the order, an integer such that satisfies the inequality

$$|t| \leq k_f$$

$n_f$  is a natural number characterizing the length of the mantissa of the floating-point number,

$k_f$  is a natural number characterizing the maximum order of representable numbers.

Table I includes positive and negative minimum and maximum numbers representable in the form (6) [11,13,17,18].

The range of representable numbers (6) is as follows:

$$\left(-q^{n_f+k_f}, q^{n_f+k_f}\right) \tag{7}$$

Consider the sum of numbers of the form:

$$S = K_1 \cdot q^{t_1} + K_2 \cdot q^{t_2} + \dots + K_n \cdot q^{t_n} \tag{8}$$

Suppose  $t_S = \min(t_1, t_2, \dots, t_n)$  then:

$$\begin{aligned} S &= q^{t_S} \left( K_1 \cdot q^{t_1-t_S} + K_2 \cdot q^{t_2-t_S} + \dots + K_n \cdot q^{t_n-t_S} \right) = \\ &= K_S \cdot q^{t_S} \end{aligned} \tag{9}$$

where

$$K_S = K_1 \cdot q^{t_1-t_S} + \dots + K_n \cdot q^{t_n-t_S}$$

Let us estimate the maximum number of digits required to describe the result of the sum (9). Using Table I, we have:

$$\begin{aligned} K_S &= K_1 \cdot q^{t_1-t_S} + \dots + K_n \cdot q^{t_n-t_S} \\ |K_S| &\leq q^{n_f+k_f} + \dots + q^{n_f+k_f} = n \cdot q^{n_f+k_f} \end{aligned} \tag{10}$$

If  $q = 10$ , the maximum number of digits required to describe the sum will be equal to:

$$\lfloor \lg n + n_f + k_f \rfloor$$

From the last expression it can be seen that to implement the addition of groups of numbers in floating-point format calculations with large numbers are required.

Consider the format for representing floating-point numbers (6) in the signed-digit numeral system as follows:

$$A = [ (\alpha_1, \alpha_2, \dots, \alpha_i, \dots, \alpha_n), t ] \tag{11}$$

where

$\alpha_i$  are the digits of the signed-digit representation.

This format will be referred to as the floating-point signed-digit format.

Consider the rule for adding two numbers and a group of numbers in this format.

TABLE I. MAXIMUM AND MINIMUM POSITIVE AND NEGATIVE NUMBERS

|                  |                 |
|------------------|-----------------|
| Maximum positive | $q^{n_f+k_f}$   |
| Minimum positive | $q^{-k_f-n_f}$  |
| Minimum negative | $-q^{n_f+k_f}$  |
| Maximum negative | $-q^{-k_f-n_f}$ |

### III. RULES FOR PERFORMING ARITHMETIC OPERATIONS IN SIGNED-DIGIT FLOATING-POINT FORMAT

Let there be given two floating-point numbers of the form (11).

$$A_1 = [ (\alpha_1, \alpha_2, \dots, \alpha_i, \dots, \alpha_n), t_1 ]$$

$$A_2 = [ (\beta_1, \beta_2, \dots, \beta_i, \dots, \beta_n), s_1 ]$$

Calculating the sum  $A_1 + A_2$

$$A_4 = A_1 + A_2$$

Suppose  $s_1 > t_1$ . This does not change the generality of reasoning. Then:

$$A_1 \pm A_2 = K_1 \cdot q^{t_1} \pm K_2 \cdot q^{s_1} = q^{t_1} \cdot (K_1 \pm q^{s_1-t_1} \cdot K_2) = [(\delta_1, \delta_2, \dots, \delta_n), t_1] \quad (12)$$

Let  $-6, -5, \dots, 0, \dots, 6$  be the digits of the signed-digit numeral system.

The time required for adding numbers does not depend on the number of digits. The addition is performed using the following steps:

1. Calculation of the sum  $ui = xi + yi - 10 \square i$ , where  $c = \begin{cases} 1, & \text{if } (xi + yi) > 6 \\ -1, & \text{if } (xi + yi) < -6 \\ 0, & \text{otherwise} \end{cases}$
2. Calculation of the final result  $\delta_i = u_i + \square_{i-1}$

Product of numbers  $A_1 \cdot A_2$  is calculated by the formula:

$$A_3 = A_1 \cdot A_2 = [(\chi_1, \chi_2, \dots, \chi_i, \dots, \chi_n), t_1 + s_1] \quad (13)$$

where coefficients  $\chi_1, \chi_2, \dots, \chi_i, \dots, \chi_n$  are determined according to the rules for multiplying numbers in the signed-digit numeral system.

In the next section, the first and second methods of summing a group of numbers are considered.

### IV. METHODS OF SUMMING GROUPS OF NUMBERS

The first method of summing groups of numbers is carried out according to the formula (12) using the addition rule presented in Section III in redundant signed-digit arithmetic in parallel over the digits and sequentially for each number of the group. This method, when implemented on GPU, requires a large number of synchronizations between cores. In Section V

the results of experimental study of the effectiveness of this method are presented.

Consider the second method of summation with fewer synchronizations.

Let the number of digits of the summed numbers equal  $d$ , and the maximum possible number of digits required to describe the sum equal  $max d$ .

Let a set of numbers for summation be given:

$$1^{st} \text{ number: } a_1^1 a_2^1 \dots a_d^1$$

$$2^{nd} \text{ number: } a_1^2 a_2^2 \dots a_d^2,$$

$$k^{th} \text{ number: } a_1^k a_2^k \dots a_d^k$$

Let us denote their sum by  $S$ .

Let us add these numbers to the left with zeros to the maximum possible number of digits  $max d$ , getting the following:

$$\begin{array}{cccccc} 0 & \dots & 0 & a_1^1 & \dots & a_d^1 \\ 0 & \dots & 0 & a_1^2 & \dots & a_d^2 \\ 0 & \dots & 0 & a_1^3 & \dots & a_d^3 \\ 0 & \dots & 0 & a_1^4 & \dots & a_d^4 \\ 0 & \dots & 0 & \dots & \dots & \dots \\ 0 & \dots & 0 & a_1^k & \dots & a_d^k \end{array}$$

Summation is carried out as follows:

1. Each one of  $max d$  threads in parallel and independently calculates the sums of the numbers of the corresponding column, for example, thread number  $max d$  calculates the sum of the numbers:

$$S_{max d} = \sum_{i=1}^k a_d^i$$

thread number  $max d - 1$  calculates the sum of the numbers:

$$S_{max d - 1} = \sum_{i=1}^k a_{d-1}^i$$

2. At the end of the first step, the calculations are synchronized.
3. By definition the value  $S$  equals:

$$S = S_{max d - d} \cdot 10^{d-1} + \dots + S_{max d - 3} \cdot 10^2 + S_{max d - 2} \cdot 10 + S_{max d - 1} \quad (14)$$

Each thread extracts digits from its result, for example, thread number  $max\ d - 1$  finds digits:

$i_1, i_2, i_3$ , thread number  $max\ d - 2$  finds digits:  
 $j_1, j_2, j_3$ .

4. Each thread forms numbers of the form:

$$\begin{matrix} 0 & \dots & 0 & i_1 & i_2 & i_3 \\ 0 & \dots & j_1 & j_2 & j_3 & 0 \\ 0 & \dots & \dots & \dots & \dots & \dots \end{matrix}$$

Thread number  $max\ d - 1$  forms the first line, thread number  $max\ d - 2$  forms the second line etc.

5. These numbers are summed sequentially one after another bit-parallel in the signed-digit numeral system.  
6. The final result is transferred from the GPU to the CPU.

Such summation method requires  $d - 1$  synchronizations in the process of summing this array.

Consider an example.

Suppose  $d = 3$ ,  $max\ d = 5$ ,  $k = 3$ .

Numbers for the summation:

|   |   |   |
|---|---|---|
| 5 | 5 | 6 |
| 1 | 7 | 9 |
| 9 | 7 | 9 |

Calculated sums:

$S_4 = 24$   
 $S_3 = 19$   
 $S_2 = 15$   
 $S_1 = 0$   
 $S_0 = 0$

Each thread forms one of the following numbers:

|   |   |   |   |   |
|---|---|---|---|---|
| 0 | 0 | 0 | 2 | 4 |
| 0 | 0 | 1 | 9 | 0 |
| 0 | 1 | 5 | 0 | 0 |

Then these numbers are summed sequentially one after another bit-parallel in the signed-digit numeral system according to the first method.

Next section considers the results of experimental studies of the efficiency of summation of groups of numbers.

#### V. EXPERIMENTAL STUDY OF THE EFFICIENCY OF HIGH-PRECISION SUMMATION OF GROUPS OF NUMBERS IN THE SIGNED-DIGIT NUMERAL SYSTEM

Numerical experiments were carried out on the addition of groups of integers of different magnitudes, with the number of integers  $k = 10000, 100000, 1000000$ . The addition was carried out according to the rules of traditional arithmetic on the CPU bitwise each number of the group sequentially and using the first method on the Nvidia GPU (1.78 GHz, 1280 cores) bit-parallel and sequentially for each number of the group.

GPU calculations were performed as follows:

- 1) Initial data were generated randomly, integers of fixed length were generated and stored in arrays.
- 2) Arrays were transferred to the GPU.
- 3) A number of threads were created matching the number of digits. Each thread carried out sequential summation of the array numbers in its corresponding digit in parallel and independently.
- 4) The result of the summation was transferred from the GPU to the CPU.

The time required for summation of numbers on the CPU and the GPU was calculated, considering the transfer of data to the GPU and in the opposite direction to the CPU. On the basis of these calculations, the absolute acceleration coefficients were determined for different numbers of digits and values of  $k$  by the formulas:

$$K_{abs} = \frac{T_{cpu}}{T_{gpu}} \quad (15)$$

Fig. 5 shows the dependence of this coefficient on the number of digits.

Fig. 5 shows that the acceleration effect provided by the first method is insignificant and is achieved for numbers exceeding 400 decimal digits.

Experiments showed that data transfer from CPU to GPU and from GPU to CPU was very fast and did not lead to delays in the computation process. One of the main reasons for the low efficiency of the first method of summation on the GPU in the signed-digit numeral system is associated with the need for synchronization after addition of each pair of numbers in the group, which slows down the computation. If the array contains  $k$  numbers, then this summation method requires  $k-1$  synchronizations in the process of summing this array.

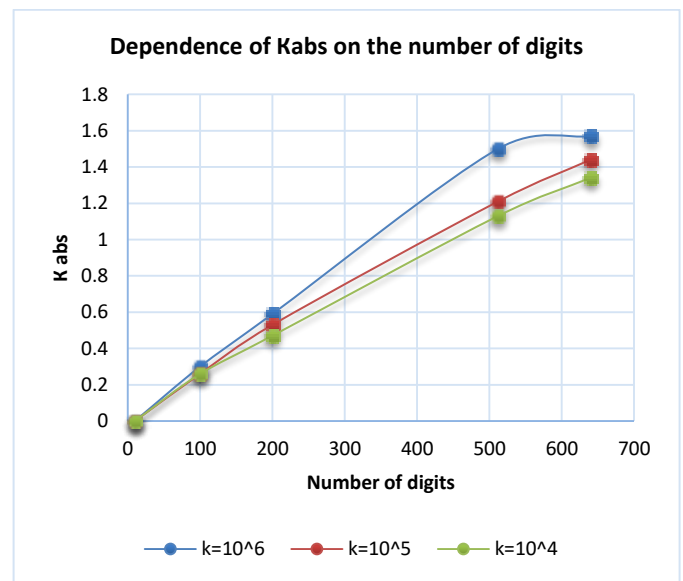


Fig. 5. Dependence of the Absolute Acceleration Coefficient on the Number of Digits for the First Summation Method.

Experiments on summing the same groups of numbers using the second method have shown that it is more efficient than the first method. For the second method the computation times on the CPU and the GPU were calculated, on the basis of which the absolute acceleration coefficient was determined by the formula (15) for different values of  $k$ .

Fig. 6 shows the dependence of the absolute acceleration coefficient on the number of digits for this summation method.

Fig. 6 shows that using the second summation method on the GPU with numbers comprising of 800 to 900 digits speeds up the computation 3 to 4 times in comparison with summation on the CPU. With further increase in value of  $k$  and the number of digits, the acceleration is supposed to be even greater.

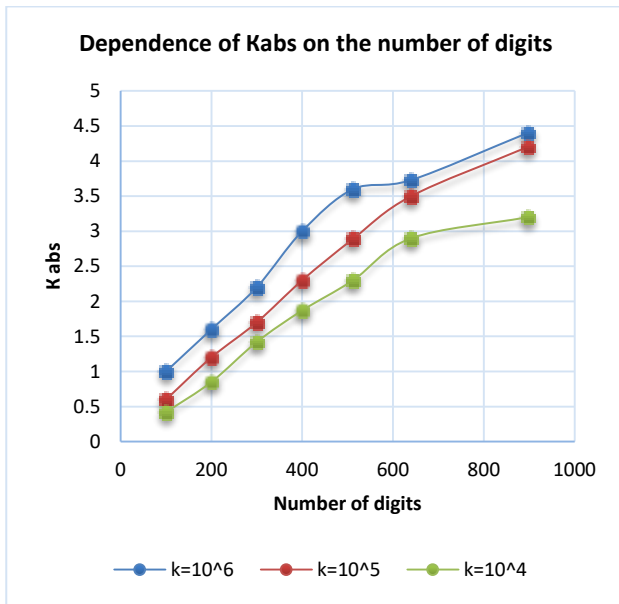


Fig. 6. Dependence of the Absolute Acceleration Coefficient on the Number of Digits for the Second Summation Method.

Next section considers a method for multiplying numbers based on redundant signed-digit arithmetic.

### VI. METHOD FOR MULTIPLYING NUMBERS

Let the number of digits required to represent initial data equal  $d$ , and the maximum possible number of digits required to represent the result of multiplication equal  $max d$ .

Let two numbers be given:

$$1^{st} \text{ number: } a_1^1 a_2^1 \dots a_d^1$$

$$2^{nd} \text{ number: } a_1^2 a_2^2 \dots a_d^2$$

Let us add these numbers to the left with zeros to the maximum possible number of digits  $max d$ , getting the following:

$$\begin{matrix} 0 & \dots & 0 & a_1^1 & \dots & a_d^1 \\ 0 & \dots & 0 & a_1^2 & \dots & a_d^2 \end{matrix}$$

Multiplication process on the GPU involves  $d$  threads.

Thread number 1 forms the results (partial products):

$$0 \quad \dots \quad 0 \quad a_1^1 \cdot a_d^2, a_2^1 \cdot a_d^2, \quad \dots \quad a_d^1 \cdot a_d^2$$

Thread number 2 forms the results:

$$0 \quad \dots \quad a_1^1 \cdot a_{d-1}^2, a_2^1 \cdot a_{d-1}^2, \quad a_d^1 \cdot a_{d-1}^2 \quad 0 \tag{16}$$

Thread number  $d$  forms the results:

$$0 \quad \dots \quad a_1^1 \cdot a_1^2, \dots, a_2^1 \cdot a_1^2, a_d^1 \cdot a_1^2 \quad \dots \quad 0$$

The results are formed in parallel and independently by each thread for each product of two numbers.

Then they are summed up using the second method described in Section IV.

Next section considers the results of experimental studies of the efficiency of calculating the dot product of vectors.

### VII. EXPERIMENTAL STUDY OF THE EFFICIENCY OF HIGH-PRECISION CALCULATION OF THE DOT PRODUCT OF VECTORS IN SIGNED-DIGIT NUMERAL SYSTEM

The dot product of vectors  $(x,y)$ , where  $x = (x_1, x_2, \dots, x_k)$ ,  $y = (y_1, y_2, \dots, y_k)$ , is calculated as follows:

- 1) The values of arrays  $x,y$  are transferred to the GPU.
- 2) For each pair  $x_i$  and  $y_i$  partial products (24-26) are calculated in parallel and independently on the GPU.
- 3) Next, the summation of the obtained partial products is carried out using the second method.
- 4) The result of the dot product is transferred to the CPU.

Numerical experiments were carried out to calculate the dot product of vectors with different numbers of coordinates  $k = 10000, 100000, 1000000$ . The coordinates were integers with the length of 50 decimal digits. The dot product was calculated according to the rules of traditional integer arithmetic on the CPU and in redundant signed-digit arithmetic bit-parallel on the GPU.

The results of the experiments are presented on the graph in Fig. 7.

The graph shows that with the increase of the number of digits required to represent initial data  $d$  and the maximum possible number of digits required to represent the result of the dot product  $max d$  the proposed method provides greater acceleration of the computation process.

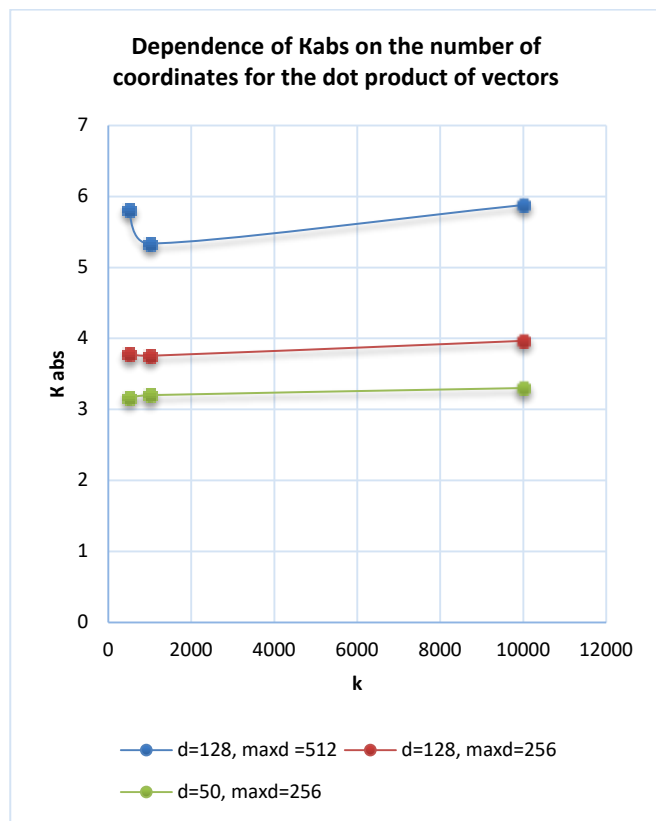


Fig. 7. Dependence of the Absolute Acceleration Coefficient on  $k$  with different Values of  $d, \max d$ .

### VIII. CONCLUSION

This article proposes an approach to software implementation of computations on a GPU, which prevents sharp loss of precision in calculations with numbers that differ greatly from each other in magnitude. The approach is based on representation of floating-point numbers in the form of decimal fractions and the use of a redundant signed-digit numeral system to speed up computations with them on the GPU.

The effect of accelerating computations was obtained and proven experimentally for the operations of summation of an array of numbers on the GPU and calculating the dot product of vectors.

The proposed approach is also applicable for the discrete Fourier transform, for the case presented in the article as well as in other cases.

### REFERENCES

- [1] Bailey D.H., Barrio R., Borwein J.M. High-precision computation: Mathematical physics and dynamics // Applied Mathematics and Computation. 2012. Vol. 218, No. 20. P. 10106-10121.
- [2] D. H. Bailey, J. M. Borwein. "High-Precision Arithmetic in Mathematical Physics", Mathematics, 3 (2015), pp. 337–367.
- [3] David H. Bailey. High-Precision Computation and Mathematical Physics Lawrence Berkeley National Laboratory, 2009
- [4] J.-M. Muller, N. Brisebarre, F. de Dinechin, C.-P. Jeannerod, V. Lefevre, G. Melquiond, N. Revol, D. Stehl'e, S. Torres. Handbook of Floating-Point Arithmetic, Birkhauser, Boston, 2010, 572 p.
- [5] Koren Israel. Computer arithmetic algorithms // University of Massachusetts, 2nd ed. 2002.
- [6] Fousse L., Hanrot G., Lefèvre V., Pélissier P., Zimmermann P. MPFR: a multiple-precision binary floating-point library with correct rounding // ACM Transactions on Mathematical Software. 2007. Vol. 33, No. 2. Article No. 13.
- [7] MParithm - package for high precision computation, 2015, www.wolfgang-ehrhart.de/mp\_intro.html
- [8] GNU Scientific Library 2.5 released — 2018, https://savannah.gnu.org/forum/forum.php?forum\_id=9175.
- [9] Operations with multi-digit real numbers of the ZReal type. http://ishodniki.ru/list/index.php?action=name&show=pascal-math&cat=11.
- [10] D. H. Bailey, X. S. Li and B. Thompson, "ARPREC: An arbitrary precision computation package," Sep 2002, http://crd.lbl.gov/~dhbailey/dhbpapers/arprec.pdf.
- [11] Otsokov Sh. A and Magomedov Sh.G, "On the Possibility of Implementing High-Precision Calculations in Residue Numeral System" International Journal of Advanced Computer Science and Applications(IJACSA), 10(11), 2019. DOI: 10.14569/IJACSA.2019.0101102.
- [12] R. Amos Omondi, Benjamin , "Residue Number Systems: Theory and Implementation," Imperial College Press, 2007.
- [13] Solovyev R.A., Balaka E.S., Telpukhov D.V. Device for calculation of vector dot product with error correction based on residue number system // Problems of the development of prospective micro- and nanoelectronic systems - 2014. Proceedings / edited by A.L. Stempkovskiy. M.: IPPM RAS, 2014. Part IV. pp. 173-178.
- [14] Jen-Shiun Chiang, Mi Lu, «Floating-point numbers in residue number systems» Computers & Mathematics with Applications, vol. 22, issue 10, pp. 127–140, 1991.
- [15] K. S. Isupov, A. N. Mal'tsev. "A parallel-processing-oriented method for the representation of multi-digit floating-point numbers", Vychislitel'nyye metody i programirovaniye, 15:4 (2014), pp. 631–643 (in Russian).
- [16] Magomedov S.G. Increasing the efficiency of microprocessors in an access control systems. International Journal of Engineering and Technology (UAE). 2018. T. 7. № 4.36. C. 80-83.
- [17] Mukunoki D., Ogita T. Performance and energy consumption of accurate and mixed-precision linear algebra kernels on GPUs //Journal of Computational and Applied Mathematics. – 2020. – T. 372. – C. 112701.
- [18] Isupov K., Knyazkov V., Kuvaev A. Design and implementation of multiple-precision BLAS Level 1 functions for graphics processing units //Journal of Parallel and Distributed Computing. – 2020. – T. 140. – C. 25-36.

# Self-Configurable Current-Mirror Technique for Parallel RGB Light-Emitting Diodes (LEDs) Strings

Shaheer Shaida Durrani<sup>1</sup>

Sustainable Energy and Power Electronics Research (SuPER), Faculty of Electrical and Electronics Engineering Technology, University Malaysia Pahang, Pekan, Malaysia

Asif Nawaz<sup>2</sup>

ETS (Electronics) Faculty of Engineering Higher College of Technology, Dubai, UAE

Muhamamd Shahzad<sup>3</sup>, Zeeshan Najam<sup>8</sup>

Department of Electrical Engineering MNS University of Engg and Tech, Multan, Pakistan

Rehan Ali Khan<sup>4</sup>

Department of Electrical Engineering University of Science and Technology, Bannu, Pakistan

Abu Zaharin Ahmad<sup>5</sup>

Sustainable Energy and Power Electronics Research (SuPER) Faculty of Electrical and Electronics Engineering Technology University Malaysia Pahang, Pekan Malaysia

Ahmed Ali Shah<sup>6</sup>

Department of Electrical Engineering Sukkur IBA University, Sukkur Pakistan

Sheeraz Ahmed<sup>7</sup>

Department of Computer Science Iqra National University, Peshawar Pakistan

**Abstract**—Traditional current-mirror circuits require buck converter to deal with one fixed current load. This paper deals with improved self-adjustable current-mirror methods that can address different LED loads under different conditions with the help of one buck converter. The operating principle revolves around a dynamic and self-configurable combinational circuit of transistor and op-amp based current balancing circuit, along with their op-amp based dimming circuits. The proposed circuit guarantees uniformity in the outputs of the circuit. This scheme of current-balancing circuits omitted the need for separate power supply to control the load currents through different kinds of LEDs, i.e. RGB LEDs. The proposed methods are identical and modular, which can be scaled to any number of parallel current sources. The principle methodology has been successfully tested in Simulink environment to verify the current balancing of parallel LED strings.

**Keywords**—Current-balancing; LED driver; current mirror

## I. INTRODUCTION

With the invention of visible phenomenon known as light-emitting diodes (LEDs), which is about a half-a-century old, they have now come into prominence for the past few years. Nowadays, LEDs are preferred choices since they are energy efficient and reliable [1]. The trend of utilizing LEDs is growing day by day and it is becoming the vital mandatory parts in some applications [2]. This has attracted the attention of many researchers, to study its applications. A lot of research has been going on to address various issued of LEDs, and one of the main issues is related to the imbalances of currents in parallel LED strings have drawn attention in recent years as LED technology is perceived as an emerging

technology for replacing incandescent and fluorescent lamps. A most conventional method of current balancing is to connect each individual LED string [3] with resistor in series, which the power loss may lower the system efficiency. Various passive methods have been implemented with higher capacitors impedance to dominate the load impedance to achieve the current balancing [3] – [11]. In [12] – [18], implementing the combined methods to obtain the current balancing; however, the designed circuit contributed to the complicated ones.

In most approach, driving RGB LED arrays, need three converters, with their allied inductors. This causes the number of counts and increases the cost significantly [19]. As observed in [20], the RGB LED driver has been treated as a single input multi-output (SIMO), which is operated by a single inductor with different controllers for running RGB LED modules, to decrease the counts' numbers and to decrease the cost. SIMO drives RGB LEDs in a sequential scanning color scheme display (SCD) to reduce power consumption with fast response.

In [21], whenever the load varies suddenly, inductor current increases to get a new balance. Such the phenomenon is not considered suitable for the other converter of the system since the earlier converter takes time to get new balance and impact the output of other converters, which is in sequence. It eventually produces delay in the other converter's output and creates problems in getting uniform illumination of RGB LED modules. It is generally desired, to have a high degree of color stability for high-end applications like medical or museum lighting [22].

Meanwhile, [23] proposed the constant current sources for each load to energize the LEDs, i.e. a LED or string of LEDs. This method sounds easy, but it is costly and needs more components to the driver circuit, which creates the whole system complex. To address the issue, the current mirror (CM) is employed to remove imbalances in the currents of the strings. This solution requires very few components and is also less expensive solution in comparison to using only current sources. The CM driver has been recommended to be used for LED applications, where there is less likelihood of high value of current [24] and enhances the circuit response time and reduces complexity [25].

Dimming is an essential factor in the effective control of lighting and in saving energy, but it faces various challenges [26]. According to [21], while implementing analog dimming it is tough to maintain higher efficiency throughout the operation of wide range of diming. Contrary to this, PWM dimming approach has been used, and it is possible to utilize the full dimming range, since it has been observed that the color scheme of RGB LEDs can be controlled by changing the ratio of pulse width modulation (PWM) [27]. As stated in [20], there are three techniques available to control dimming phenomenon on system-level, and two techniques are available to control the dimming operation on string level. These methods need three different controllers as well as specific techniques like instant-duty-restoration (IDR), which makes the system costly and complicated.

Therefore, this article is meant to develop a CM circuit that facilitates different kinds of current controllers to be used efficiently. It also improves the time response of the entire system, under new desire output of dimming phenomenon. Secondly, to develop the RGB LED driver with single controller and inductor, which could overcome the issues related to the cross regulations. Thirdly, to build dimming circuit to exploit the full range with the technique of random PWM without using the method of IDR on string level as well as on the system level. This paper mainly discusses self-configurable CM methods for improving current imbalance in different parallel loads of color LEDs. The design method does not require a separate power supply for powering the CM circuit as well as their associate's circuits. Simulations have been conducted in the environment of Matlab/Simulink, to validate the proposed RGB LED driver.

## II. CIRCUIT DESCRIPTION

The circuit configuration is for driving three string loads, in which each string connecting nine LEDs in parallel for red, green and blue LEDs. The loads have been designed to operate with buck converter, with the intent to provide PWM signals for the LED loads. The proposed particular LED driver system consists of a DC-DC buck converter or current generator, current mirror, LED loads and current regulator as depicted in Fig. 1.

The modified super diode in [28] for CM circuit is employed, while by applying some modifications i.e. diode of the super diode circuit with a resistor, along with the combination of op-amp based circuit for dimming with a properly biased transistor, to reduce power losses across

transistor. The load arrays of RGB LEDs, along with their associates' circuits are shown in Fig. 2.

The proposed approach has a dynamic and self-configurable current-balancing circuit structure that allows the best current source (i.e., the smallest current source in the case of current balancing of parallel LED strings) to be selected. The proposed CM based current balancing circuit (refer Fig. 2) not requires 1) external power supply and 2) associated control circuit. Three parallel loads (LED strings, consisting of red, green, and blue LEDs) are connected to a self-configurable CM circuit. The transistors Q1 to Q3 (also called Q-transistors) represent the transistors used in a proposed CM circuit, such as the one shown in Fig. 2. Extra resistors that may be required for proper biasing of the circuit's transistors are not shown in Fig. 2 for the sake of simplicity.

Three transistor-based differential circuits are required along with super diode circuits for three difference loads connected in parallel as shown in Fig. 2. The differential circuit, which has the lowest load current, is selected as a current reference in a CM circuit to prevent the saturation of the transistors in the remaining differential circuits.

In analysis purposes, for the sake of simplicity, all the transistors are considered matched with the same current gain,  $\beta$ , and same resistive loads. If the current imbalances among the parallel current sources are not too significant (i.e., current inequality has been reduced), the transistor of the red LED with the smallest  $V_{CE}$ , drives the bases transistors of other differential circuits through current mirrors, forcing all the associated transistors with the super diode to work linearly.

Furthermore, the CM make changes in  $V_{CE3}$  and  $V_{CE2}$  to reduce  $I_1$  and  $I_2$  to follow the current reference  $I_1$ . This operating mode is based on the CM concept, except that there is a newness of a self-configurable feature that allows the best current source to be dynamically chosen as the reference current source for the CM operation.

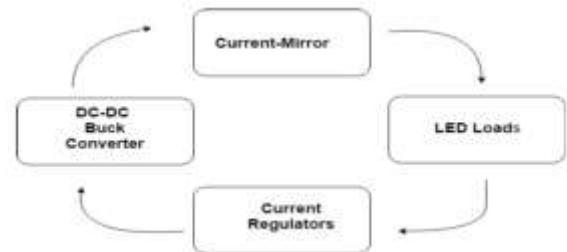


Fig. 1. Block Diagram of the Proposed LED Driver.

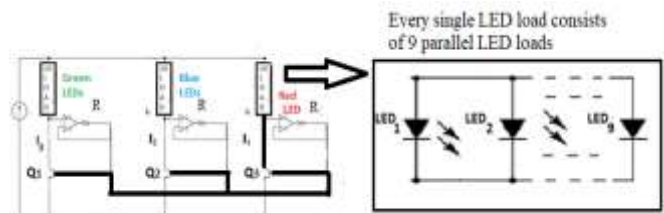


Fig. 2. Self-Configurable RGB LED Current Mirror Circuit.



Using the assumption of  $I_1 > I_2 > I_3$ , and  $V_{CE3} > V_{CE2} > V_{CE1}$ , the self-configurable principle can be illustrated with the aid of Fig. 2. The  $V_{CE3}$  being the smallest voltage across Q3 and the critical conducting path is highlighted with the bolded line in Fig. 2.

### III. DIMMING CONTROL MECHANISM-A HYBRID FUSION OF VOLTAGE DIVIDER BIAS CIRCUIT AND COLLECTOR FEEDBACK RESISTOR

A proper biasing is mandatory to operate transistor and to prevent it from going into saturation mode. It is a phenomenon related to the arrangement of dc collector current at a specific dc voltage by setting up an appropriate quiescent point. The biasing scheme is adopted by putting the base resistor  $R_B$  in between the collector and the base terminals of the transistor as shown in Fig. 3. The resistor  $R_A$  is connected between the emitter and the base terminals of the transistor, for sufficient biasing condition provided by the voltage drop across  $R_A$ . To lessen the losses, the kind of virtual resistance is proposed, which is not presented physically between the dimming circuit and the rest of the biasing circuit of the transistor. The dimming circuit is set to create different ground voltage references for the flow of load current from the biased transistor. Thus, by varying different voltage ground references, it becomes possible to create a kind of virtual resistance in the path of flow of load current, without having a real physical resistance over there.

All these LEDs need an individual dimming system. The applied dimming frequency to the dimming switching circuit is 4Khz. The switching phenomenon has been implemented in each leg of the CM circuit, as shown in Fig. 4 and completed three string dimming is shown in Fig. 5.

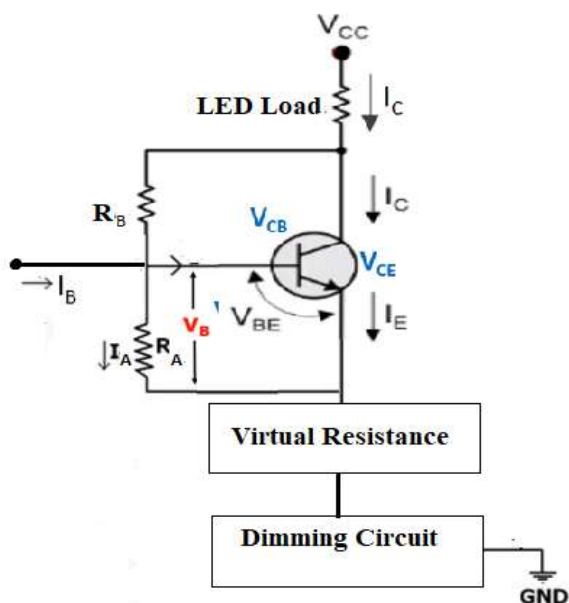


Fig. 3. Dimming Control Mechanism.

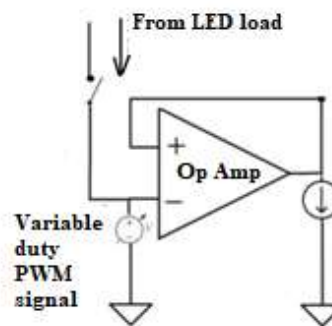


Fig. 4. Dimming/Switching Circuit and its Implementation in the Legs of CM.

### IV. MATHEMATICAL ANALYSIS FOR LED PARAMETER

The total load consists of three loads i.e., the string of red LEDs, string of green LEDs and the string of blue LEDs. The parameters of RGB LEDs have been taken from [19]. The  $R_{EQU}$  represents the equivalent resistance of the red LED, with respect to its rated forward voltage and load current, whereas  $R_{OUT}$  shows the parallel combination of the nine parallel-connected red LEDs in equation (2), same calculations have been done with respect to blue or green LEDs in equation (4). In the proposed scheme, there exists a total of three parallel strings, each individual string carries the load of parallel connected of nine LEDs with respect to their associated colors i.e. red, green and blue. The DC equivalent resistance ( $R_{EQU}$ ) of the red LED has been found at their DC operating points [19] as follows:

$$R_{EQU} = \frac{V_{FWD}}{I_{LED}} = \frac{1.9}{20 \times 10^{-3}} = 95 \Omega \quad (1)$$

$$R_{OUT} = R_{EQU1} // \dots // R_{EQU9} = 10.56 \Omega \quad (2)$$

The currents forward voltages of the green and blue LEDs are similar. Their DC equivalent resistances of the green and blue LED can be found at their DC operating points as follows;

$$R_{EQU} = \frac{V_{FWD}}{I_{LED}} = \frac{3.2}{20 \times 10^{-3}} = 160 \Omega \quad (3)$$

$$R_{OUT} = R_{EQU1} // \dots // R_{EQU9} = 17.78 \Omega \quad (4)$$

Where, the  $V_{FWD}$  denoted a forward voltage of LED,  $I_{LED}$  is the current pass through the LED string and  $R_{out}$  is represented the total resistance of 9 LEDs in each string. It is noted that the red LED consumed less resistive than green and blue LEDs.

In the proposed dimming control circuit, it can introduce a new dimming level across a single and whole load array of LEDs. In principle, when dimming, the associated transistor of the string reduces the quantity of current through the load by raising the emitter's voltage, which eventually increases the collector's voltage of the transistor. It gives a freedom to LED running with a constant source of voltage with emitter coupled logic (ECL) topology, whereas avoiding transistors go to saturation. In which, the ECL is also known to be high speed integrated bipolar transistor logic. Hence, the combinational circuit of super diode, ECL topology and dimming control

circuit provide a kind of controlling mechanism in providing a constant load current.

Refer to Fig. 5,  $R_e$  is added in the designed as thermal runaway resistor, which has a small value to reduce the conduction losses and omitting the thermal runaway of transistors. The dimming on system level can be done easily while activating all the dimming circuits, which the  $R_l$  is used to change the output voltage of op-amp of super diode to current. Meanwhile, the  $R_B$  and  $R_A$  are the feedback resistors and biasing resistors that used to reduce the power losses across the transistor. The efficiency ( $\eta$ ) of the LED driver can be computed by the ratio of the power dissipated by the LED to the total power dissipated in the string [29] as in following equation.

$$\eta = \frac{P_{LED}}{P_{LED} + P_{comp} + P_{tran}} \times 100 \quad (5)$$

### A. Biasing Scheme

In this approach, a proper biasing scheme is needed to do dimming uniformly throughout the whole system comprises of different resistances. The red LED has different resistive nature as compared to the green and blue LEDs, which leads to a differences of load current in red LED string. To make the system efficient, the load of red LED needs to be equal with other loads string, which can be made possible by adding external resistance (compensating resistance) in series with total red LEDs. However, such arrangements, contribute certain power losses as well as degrading the factor of efficiency  $\eta$  as stated in equation (5).

To make the system comparatively efficient, a proper biasing scheme is needed. The biasing scheme adopted to all transistors  $Q_1$ ,  $Q_2$  and  $Q_3$  using dual feedback transistor biasing scheme, as shown in Fig. 6. This sort of biasing method is beta ( $\beta$ ) dependent. The resistor  $R_B$  is utilized for the collector to base feedback configuration to ensure transistor to always remain biased in the active region regardless of any value of the beta ( $\beta$ ) factor. The base bias voltage is dependent on the collector voltage, thus ensures good stability.  $R_A$  has been added to the system to improve stability even more with respect to variations in beta ( $\beta$ ) by increasing the current flowing through the base biasing resistors. Thus, it increases the possibilities of a wide range of dimming spectrum.

A modification for biasing circuit of  $Q_3$  transistor is carried out due to difference resistive load of red LED, by putting a base limiting resistance ( $R_{bl}$ ) to limit the base signal of  $Q_3$  transistor (refer Fig. 6). The purpose of its introduction in the circuit is to improve the stability of the overall system by making  $Q_3$  less responsive as compared to the response behavior of the other two  $Q_2$  and  $Q_1$  transistors by making beta ( $\beta$ ) of  $Q_3$  less sensitive as compared to the transistors of other loads. Thus, the load current flowing through the  $Q_3$ , becomes almost identical to the load currents flowing through the transistors  $Q_2$  and  $Q_1$ .

## V. VALIDATION

A dimming approach for RGB LED driver has been set up to evaluate the performance of the self-configurable CM

circuits. A DC-DC buck converter circuit has been set up to act as a current source for the proposed circuit.

### A. Self-Configurable Current-Mirror Circuit

Since the red LED has less resistive value than other color LEDs, the modification of the dimming circuit is proposed so that the dimming can be conducted for all three consecutive string of RGB LEDs as shown in Fig. 5 and the designed parameters is depicted in Table I. The power of each string is always operated under the full load condition when the LED string is turned on, and no-load condition when LED string off. A SIMO driver for RGB LED with PWM dimming mechanism producing a relatively constant power efficiency at any dimming ratio, but in practice, there exists a small variation in their values as shown in the sets of calculations with references to Table II.

It has been observed that in a situation, where  $R_{BL}$  is kept equal to  $3\Omega$ , and alternately turning on and off of red, green and blue LEDs. It is found that a significant differences of load current red and blue/green LEDs as shown in Fig. 6 to Fig. 8. The efficiency of the red LED drop significantly and return increasing when rising the  $R_{BL}$ . However, the load current increases accordingly as listed in Table II.

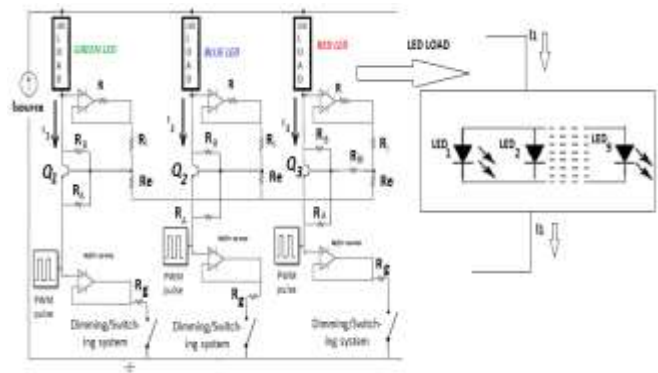


Fig. 5. Modified Self-Configurable RGB LED Current Mirror Circuit.

TABLE I. DESIGN SPECIFICATIONS OF AN OFF-LINE SIMO LED DRIVER

| Parameter                | Value        |
|--------------------------|--------------|
| Input voltage ( $V_g$ )  | 5 volts      |
| Maximum current          | 22.17 mA     |
| Minimum current          | 2.453 mA     |
| Main switching frequency | 300 KHz      |
| Dimming frequency        | 4KHz         |
| Inductor load current    | 540 mA       |
| $R$                      | $0.6 \Omega$ |
| $R_B$                    | $9 \Omega$   |
| $R_A$                    | $0.5 \Omega$ |
| $R_e$                    | $2 \Omega$   |
| $R_{BL}$                 | $3 \Omega$   |
| $R_l$                    | $20 \Omega$  |
| $R_g$                    | $1 \Omega$   |

TABLE II. DESIGN SPECIFICATIONS OF AN OFF-LINE SIMO LED DRIVER

|    | $R_{BL}$      | Load Current through Red LED | Load Current through Green/Blue LED | $\Delta$ Current | Red LED $\eta$ | Blue/Green LED $\eta$ |
|----|---------------|------------------------------|-------------------------------------|------------------|----------------|-----------------------|
| 1. | 3 $\Omega$    | 210 mA                       | 220.4 mA                            | 10.4             | 21%            | 57%                   |
| 2. | 50 $\Omega$   | 221.5 mA                     | 211.5 mA                            | 10               | 45%            | 51%                   |
| 3. | 100 $\Omega$  | 222.4 mA                     | 210.8 mA                            | 11.6             | 48%            | 51%                   |
| 4. | 200 $\Omega$  | 222.9 mA                     | 210.5 mA                            | 12.4             | 52%            | 51%                   |
| 5. | 500 $\Omega$  | 223.2 mA                     | 210.5 mA                            | 12.7             | 52%            | 51%                   |
| 6. | 1000 $\Omega$ | 223.5 mA                     | 210.5 mA                            | 13               | 52%            | 51%                   |

### VI. DIMMING

There are two options of dimming available for the LEDs, based on the operations of the dimming circuits, associated with their LED loads. First is associated with the dimming available for every individual LED load, which is by controlling their associated dimming circuits. Second option is to facilitate the user to do dimming operations simultaneously for all string. The proposed LED driver has a 5V DC voltage from buck DC-DC converter. At first, the validation is carried out through the circuit simulation.

#### A. String-Level Dimming Evaluation

In this dimming mode, the technique of employing super diode along with its dimming circuit is purposed to evaluate the circuit ability for controlling three separate independent different load currents. The load currents of  $I_1$ ,  $I_2$ ,  $I_3$  and dimming gate signal are shown in Fig. 6. The dimming signal is set from 10% to 90% of range. Fig. 6 to Fig. 8 show the current waveforms ( $I_1$ ,  $I_2$  and  $I_3$ ), as well as the corresponding gate signals for the dimming circuits, associated with their load currents. The reference output currents for the three LED strings are respectively set at 210 mA, 220 mA and 220 mA, and the dimming duty cycles is set at 10%. For further analysis, the 10% dimming signal is equal to on-state while 90% dimming signal is off-state.

The result showed that during the on-state, the output currents are regulated to the reference. In comparison, the three LED strings are set identical of dimming parameters, i.e., the same on-state current references and dimming ratios frequency. The dimming operation is generated at differences timing for showing the robustness of the system. These waveforms are almost identical to the voltage drops across their associated loads. The behaviors waveforms show that the loads are operating according to their power ratings.

In Fig. 6, illustrating the red LED behavior of load current, voltage across the load and dimming signal, respectively. It is

clearly visible that the voltage across the red LEDs following the load current. Referring to dimming signal, the on and off states are implemented. For 10% of dimming, the maximum amount of current is found approximately to 210 mA, while during 90% of dimming, the minimum amount of current is found drop to around 25 mA. The loading effects could be seen during off-state for the individual LEDs, which is not more than 2% of the full load current. For others string, the green and blue LEDs are shown in Fig. 7 and 8, respectively. It is clearly seen that the signal behavior comparable with red LED load. During 10% dimming, the current flow comparable to the reference of 220 mA while during off-state drop to around 25 mA and 30 mA for green and blue respectively. The loading effects also could be seen at off-state period.

On the other hand, through the proposed dimming mechanism, the buck dc-dc converter does not require closed-loop control to control the dimming mechanism as shown in Fig. 9. As a result, the range of the dimming signals could be extended from 10% up to 90% dimming, which is almost to the full range dimming.

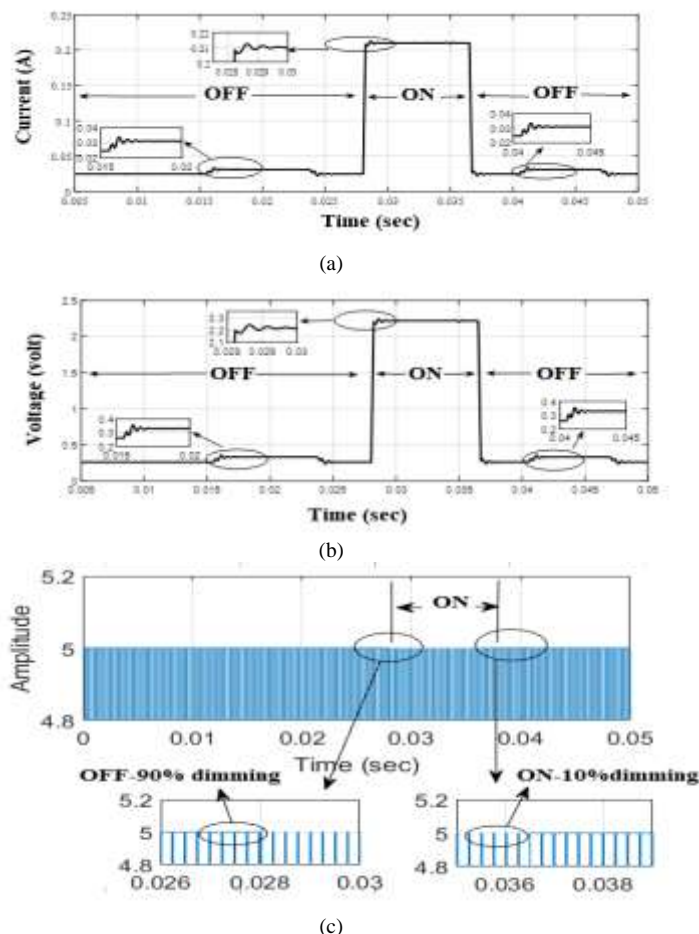
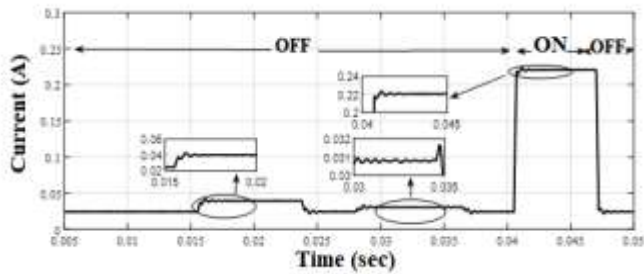
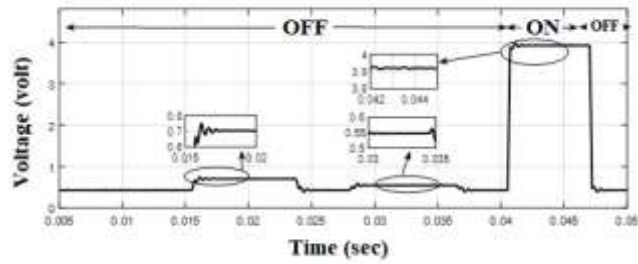


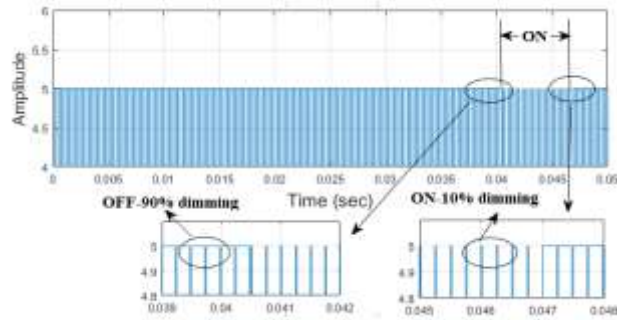
Fig. 6. Waveforms of the Output with PWM Dimming of 10% and 90% (a) Load Current Response of Red LEDs (b) Voltage Across the Red LEDs (c) Close-up view of the widths of the PWM Signals for Turning off and on the red LEDs.



(a)

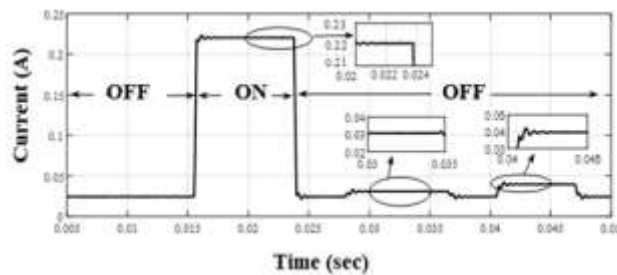


(b)

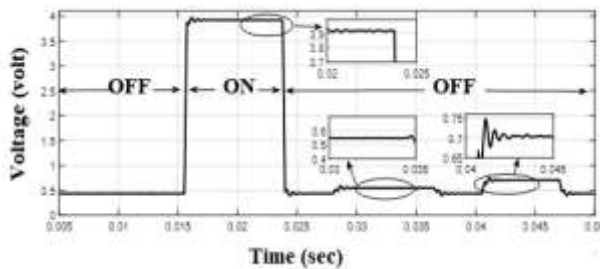


(c)

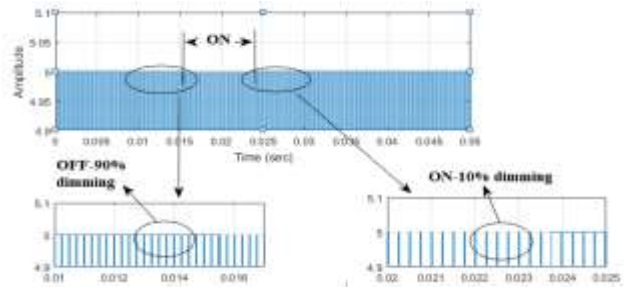
Fig. 7. Waveforms of the Output with PWM Dimming of 10% and 90% (a) Load Current Response of Green LEDs (b) Voltage across the Green LEDs (c) Close-up view of the widths of the PWM Signals For turning off and on the Green LEDs.



(a)



(b)



(c)

Fig. 8. Waveforms of the Output with PWM Dimming of 10% and 90% (a) Load Current Response of Blue LEDs (b) Voltage Across the Blue LEDs (c) Close-up view of the widths of the PWM Signals for Turning off and on the Blue LEDs.

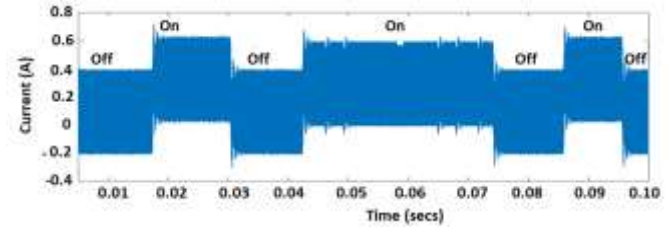


Fig. 9. Inductor's Load Current Variations Corresponding to the Turning off and on of the LEDs through Dimming Circuits.

### B. System Level Dimming Evaluation

Second scenario, the combinational circuits of super diode and dimming circuit for all the load strings are analyzed. Simultaneously operating in the range of 10% to 90% of dimming signal is implemented. In Fig. 10, showing the results for load current, voltage across the load and dimming signal respectively. The behavior is comparable according to the on and off states conditions. Nevertheless, in the system level dimming, all the loads take current spontaneously and share their currents with other loads. Since the current at green and blue LEDs are slightly more than red LEDs, hence the surplus current goes to the red LEDs due to the proposed circuit act to balance current through the string with differences loads. As the result, the load current at red LEDs increases significantly. It is showed the effectiveness of the proposed mirroring circuit in managing the current sharing. It is different with the string level dimming, which is the green and blue LEDs are not activated when the red LEDs on. The process same with other LEDs string turning on. Thus, avoiding the phenomenon of current balancing from other string LEDs.

In Tables III and IV, listing the computed power losses through the LED string loads (i.e.; red, green and blue LEDs). The condition is activating the dimming simultaneously for three string LEDs loads to figure out the efficiency. Through the measurement, the maximum load currents for red and green/blue LEDs are 221.8 mA and 219.0 mA respectively. For red LEDs, the measured voltage drop across the transistor and current flow through the transistor are 2.111 V and 229.6 mA, respectively. Meanwhile, 0.56 V and 243.7 mA are measured through green/blue LEDs. Hence, the power losses through transistor for red and green/blue LEDs could be computed accordingly as shown in Tables III and IV.

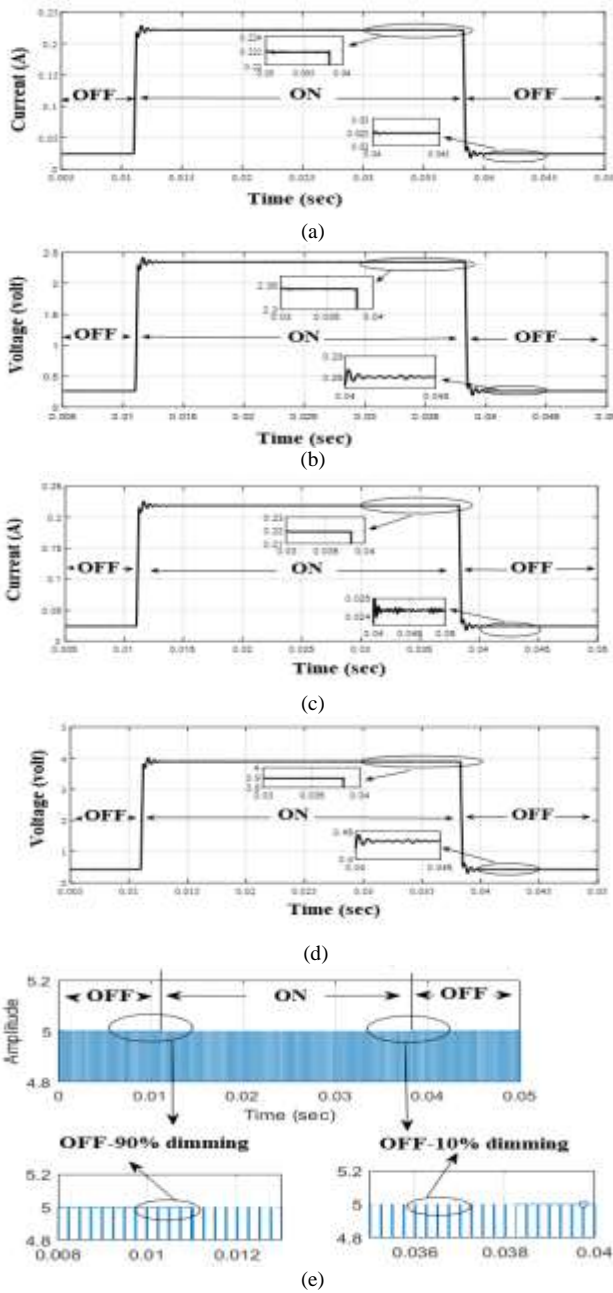


Fig. 10. Waveforms of the System Output with PWM Dimming of 10% and 90% (a) Load Current Response of Red LEDs (b) Voltage Across the Red LEDs (c) Load Current Response of Green/Blue LEDs (d) Voltage Across the Green/Blue LEDs (e) Close-up view of the widths of the PWM Signals for Turning off and on the LEDs.

Thus, the computation efficiency for red ( $\eta_{r\_max}$ ) LEDs is 52%, and 86% efficiency ( $\eta_{gb\_max}$ ) for green/blue LEDs.

On the other hand, the measurement of power losses at minimum load currents for red and green/blue LEDs are 24.58 mA and 24.27 mA. In red LEDs string, the measured transistor voltage drops and current pass through the transistor are 0.234 V and 25.44 mA respectively. Meanwhile, for green/blue LEDs string, 62.04 mV and 27.0 mA are measured for voltage drop and current through the transistor. Thus, the power losses can be estimated as listed in Tables V and VI.

TABLE III. POWER LOSSES FOR RED LEDs AT MAXIMUM CURRENT

|                                     |        |
|-------------------------------------|--------|
| Red LEDs ( $P_{LED}$ )              | 520 mW |
| Transistor terminal ( $P_{trans}$ ) | 480 mW |

TABLE IV. POWER LOSSES FOR GREEN/BLUE LEDs AT MAXIMUM CURRENT

|                                     |          |
|-------------------------------------|----------|
| Green/Blue LEDs ( $P_{LED}$ )       | 853 mW   |
| Transistor terminal ( $P_{trans}$ ) | 136.3 mW |

TABLE V. POWER LOSSES FOR RED LEDs AT MINIMUM CURRENT

|                                     |        |
|-------------------------------------|--------|
| Red LEDs ( $P_{LED}$ )              | 6.4 mW |
| Transistor terminal ( $P_{trans}$ ) | 6 mW   |

TABLE VI. POWER LOSSES FOR GREEN/BLUE LEDs AT MINIMUM CURRENT

|                                     |         |
|-------------------------------------|---------|
| Green/Blue LEDs ( $P_{LED}$ )       | 10.5 mW |
| Transistor terminal ( $P_{trans}$ ) | 2 mW    |

Therefore, the estimation efficiency at minimum load current for red ( $\eta_{r\_min}$ ) LEDs is 52%, and for green/blue ( $\eta_{gb\_min}$ ) LEDs is 85%.

## VII. CONCLUSION

In this paper, an investigation has been carried out to determine the effectiveness of the current mirroring circuit to regulate the different LED loads, while keeping and regulating the minimum power losses through the string module. Though green and blue LEDs have same resistive in nature, which is slightly difference with red LED, the current mirror has successfully managed to bring identical currents passing through the loads. The proposed circuit could be implementing two option of dimming, which are individual string level and system level dimming. It is show flexibility of the dimming purposes. Another advantage, it has found that, the analyses has done the measurement of precisely dimming from 10% to 90% of dimming range. It is showed the effectiveness of the proposed LED driver and dimming circuit for color LED string.

## REFERENCES

- [1] K.H. LiW.Y. FuH.W. Cho.: Chip-scale GaN integration, Progress in quantum electronics 70 (2020) 100247.
- [2] Durrani S.S., Zaharin A., Hassan B., Ishak R.B. (2020) Comparative Analysis for LED Driver with Analog and Digital Controllers. In: Kasruddin Nasir A. et al. (eds) InECCE2019. Lecture Notes in Electrical Engineering, Vol 632. Springer, Singapore. [https://doi.org/10.1007/978-981-15-2317-5\\_73](https://doi.org/10.1007/978-981-15-2317-5_73).
- [3] Mudassar Khatib.: Ballast Resistor Calculation-Current Matching in Parallel LEDs, Texas Instruments, April 2009.
- [4] W. K. Lun, ., K. H. Loo, ., S. C. Tan, ., Y. M. Lai, ., C. K. Tse.:Bi level Current Driving Technique for LEDs, IEEE Trans. Power Electronics, vol. 24, no. 12, pp. 2920-2930, 2009.
- [5] C. H. Lin, ., T. Y. Hung, C. M. Wang, ., d K. J. Pai.:A Balancing Strategy and Implementation of Current Equalizer for High Power LED Backlighting, IEEE Power Electronics and Drive Systems PEDS '07. 7th International Conference, pp. 1613-1617, Nov. 2007.
- [6] K. H. Jung, J. W. Yoo, ., C. Y. Park.:A Design of Current Balancing Circuit for Parallel Connected LED strings using Balancing Transformers, IEEE Power Electronics and ECCE Asia (ICPE & ECCE), pp. 528-535, 2011.

- [7] K. I. Hwu., S. C. Chou: A Simple Current-Balancing Converter for LED Lighting" IEEE Applied Power Electronics Conference (APEC) Proc, pp. 587-590, February 2009.
- [8] Y. Hu, M. M. Jovanovic: A New Current-Balancing Method for Paralleled LED Strings, IEEE Applied Power Electronics Conference (APEC), 26th Annual, Mar. 2011.
- [9] J. Zhang, ., J. Wang, . X. Wu,:A Capacitor-Isolated LED Driver with Inherent Current Balance Capability, IEEE Trans. Industrial Electronics, vol. 59, no. 4, pp. 1708-1716, 2011.
- [10] S. Choi, P. Agarwal, ., T. Kim, ., J. Yang, ., ., B. Han, Symmetric Current Balancing Circuit for Multiple DC loads, IEEE Applied Power Electronics Conference (APEC) Proc., pp. 512-518, 2010.
- [11] S. M. Baddela, D. S. Zinger, "Parallel Connected LEDs Operated at High Frequency to Improve Current Sharing," IEEE Industry Applications Conference, 2004.
- [12] J.Wang, J. Zhang, Y. Shi., Z. Qian: A Novel High Efficiency and Low-Cost Current Balancing Method for Multi-LED Driver, IEEE Energy Conversion Congress and Exposition (ECCE), pp. 2296-2301, 2011.
- [13] S. H. Cho, S. H. Lee, ., S. S. Hong, D. S. Oh . S.K. Han, "High-Accuracy and Cost-Effective Current-Balanced Multichannel LED Backlight Driver Using Single-Transformer, IEEE Power Electronics and ECCE Asia 8th International Conference, pp. 520-527, 2011.
- [14] S. Zhang, Q. Chen, J. Sun, M. Xu, Y. Qiu, High-Accuracy Passive Current Balancing Schemes for Large-Scale LED Backlight System, IEEE Applied Power Electronics Conference and Exposition (APEC), pp. 723-727, 2011.
- [15] J. Wang, J. Zhang, ., X. Huang, . L. Xu,:A Family of Capacitive Current Balancing Methods for Multi-Output LED Drivers, IEEE Applied Power Electronics Conference and Exposition (APEC), pp. 2010-2046, 2011."A Dimmable Light-Emitting Diode ".
- [16] W. Chen, S. Y. R. Hui,:A Dimmable Light-Emitting Diode (LED) Driver with cascaded mag-amp postregulators for multistring applications, IECON Conference on IEEE Industrial Electronics Society, pp. 2523-2528, 2010.
- [17] T. Werner, J. Pforr.: A Novel Low-Cost Current-Sharing Method for Automotive LED Lighting System, Power Electronics and Applications European Conference, pp. 1-10, 2009.
- C. Zhao, X.Xie, ., S. Liu,:A Precise Passive Current Balancing Method for Multi-Output LED Drivers, IEEE Trans. Power Electronics, 2011.
- [18] J. Hasan, S. S. Ang,:A high-efficiency digitally controlled RGB driver for LED pixels, IEEE Trans. Ind. Appl., vol. 47, no. 6, pp. 2422–2429, Nov./Dec. 2011.
- [19] S. Li, Y. Guo, A. Lee, S. Tan, ., S. Y. Hui,:An off-line single-inductor multiple-outputs LED driver with high dimming precision and full dimming range, IEEE Trans. Power Electron., vol. PP, no. 99, pp. 1–1, 2017.
- [20] Patra, P., Patra, A., Misra, N.: A single-inductor multiple-output switcher with simultaneous buck, boost, and inverted outputs, IEEE Trans. Power Electron., 2012, 27, (4), pp. 1936–1951
- [21] Kwai Hei Li, Member, IEEE, Yuk Fai Cheung , Weijian Jin, Wai Yuen Fu , Albert Ting Leung Lee, Siew Chong Tan , Shu Yuen Hui , and Hoi Wai Choi ,” InGaN RGB Light-Emitting Diodes With Monolithically Integrated Photodetectors for Stabilizing Color Chromaticity”. IEEE Transactions on Industrial Electronics, vol. 67, no. 6, June 2020.
- [22] Muhinthan Murugesu, Osram Opto Semiconductors Vector, Current distribution in parallel LED strings, Vector – journal of the Institution of Certificated Mechanical and Electrical Engineers (ICMEESA), and the Illumination Engineering Society of South Africa (IESSA), January 2013.
- [23] H.-J. Chiu, ., S.-J. Cheng.: LED backlight driving system for largescale LCD panels, IEEE Trans. Industrial Electronics, Vol. 54, No. 5, pp. 2751-2760, Oct. 2007.
- [24] Shaheer Shaida Durrani, abu Zaharin ahmed,: An Efficient Digitally Controlled for RGB LED Driver, IEEE ICETAS Conference Bahrain , 2017.
- [25] Vili Väinölä, ., Sina khamehchi , ., Hassan rouhi, ., Tapio kukkonen, ., Visnukumar Murugesan,: Project #21 Illumination and colour control in flicker-free LED lighting, Aalto University, School of Electrical Engineering Automation and Electrical Engineering (AEE) Master's Programme ELEC-E8002 & ELEC-E8003 Project work course Year 2017.
- [26] Xin Yu Guo, Guo Chun Wan, and Mei Song Tong: An Intelligent Control System of Music Rhythms by RGB-LED Lamp. 2019 Photonics & Electromagnetics Research Symposium | Fall (PIERS | FALL), Xiamen, China, 17{20 December.
- [27] Sung-Jin Choi, : Adaptive Current-Mirror LED Driver employing Superdiode Configuration, IEEE International Conference on Industrial Technology.2014.
- [28] J. Falin. (2008, 4Q). Compensating and measuring the control loop of a high-power LED driver. Analog Appl.J. [Online]. pp. 14–17. Available: <http://www.ti.com/lit/an/slyt308/slyt308.pdf>.

# Implementation of a Clinical Decision Support Systems-Based Neonatal Monitoring System Framework

Sobowale A. A<sup>1</sup>, Olaniyan O. M<sup>2\*</sup>, Adetan. O<sup>3</sup>, Adanigbo. O<sup>4</sup>  
Esan. A<sup>5</sup>, Olusesi. A.T<sup>6</sup>, Wahab. W.B<sup>7</sup>, Adewumi. O. A<sup>8</sup>

Department of Computer Engineering, Federal University Oye-Ekiti, Nigeria<sup>1,2,4,5</sup>

Department of Electrical and Electronic Engineering, Ekiti State University Ado-Ekiti, Nigeria<sup>3</sup>

Department of Electrical, Electronics & Computer Engineering, Bells University of Technology Ota, Nigeria<sup>6</sup>

Department of Computer Engineering, Ladoke Akintola University of Technology Ogbomosho, Nigeria<sup>7</sup>

Department of Computer Science and Information Technology, Bells University of Technology Ota, Nigeria<sup>8</sup>

**Abstract**—A Clinical Decision Support-based information systems to monitor the vital signs of the neonate's conditions in prematurely born babies placed in infant incubators of Neonatal Intensive Care Unit (NICU) is developed in this work. A DMS was developed consisting of a supervisory microcomputer and sensitive sensors for measuring the vital signs. The Conventional Monitoring System (CMS) was used simultaneously with the DMS to collect the vital sign readings of thirty (30) neonates, over a period of one week. Fuzzy Inference System CDSS (FIS-CDSS) was developed for the three inputs: Temperature, Heart rate and Respiration rate (THR) based on their membership functions' value (low, medium, high) and twenty-seven (27) IF-THEN fuzzy rules using fuzzy logic toolbox in Matrix Laboratory 8.1 (R2014a). The FIS-CDSS maps the THR to an output status (Normal, Abnormal and Critical). The performance of the FIS-CDSS was evaluated using confusion matrix. The results showed that the system yielded sensitivity ranges of 90 - 100, 80 - 89, 70 - 79, 60 - 69 and 50 - 59% for five, eleven, seven, six and one neonates, respectively with an average sensitivity of 77.92%. The specificity of the system ranged from 5.00 to 66.67% with an associated average specificity of 35.10%. The accuracy of the FIS-CDSS ranged from 70 to 100, 60 to 69, 50 to 59 and 0 to 49% for nine, nine, eight and four neonates, respectively with an average accuracy of 60.94%. The developed system provides adequate and accurate information for on-the-spot assessment of neonates for decision making that improves the mortality rate and recovery period of neonates.

**General Terms:** Neonatal Monitoring

**Keywords**—Clinical Decision Supports Systems (CDSS); Fuzzy Inference System (FIS); Neonatal Intensive Care Unit (NICU); vital signs; neonates

## I. INTRODUCTION

Decision Support Systems (DSS) are increasing in coverage of different sections of life which includes academic, engineering, business, military and medicine [1]. Any automated program that helps specialists in settling on clinical choice is categorized as Clinical Decision Support Systems (CDSS) [2]. CDSS provide clinicians, staff, patients and other individuals with knowledge and person-specific information, wisely separated and displayed at proper times, to upgrade wellbeing, medicinal services and reduce medical errors [3][4][5]. CDSS does not decide; It just gives direction to provide current and pertinent knowledge to clinicians to aid

patient care at the exact time of care delivery [6][7] It is a major technology application to make the right decision at the right time which aids in building an intelligent system for monitoring neonatal vital parameters [8]. The CDS Systems are computer-based information systems used to integrate clinical and patient information to provide support for decision-making in patient care. A category of such patients are the prematurely born babies, which are placed in infant incubators of Neonatal Intensive Care Unit (NICU) for continuous monitoring of their body vital signs (temperature, heart rate and respiration). [19].

Neonates born before thirty-seven (37) weeks gestation are considered premature and are usually in a fragile condition and may be at risk of complications, such babies therefore require special monitoring and intensive care involving treatment in an incubator at an NICU [9][10][11]. Neonatal monitoring refers to the monitoring of vital physiological parameters of premature infants [12]. The survival rate of premature infants is dependent on the continuous monitoring of vital signs; this provides a lot of information about a baby's state of health [18].

In the last decades the advances in sensor technologies and wireless communications technologies have resulted in the possibility of developing intelligent systems for monitoring neonatal vital parameters [13]. Technology therefore provides easy data collection from the neonates monitoring system and aids the neonatologists' in taking appropriate decision.

However, the quality of neonatal care provided by Nigerian hospitals is not uniform and mostly manual, which creates difficulty of interpretation for inexperienced staff [14][15][16] More so, despite the impact of CDSS applications in various sectors of the health system, its application to monitoring of vital signs of preterm babies in the NICU is limited [5], [17].

This paper therefore developed a CDSS that can be used to efficiently monitor the neonate's condition in the incubators of NICU. The paper has five (5) sections in all. Section I is the introduction to the work. Section II gives the architectural framework of CDS systems. The methodology adopted in this work is discussed under Section III. Results obtained in this work are discussed under Section IV while the conclusion is given under Section VI.

\*Corresponding Author

## II. THE CDSS-BASED ARCHITECTURAL FRAMEWORK

The CDSS architectural framework is made up of three components (knowledge base, inference engine and interface) as shown in Fig. 1. This is made up of a set of functional and informational units. The functional unit is divided into the reasoning engine and the connection component. The informational unit comprises the data source and the knowledge base. The knowledge base consists of decision rules, low, medium and high boundary values, diagnosis terms, and clinical recommendation contents. The reasoning engine takes the readings of the vital signs as its data source. After the execution of the decision rules on the data source, the reasoning engine generates the output result, which is displayed on the monitor of the CDSS system and printed from the CDS located at the nursing stand. The clinicians take informed, on the spot decision based on the printed results. This enhances decision making and general performance as the manual routine checks by the nurses is no more the only basis of attending to neonates.

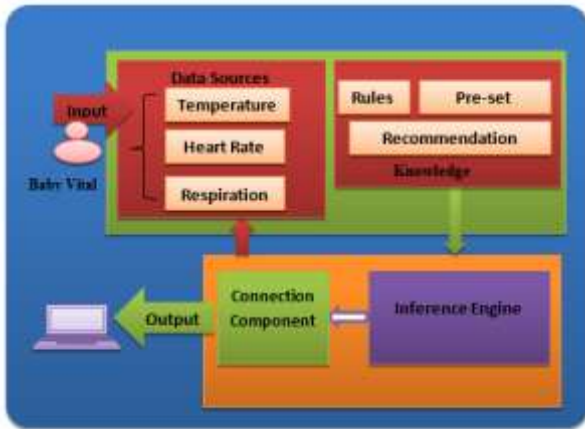


Fig. 1. The CDSS Architecture of the Developed Monitoring System.

## III. METHODOLOGY

The Fuzzy Inference System (FIS) was developed using MATLAB R2014a to implement the CDSS architectural framework. The FIS uses fuzzy logic to map the vital signals Temperature, Heart rate and Respiration rate (THR) to a status (*Normal, Abnormal and Critical*). The output is used to decide on the appropriate treatment for a particular preterm. The FIS decisions are made by the use of membership function and If-Then rules. FIS performs fuzzification on the inputs and defuzzification of the result of fuzzy logic rule to determine the output. Aggregation is used to combine the output of all the rules into a single fuzzy set. The developed FIS takes the vital signs as the inputs and gives "Normal", "Abnormal" or "Critical" as the output. It also consists of the membership functions (MF), antecedents (or premise), consequents (conclusion), weight and connective. A membership function defines the degree to which the value of a vital sign falls within a boundary (or degree of membership). Antecedents are the MF values of the inputs while the consequents are the MF values of the output. A weight determines the level of importance of a rule relative to the others, and the maximum weight a rule can take is 1. A connective takes either "AND" or "OR". The connective "AND" implies that the values of

two antecedents determine the consequents while the connective "OR" implies that any of the antecedents can determine the consequents. The Graphics User Interface (GUI) of the developed FIS is shown in Fig. 2.

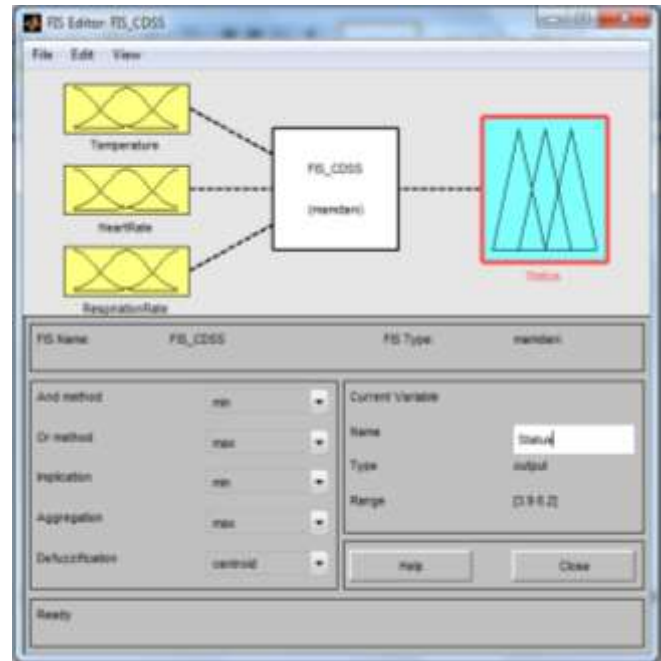


Fig. 2. GUI of the Developed Fuzzy Inference System in MATLAB Environment.

### A. Design of the Membership Function

Three linguistic terms (Low, Medium and High) were used to define the membership function of each of the input variables Temperature, Heart Rate and Respiration rate (THR).

Similarly, three linguistic values Normal (N), Abnormal (A) and Critical (C) were used to define the membership function of the Status or Output (Out) of the inference engine. The value range of the vital signs readings used in the Children Intensive Care Unit (CICU) of Ladoke Akintola Teaching Hospital (LAUTECH) Osogbo, were used to set the range used in the FIS and were classified as Low for readings below the normal range, medium for normal range and high for readings above the normal range, this is discussed below.

**Temperature:** The normal range for Temperature is 36.5-37.5°C; if the input temperature value is more than this range then its MF is High, and if it is below this range then its MF is Low. The classification of Temperature is presented in Table I(A). The MF for the fuzzy set for Temperature (Tmp) is defined as:

$$\text{Low (Tmp)} = \left\{ \begin{array}{ll} 1 & Tmp \leq 32.5 \\ \frac{38 - Tmp}{1} & 32.5 < Tmp < 36.5 \end{array} \right\} \quad 1a)$$

$$\text{Medium (Tmp)} = \left\{ \begin{array}{ll} \frac{Tmp - 32.5}{2} & 35 \leq Tmp < 37.5 \\ 1 & Tmp = 37 \\ \frac{38 - Tmp}{1} & 37 < Tmp < 38 \end{array} \right\} \quad 1b)$$



$$High(Tmp) = \begin{cases} \frac{Tmp - 37.5}{2} & 37.5 \leq Tmp < 39.5 \\ 1 & Tmp \geq 39.5 \end{cases} \quad (1c)$$

**Heart Rate:** The normal range for Heart Rate (Hr) is 130-160 bpm; if the input heartbeat rate value is more than this range then its MF is High, and if it is below this range then its MF is Low. The classification of Heart Rate is presented in Table I(B). The MF for the fuzzy set for heart rate is:

$$Low(H_r) = \begin{cases} 1 & H_r < 125 \\ \frac{132 - H_r}{7} & 125 < H_r < 132 \end{cases} \quad 2a)$$

$$Medium(H_r) = \begin{cases} \frac{H_r - 128}{17} & 128 \leq H_r < 145 \\ 1 & H_r = 145 \\ \frac{162 - H_r}{17} & 145 < H_r < 162 \end{cases} \quad 2b)$$

$$High(H_r) = \begin{cases} \frac{H_r - 158}{12} & 158 \leq H_r < 170 \\ 1 & H_r \geq 170 \end{cases} \quad (2c)$$

**Respiration Rate:** The normal range for Respiration Rate (Rr) is 40-60 cm; if the input heartbeat rate value is more than this range then its MF is High, and if it is below this range then its MF is Low. The classification of Heart Rate is presented in Table I(C). The MF for the fuzzy set for respiration is:

$$Low(R_r) = \begin{cases} 1 & R_r \leq 35 \\ \frac{42 - R_r}{7} & 35 < R_r < 42 \end{cases} \quad 3a)$$

$$Medium(R_r) = \begin{cases} \frac{R_r - 38}{12} & 38 \leq R_r < 50 \\ 1 & R_r = 50 \\ \frac{62 - R_r}{12} & 50 < R_r < 62 \end{cases} \quad 3b)$$

$$High(R_r) = \begin{cases} \frac{R_r - 60}{10} & 60 \leq R_r < 70 \\ 1 & R_r \geq 70 \end{cases} \quad (3c)$$

**Status:** This is the output variable of the FIS. The normal range for Status (Out) is 4-6; if the output value is more than this range then its MF is Critical, and if it is below this range then its MF is Abnormal. The classification of Status is presented in Table I(D).

$$Abnormal(Out) = \begin{cases} 1 & Out \leq 3.5 \\ \frac{4 - Out}{0.5} & 3.5 < Out < 4 \end{cases} \quad 4a)$$

$$Normal(Out) = \begin{cases} \frac{Out - 3.8}{1.2} & 3.8 \leq Out < 5 \\ 1 & Out = 5 \\ \frac{6.2 - Out}{1.2} & 5 < Out < 6.2 \end{cases} \quad 4b)$$

$$Critical(Out) = \begin{cases} \frac{Out - 6}{0.2} & 6 \leq Out < 6.2 \\ 1 & Out \geq 6.2 \end{cases} \quad (4c)$$

The MF plots for Temperature, Respiration rate, Heart rate and Status are shown in Fig. 3 to Fig. 6.

TABLE I. A: CLASSIFICATION OF TEMPERATURE

| Vital Sign  | Range       | Linguistic Term |
|-------------|-------------|-----------------|
| Temperature | < 36.5      | Low             |
|             | 36.5 – 37.5 | Medium          |
|             | > 37.5      | High            |

TABLE I (B): CLASSIFICATION OF HEART RATE

| Vital Sign | Range     | Linguistic Term |
|------------|-----------|-----------------|
| Heart Rate | < 130     | Low             |
|            | 130 – 160 | Medium          |
|            | > 160     | High            |

TABLE I (C): CLASSIFICATION OF RESPIRATION RATE

| Vital Sign       | Range   | Linguistic Term |
|------------------|---------|-----------------|
| Respiration Rate | < 40    | Low             |
|                  | 40 – 60 | Medium          |
|                  | > 60    | High            |

TABLE I(D): CLASSIFICATION OF STATUS

| Output | Range | Linguistic Term |
|--------|-------|-----------------|
| Status | < 4   | Abnormal        |
|        | 4 – 6 | Normal          |
|        | > 6   | Critical        |

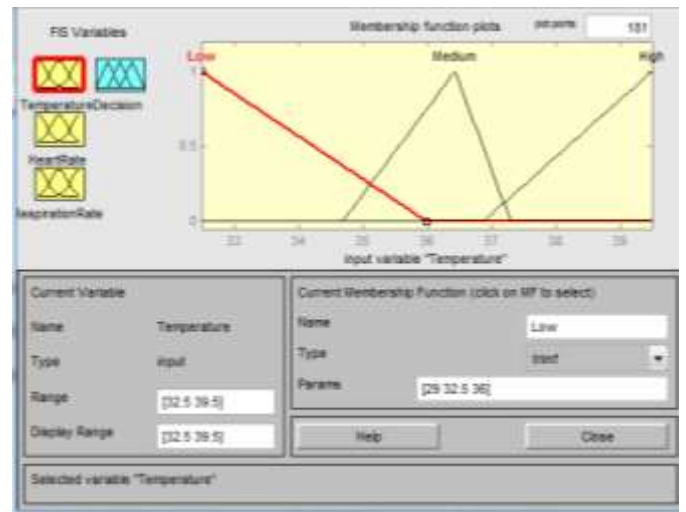


Fig. 3. Membership Functions for Temperature.

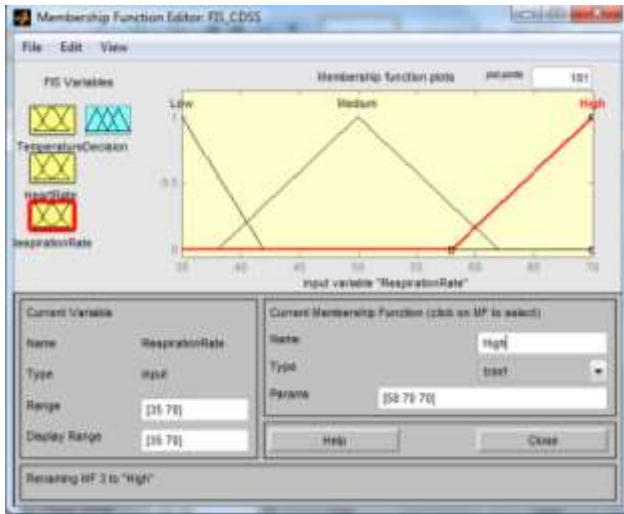


Fig. 4. Membership Functions for Respiration Rate.

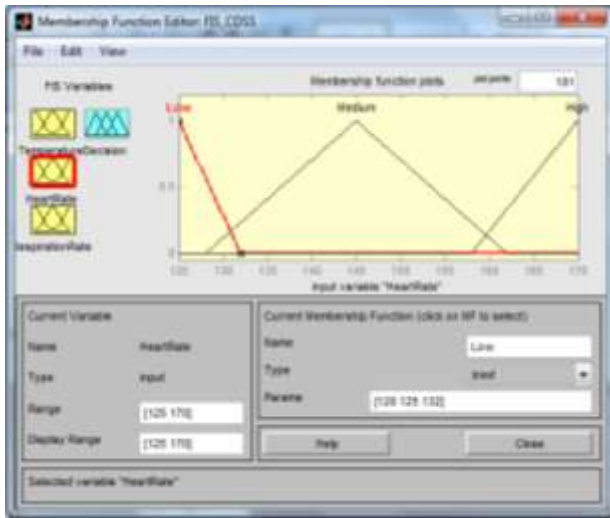


Fig. 5. Membership Functions for Heart Rate.

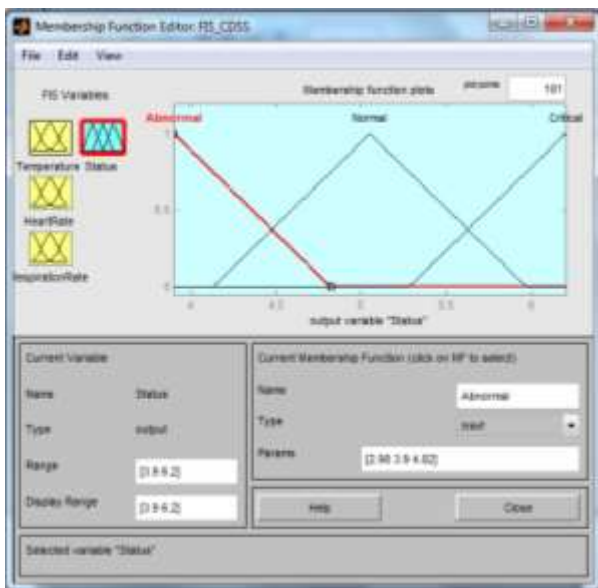


Fig. 6. Membership Functions for Status.

#### IV. RESULTS AND DISCUSSION

An interactive Graphic User Interface (GUI) application was developed using MATLAB R2014a as the frontend and MYSQL 5.1 as the backend to implement the CDSS architectural framework. The developed system named Fuzzy Inference System Clinical Decision Support System (FIS-CDSS) was copied in a folder into the Clinical Database Server (CDS) with a Matlab file (FIS-CDSS\_gui.m); the CDS contains database of the vital signs readings collected from the measuring sensors attached to each neonate. The FIS-CDSS GUI window (Fig. 7) appeared as the filename was executed. The vital signs (Temperature, Heart rate and Respiration) readings from the DMS were loaded into the developed FIS-CDSS as shown in Fig. 8. The loaded readings were run through the FIS-CDSS for classification as shown in Fig. 9 and Fig. 10. The CDSS\_FIS classified the status of the baby (developed system prediction) as *Normal*, *Abnormal* or *Critical* based on the readings and the fuzzy logic rules in the knowledge base of the system; this is shown in Fig. 11, the developed system's prediction can be saved into the CDS as shown in Fig. 12.

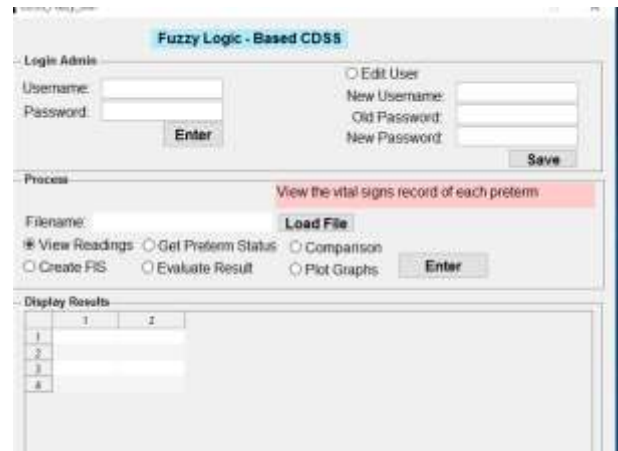


Fig. 7. The Developed FIS-CDSS Window.



Fig. 8. Viewing of the Recorded Vital Signs on the GUI.

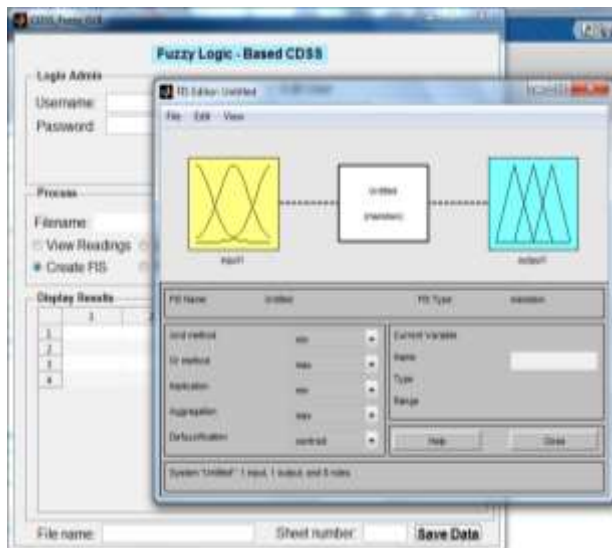


Fig. 9. Creation of a Fuzzy Inference System (FIS) Model for the CDSS.



Fig. 12. Saving of the FIS Classification Results.

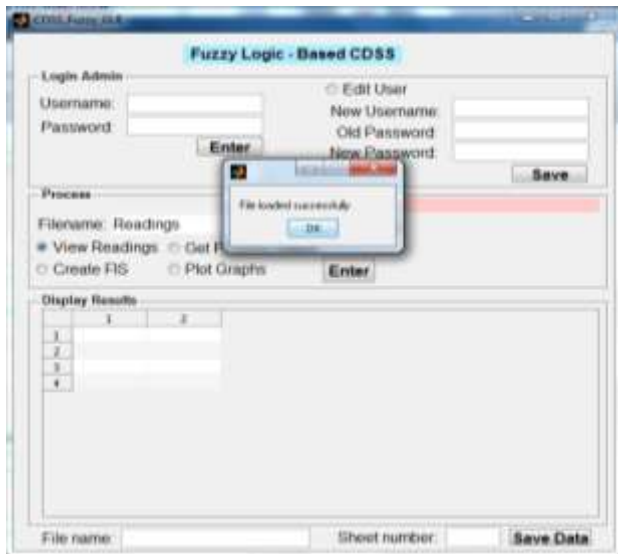


Fig. 10. Acquisition of the Recorded Vital Signs for the FIS Classification.



Fig. 11. Display of the FIS Classification Results.

The predictions of the developed system for the thirty (30) neonates was taken 4 times daily (6:00am, 10:00am, 2:00pm and 6:00pm) for seven (7) days, giving a total of twenty-eight (28) predictions per neonate as shown in Table II.

TABLE II. PREDICTIONS MADE FOR A NEONATE BY THE DEVELOPED SYSTEM (FIS-CDSS)

| PERIOD | FIS-CDSS Prediction |
|--------|---------------------|
| 1      | Normal              |
| 2      | Abnormal            |
| 3      | Normal              |
| 4      | Normal              |
| 5      | Normal              |
| 6      | Abnormal            |
| 7      | Normal              |
| 8      | Normal              |
| 9      | Normal              |
| 10     | Normal              |
| 11     | Normal              |
| 12     | Normal              |
| 13     | Normal              |
| 14     | Normal              |
| 15     | Normal              |
| 16     | Normal              |
| 17     | Abnormal            |
| 18     | Normal              |
| 19     | Normal              |
| 20     | Normal              |
| 21     | Normal              |
| 22     | Normal              |
| 23     | Normal              |
| 24     | Abnormal            |
| 25     | Abnormal            |
| 26     | Normal              |
| 27     | Normal              |
| 28     | Abnormal            |

## V. CONCLUSION AND FURTHER WORK

In this research, a CDSS based architecture for monitoring neonates in the NICU has been implemented. The developed system collects readings of the vital signs of neonates from measuring sensors attached to the wrist of the neonates. Fuzzy Inference System CDSS (FIS-CDSS) was developed for the three inputs: Temperature, Heart rate and Respiration rate (THR) based on their membership functions' value (low, medium, high) and twenty-seven (27) IF-THEN fuzzy rules using fuzzy logic toolbox. The FIS-CDSS maps the THR to an output status (Normal, Abnormal and Critical). The vital signs' readings were fed into the FIS-CDSS, which fuzzifies them and outputs the health status of the neonates.

The research work could be extended to measure or include more factors than the three basic vital signs temperature, heart beat rate and respiration. Other factors being observed by the specialist nurses such as transient clinical death, feeding rate and wavering weather could be included. The research work could also be extended to cover adults and other areas of health could be monitored and remotely reported to the physicians anywhere, anytime.

### REFERENCES

- [1] J.D. Marek and R.F. Roger, "Decision support systems" Decision systems laboratory, School of Information Sciences and Intelligent Systems Program University of Pittsburgh Pittsburgh, pp.3, 2002.
- [2] M.M Abbasi and S. Kashiyarndi, "Clinical decision support systems: A discussion on different methodologies used in health care; Proceedings of the International Conference on Frontiers of Intelligent Systems, pp.1-15, 2020.
- [3] J.A. Osheroff, J.M. Teich and B.F. Middleton, A roadmap for national action on clinical decision support. American Medical Informatics Association, 2006.
- [4] C. Catley, and M. Frize, "A prototype XML-based implementation of an integrated intelligent neonatal intensive care unit" Research work published by University of Ottawa, Canada, 2003.
- [5] K. Tan K, P.R.F. Dear and S.J. Newell, "Clinical decision support systems for neonatal care" published in The Cochrane Library, Issue 2; pp. 2-20, 2009.
- [6] M. Prabhu, P.N. Senthil and K. Lakshmi, "Clinical decision support systems" Computer Sciences Corporation (CSC), pp. 1-19, 2014.
- [7] J. Avansino and M.G. Leu, "Effects of CPOE on provider cognitive workload: a randomized crossover trial, pp. 547-552, 2012.
- [8] M.J. Ball, J.V. Douglas, J. Lillis, "Health informatics: managing information to deliver value", PubMed indexed for MEDLINE, Stud Health Technol Inform.;84(Pt1):305-8., 2001.
- [9] J.A. Quinn: "Bayesian Condition Monitoring in Neonatal Intensive Care" PhD thesis submitted to Institute for Adaptive and Neural Computation School of Informatics University of Edinburgh, 2007.
- [10] S. Nicklin, Y.A. Wickramasinghe and S.A. Spencer, "Neonatal intensive care monitoring current paediatrics"; Vol 14(1), pp. 1-7, 2004.
- [11] STELLA Newsletter, "2nd International Workshop on Flexible and Stretchable Electronics" Content Issue No. 6, 2010.
- [12] L. Suresh, A.N. Latha A.N, R.B. Murthy, K.T. Alam, and J.K. Babu, Neonatal monitoring system Int. Journal of Engineering Research and Applications, ISSN : 2248-9622, Vol. 4(7), pp. 12-15, 2014.
- [13] P.J. Pawar, B.V. Val, F. Bert-Jan and H. Hermens, "A framework for the comparison of mobile patient monitoring systems". Biomedical Journal, 2012.
- [14] I.R. Okonkwo, B.I. Abhulimhen-Iyoha and A.A. Okolo, "Scope of neonatal care services in major Nigerian hospitals", Niger J Paed Vol 43(1), pp.8-13, 2016.
- [15] A.A. Sobowale, S.O. Olabiyisi, and T.A. Abdul-Hameed, Development of a Framework for Computerized Health Management Information Systems in Nigeria: International Journal of Information and Communication Technology Research, 2011.
- [16] Mednax Services, "Pediatrx Medical Group: For Parents, Your Baby and the NICU" Important Information from Your Health Care providers through The Centre for Research, Education and Quality. www.pediatrx.com/forparents , 2011.
- [17] Ye Y. and Tong S. J A Knowledge-Based Variance Management System for Supporting the Implementation of Clinical Pathways.", Management and Service science, 2009, IEEE, pages 1 -4 (2009).
- [18] Warren J, Beliakov G and Zwaag B. " Fuzzy logic in clinical practice Decision support system", Proceedings of the 33rd Hawaii International Conference on System Sciences. (2000).
- [19] Prabhu M., Senthil P. N. and Lakshmi K., Clinical Decision Support Systems" Computer Sciences Corporation (CSC), pages 1-19 accessed in June, 2016.

# Machine Learning-Based Phishing Attack Detection

Sohrab Hossain<sup>1</sup>, Dhiman Sarma<sup>2\*</sup>, Rana Joyti Chakma<sup>3</sup>

Department of Computer Science and Engineering, East Delta University, Chittagong, Bangladesh<sup>1</sup>

Department of Computer Science and Engineering, Rangamati Science and Technology University, Rangamati, Bangladesh<sup>2,3</sup>

**Abstract**—This paper explores machine learning techniques and evaluates their performances when trained to perform against datasets consisting of features that can differentiate between a Phishing Website and a safe one. This capability of telling these sites apart from one another is vital in the modern-day internet surfing. As more and more of our resources shift online, one vulnerability and a leak of sensitive information by someone could bring everything down in a connected network. This paper's objective through this research is to highlight the best technique for identifying one of the most commonly occurring cyberattacks and thus allow faster identification and blacklisting of such sites, therefore leading to a safer and more secure web surfing experience for everyone. To achieve this, we describe each of the techniques we look into in great detail and use different evaluation techniques to portray their performance visually. After pitting all of these techniques against each other, we have concluded with an explanation in this paper that Random Forest Classifier does indeed work best for Phishing Website Detection.

**Keywords**—Phishing attack; phishing attack detection; phishing website detection; machine learning; random forest classifier

## I. INTRODUCTION

Phishing Attacks are the most common ways of attack in the digital world these days. Any method of communication can be used to target an individual to trick them into leaking confidential data in a fake environment, which can later be used to harm the sole victim or even an entire business depending on the attacker's intent and the type of data leaked.

Phishing attacks, while dangerous, can be avoided by simply creating awareness and developing habits of staying alert and continuously being on the lookout when surfing through the internet and only clicking links after verifying if the source of the links is trustworthy at all. There are also tools such as browser extensions that notify users when they have entered their credentials on a fake site, possibly having their credentials transferred to a user with malicious intent. Other tools can also allow networks to lock down everything and allow access to whitelisted sites to provide extra security while compromising some convenience on the user side [1].

In a related study, five main reasons have been stated behind users falling into traps of phishing attack schemes:

- Lack of knowledge about URLs.
- Lack of knowledge about trusted websites.
- Lack of visibility of full web addresses due to the redirection or hidden URLs.

- Lack of time for analyzing URLs, and accidental entries of some web pages.
- Lack of capability of telling phishing web pages apart from legitimate ones.

One example of such an attack would be the attack in 2016, known as the Bangladesh Bank Cyber Heist. Security Hackers issued thirty-five fraudulent instructions via the SWIFT network to illegally transfer almost 1 billion US dollars from the Federal Reserve Bank of New York account that belonged to Bangladesh Bank. Out of these 35 instructions, 5 of them successfully transferred 101 million dollars, with 20 million traced to Sri Lanka and 81 million traced to the Philippines. Fortunately, the Federal Reserve Bank of New York was able to block the remaining thirty transactions. Without this block, another 850 million dollars would have been lost. And it was possible all thanks to noticing a misspelled instruction that raised suspicions among the authorities. The money transferred to Sri Lanka was all recovered, but from the US\$ 81 million transferred to the Philippines, only US\$ 18 million was recovered. Most of the money transferred to the Philippines were collected into four personal accounts [2].

The method of this attack has been suspected to be a Dridex malware. It specializes in stealing bank credentials by using macros set up in a Word or Excel document. Windows users can fall victim to such an attack if they open email attachments in Word or Excel, containing such a macro, which once activated on opening these documents, begin downloading Dridex, which then infects computers and sets up the stage for a banking theft. A knowledgeable and alert employee or a software aiding in detecting such an attack would have helped immensely in this event [3].

Machine learning algorithms are widely used to detect hidden patterns in the dataset. The most common algorithms are K-nearest neighbor, decision trees, random forest, and support vector machine [4]. In addition, belief rule-based expert system can mine rules from the dataset [5] [6].

In this paper, we focus on training machine learning models that can detect phishing web pages apart from real web pages. We analyze each of these models and state our findings and research in this paper to allow for others to have a clear understanding of the performance of these models when trained for this purpose. Of course, data preprocessing is very crucial for the models to work as they did in our case, and that is an essential part of the procedure. Papers from other researchers contributed immensely to our research, and we hope our paper will do the same by providing a collection of our findings regarding Phishing Detection using Machine Learning in this paper.

\*Corresponding Author

The remaining of the paper is organized as follows. In Section II, we reviewed the literature, followed by presenting the proposed methodology in Section III. The empirical results of the proposed approach are explained in Section IV, followed by Section V where a conclusion and further research scopes are discussed.

## II. LITERATURE REVIEW

### A. Types of Phishing Attacks

1) *Algorithm-Based phishing*: Attackers access sensitive information from a website's database by employing different algorithms. V. Shreeram, M. Suban, P. Shanthi, K. Manjula proposed an anti-phishing detection method that would detect phishing hyperlinks with the help of the rule-based system that is formulated from the genetic algorithm (GA). A phishing link is detected if it matches the ruleset that is created by GA, which is stored in a database [7].

2) *Deceptive phishing*: This technique involves supplying clients with malicious links via emails and redirecting them to malicious websites where they are likely to enter sensitive information. Huajun Huang, Junshan Tan, Lingxi Liu gives a thorough overview of a deceptive phishing attack and different anti-phishing techniques. They present the different methods used by phishers and the advantages and disadvantages of the different countermeasures used [8].

3) *URL phishing*: Hackers can inject hidden links that redirect to malicious pages into the URL, where one may not expect to find one. Mohammed Nazim Feroz, Susan Mengel, proposes a method to detect URL phishing with URL ranking. They classify the URLs by their lexical and host-based features and categorizes and rank the URLs using the online URL reputation services [9].

4) *Hosts file poisoning*: Replacing hostnames in the host records can override the usual process of DNS servers trying to retrieve actual IP addresses from beyond the network. This technique can poison the records and allow valid URLs that are meant to lead to secure sites lead to malicious pages instead, due to compromised IP associations in the server. Saeed Abu-Nimeh, Suku Nair, proposes a new attack that can bypass security toolbars and phishing filters by using DNS poisoning. They use spoofed DNS cache entries to create fake results and successfully attack four renowned security toolbars and the phishing filters of three popular browsers without being detected [10].

5) *Content injection phishing*: Data collection is achieved in this technique by concatenation of malicious sections within a real website. Jussi-Pekka Erkkil presents the different methods by which phishing techniques can trick a person. A list of several strategies is listed that can detect phishing. The paper proposes that the company adapt effective protocols to keep their security features up to date [11].

6) *Clone phishing*: Duplicating already sent emails and attaching a malicious link into it can allow for a successful attack on an unsuspecting user. Ahmad Alamgir Khan proposed a new method where websites use One Time

Password and User-machine Identification system to combat phishing attacks. Webservers will send a one time password to a user by SMS or email and create an encrypted token for the device after the user inputs the password [12].

### B. Phishing Website Detection Techniques

1) *Blacklist filter*: Blacklists can be maintained to block recorded unwanted sites from reaching the client's machine. These filters can be applied in different security measures like DNS servers, firewalls, email servers, etc. A blacklist filter maintains a list of elements like IP addresses, domains, IP netblocks that are commonly used by phishers. Adam Oest, Yeganeh Safaei, Adam Doupé, Gail-Joon Ahn, Brad Wardman, Kevin Tyers uses a scalable framework to test the effectiveness of browser blacklist filters. Their study concluded that most blacklist filters in mobile browsers failed to combat phishing attacks and are more vulnerable [13]. Mohsen Sharifi, Seyed Hossein Siadati, proposes a new method that will create a blacklist generator and keep a timely track of phishing website blacklists. Their techniques yield an accuracy of 91% and 100% in detecting real pages and phishing websites, respectively [14].

2) *Whitelist filter*: Unlike a Blacklist, Whitelist filters allow recorded website URLs, schemes, or domains to make it through to the client machine and block all other unrecorded sites. A whitelist, contrary to a blacklist, maintains a list of all legitimate websites. A. Belabed, E. Aïmeur, A. Chikh proposes a method that combines the whitelist approach with machine learning. A support vector classifier is used to filter further the websites that are not blocked by the whitelist filter [15]. Linfeng Li, Marko Helenius, and Eleni Berki conducted tests that compared the effectiveness of blacklist and whitelist anti-phishing toolbars. Their study did not find a significant difference in performance between both toolbars but encourages that toolbars be more instructive in helping users identify phishing websites [16].

3) *Pattern matching filter*: Checks whether or not individual tokens or sequences of data is contained within a given list of data by using a pattern matching technique. Rahamathunnisa Usuff, N. Manikandan, U.S. Kumaran, and C. Niveditha propose a method that uses pattern matching to detect phishing websites. A database of blacklist and whitelist that contains malicious URL patterns and original URL patterns is used to match with the user requested URL [17].

### C. Machine Learning-Based Methods

1) *Malicious domain detection*: Machine Learning models are being trained to optimize their capabilities of detecting Phishing pages, one of the most common forms of cyberattacks. Nitay Hason, Amit Dvir, and Chen Hajaj propose a robust feature selection mechanism that creates better malicious domain detection models. All of the data are collected from 5000 legitimate URLs and 1350 harmful URLs. The models created are robust to different malicious abnormalities and show the effectiveness of models trained on features [18]. Hossein Shirazi, Bruhadeshwar Bezawada,

Indrakshi Ray shows concern about the large number of training features and types of datasets used and suggests that the domain name is much better and useful detecting phishing websites. Their learning model detects unknown live phishing URLs with an accuracy of 99.7% [19]. Krzysztof Lasota, Adam Kozakiewicz proposes a study that shows the similarity of different malicious domain name creations. The main task for detecting malicious behaviors was to detect similarity based on sets of domain names, URL names, and hostnames [20].

2) *Email spam filtering*: Emails are screened through various scoring techniques based on thousands of rules set to predict their probability of being an actual spam email. If the evaluated probability is beyond the acceptable range, then the email is blocked via the spam filter. Phishers use spam emails to direct a client to their malicious webpage and steal data. Andronicus A. Akinyelu and Aderemi O. Adewumi research about the effectiveness and use of random forest classifier in

developing a phishing email classifier by extracting pertinent phishing email features from a dataset of 2000 phishing and ham emails. The proposed machine learning models shows a classification accuracy of 99.7% with low false positives and negatives [21]. Tushaar Gangavaraapu, C.D. Jaidhar, and Bhabesh Chanduka focus on the proper ways of extracting features from spam email content and behavior-based features, the features necessary in detecting spam emails, and on the selection of an important feature set. Their proposed machine learning model based on their selected features yields a constant accuracy of 99% in spam emails [22]. Table I illustrate the advantages and limitations of existing phishing detection researches. In Table I, we observed that most of the researches consider a small number of features and datasets. In this research, we try to overcome the limitations observed from Table I by increasing the number of features and dataset volume.

TABLE I. COMPARISON OF MACHINE LEARNING BASED PHISHING DETECTION SYSTEMS

| Description  | Pros  | Cons  | Ref. |
|--|---|---|------|
| Detects phishing attacks by using a whitelist filter.  | <ul style="list-style-type: none"> <li>* Pages that bypass the whitelist filter are filtered again by Support Vector Machines.</li> <li>* Maintains accuracy of whitelist filter by using a personalized whitelist.</li> </ul>                                | <ul style="list-style-type: none"> <li>* Limited dataset of 850 pages.</li> <li>* Unable to detect the attachment of DNS spoofs to legitimate web pages.</li> <li>* High False positive rate.</li> </ul>  | [23] |
| Implement a comment spam detection mechanism that can be used as a browser plugin and remove spam comments.  | <ul style="list-style-type: none"> <li>* Balances dataset by applying WEKA filters to get the best suitable features.</li> <li>* Spam detection classifier can accommodate new features and detect new classes of spam content.</li> </ul>                    | <ul style="list-style-type: none"> <li>* Does not do well with a random dataset without applying a supervised resample filter.</li> </ul>   | [24] |
| Proposes a machine learning-based method that can detect whether a web page exhibits phishing attacks.   | <ul style="list-style-type: none"> <li>* Proposed method is based on an easy to acquire feature vector that does not require additional computation.</li> </ul>   | <ul style="list-style-type: none"> <li>* Only uses 10 features for detection.</li> <li>* Limited dataset of 1353 instances.</li> </ul>  | [25] |
| Uses feature selection to identify important features that categorize phishing and legitimate websites.  | <ul style="list-style-type: none"> <li>* Feature selection highly improves the accuracy score after implementation.</li> <li>* Use of feature selection reduces computational time.</li> </ul>  | <ul style="list-style-type: none"> <li>* 14 features.</li> <li>* limited dataset (200 legitimate URL and 1400 phishing URL)</li> <li>* May not work properly with datasets of equal URLs of legitimate and phishing web pages.</li> </ul>           | [26] |
| Builds a system using machine learning that can classify websites using URLs.  | <ul style="list-style-type: none"> <li>* Can be used to build a rule-based system with associative rules to classify URLs.</li> </ul>   | <ul style="list-style-type: none"> <li>* 9 features for each URL</li> <li>* All features are discrete.</li> <li>* Limited dataset (1353 URLs)</li> </ul>  | [27] |
| Proposes a learning-based aggregation analysis mechanism to decide page layout similarity, which is used to detect phishing pages.                         | <ul style="list-style-type: none"> <li>* Automatically trains classifiers to determine web page similarity from CSS layout features, which does not require human expertise.</li> </ul>   | <ul style="list-style-type: none"> <li>* Method is lightweight as it only takes one class of features, CSS structure.</li> <li>* Limited by the size of the dataset and distribution of samples.</li> </ul>   | [28] |
| This research uses a new attribute called the "domain top page similarity" to improve the efficiency of a machine learning-based phishing detection model. | <ul style="list-style-type: none"> <li>* Increases f-measure and reduces the error rate.</li> <li>* Proves that with better features, the detection rate is much higher and can be implemented in future works.</li> </ul>                                    | <ul style="list-style-type: none"> <li>* The model is highly dependent on the accuracy of the features.</li> </ul>  | [29] |
| This paper proposes a real-time anti-phishing system that uses seven classification algorithms and natural language processing-based features (NLP)        | <ul style="list-style-type: none"> <li>* Independence from language and third party services.</li> <li>* Huge dataset of legitimate and phishing data.</li> <li>* Real-time execution.</li> <li>* Can detect new websites because of NLP features.</li> </ul> | <ul style="list-style-type: none"> <li>* Machine learning-based systems cannot correctly utilize such a vast dataset.</li> </ul>  | [30] |
| Performs an extensive measurement of squatting phishing, where the phishing pages impersonate target brands at both the domain and content level.          | <ul style="list-style-type: none"> <li>* Uses features from visual analysis and optical character recognition.</li> <li>* Open sourced tool.</li> <li>* Uses evasive behaviors of phishing pages to build classifiers.</li> </ul>                             | <ul style="list-style-type: none"> <li>* Unable to detect phishing pages that use cloaking.</li> <li>* Only focuses on popular brands.</li> <li>* The classifier cannot be compared with other phishing tools like CANTINA and CANTINA+.</li> </ul> | [31] |
| Uses features from HTML content, JavaScript code and URLs to build a classifier that can detect malicious web pages and threat types.                      | <ul style="list-style-type: none"> <li>* Diverse features.</li> <li>* High accuracy score.</li> <li>* Highlights features that are necessary to extract.</li> </ul>   | <ul style="list-style-type: none"> <li>* Limited dataset (2500 URLs)</li> <li>* Classifier may not do well with large datasets.</li> </ul>  | [32] |

### III. PHISHING WEBSITE DETECTION

In this section, we explain our proposed data-driven phishing website detection system—the dataset obtained from the online repository of Mendeley. Parallel coordinates, pearson and shapiro ranking, and principal component analysis are used for feature extraction. We use KNN, decision trees, random forest, SVM, and logistic regression to detect phishing websites.

#### A. Dataset

The phishing webpage dataset contains 48 features that are obtained from the online repository of Mendeley. The total number of websites is 1000, where 5000 phishing and 5000 legitimate websites. The class label 0 indicates a phishing website and 1 a legitimate website.

#### B. Feature Extraction and uses

For feature extraction, we used parallel coordinates, pearson and shapiro ranking, and principal component analysis. We used parallel coordinates to visualize and analyze our dataset and PCA to reduce the dimensionality of our dataset. We have explained our features in Table I, Table III, and Table IV. In Table II, a total number of 27 lexical features are described like NumDots, SubdomainLevel, PathLevel, and so on. A total number of 15 host-based features are explained in

Table III. In Table IV, a set of 8 correlation features are shown with data types and description.

#### C. Classifiers

We deploy KNN, decision tree, random forest, extra trees, SVM, and logistic regression in our system.

1) *K-Nearest Neighbors (KNN)*: We calculated the distance using the Euclidean method from equation (1),

$$d(x, x') = \sqrt{(x_1 - x'_1)^2 + \dots + (x_n - x'_n)^2} \quad (1)$$

Our KNN model is based on equation (2),

$$P(y = j | X = x) = \frac{1}{K} \sum_{i \in A} I(y^{(i)} = j) \quad (2)$$

Our dataset has 48 features and a Class label where 0 indicates a phishing website, and 1 indicates a legitimate website. When given an unknown sample, KNN will first measure the distance of the unknown sample with its neighbors by using Euclidean distance. The number of neighbors that it will check will be the value of K that can be chosen by setting the value of "n\_neighbors." The distances will be measured by taking in the features of the samples that are in the dataset. The majority class of the neighbors that are the closest will be then assigned to the unknown sample.

TABLE II. LIST OF URL FEATURES IN LEXICAL FEATURE GROUP

| Feature             | Data Type | Description  |
|---------------------|-----------|--|
| NumDots             | Numeric   | The number of dots In the URL.   |
| SubdomainLevel      | Numeric   | Determines the number of subdomain levels.   |
| PathLevel           | Numeric   | Determining the level of the path in the URL.  |
| UrlLength           | Numeric   | Length of each URL used in the dataset. The length contains the number of letters or symbols used to create the URL. |
| NumDash             | Numeric   | Total number of dash in a URL.   |
| NumDashInHostname   | Numeric   | The number of dashes in a hostname   |
| AtSymbol            | Boolean   | Total number of '@' symbol in the URL.   |
| TildeSymbol         | Boolean   | Total number of tilde '~' symbol in the URL.   |
| NumUnderscore       | Numeric   | Number of underscores '_' used in the URL.   |
| NumPercent          | Numeric   | Total number of percent symbol present in the URL.   |
| NumQueryComponents  | Numeric   | Total number of query components.  |
| NumAmpersand        | Numeric   | Total number of '&' character.   |
| NumHash             | Numeric   | Total number of '#' character.   |
| NumNumericChars     | Numeric   | The total number of numeric characters.  |
| NoHttps             | Boolean   | Check if there is a HTTPS in the URL.  |
| RandomString        | String    | Set of Characters that are random.   |
| IPAddress           | Boolean   | Check if the hostname of the URL uses the IP address.  |
| DomainsInSubDomains | Boolean   | Determines if TLD or CCTLD is in the subdomain of URL.   |
| DomainsInPaths      | Boolean   | Determines if the website link has used TLD or CCTLD.  |
| HttpsInHostname     | Boolean   | Determines if HTTPS is disorderly in the hostname of the URL.  |
| HostnameLength      | Numeric   | Length of hostname which includes all the characters and symbols.  |
| PathLength          | Numeric   | Length of all paths in each URL.   |
| QueryLength         | Numeric   | Length of query in the URL.  |
| DoubleSlashInPath   | Boolean   | Checks if there is a double slash in the path.   |
| NumSensitiveWords   | Numeric   | Checks if there are any sensitive words like secure, sign in, login, etc.  |
| EmbeddedBrandName   | Boolean   | Checks if there is the name of a brand in the domain.  |
| PctExtHyperLinks    | Float     | Checks the percentage of external hyperlinks in the source code.   |



TABLE III. LIST OF URL FEATURE IN THE HOST-BASED FEATURE GROUP

| Feature                       | Data Type | Description   |
|-------------------------------|-----------|---|
| PctExtResourceUrls            | Float     | Checks the percentage of URL external resources in the source code.                                     |
| ExtFavicon                    | Boolean   | Checks if the favicon is installed from a different hostname.   |
| InsecureForms                 | Boolean   | Will see if the action in forms follow the HTTPS protocol.  |
| RelativeFormAction            | Boolean   | Checks if the action form contains a relative URL.  |
| ExtFormAction                 | Boolean   | Checks if the action form contains an external URL.   |
| AbnormalFormAction            | Boolean   | Checks if the action form contains an abnormal URL.   |
| PctNullSelfRedirectHyperlinks | Float     | Check the percentage of hyperlinks that have an empty value and also if it has an auto directing value. |
| FrequentDomainNameMismatch    | Boolean   | Checks if the URL, when accessed, shows a mismatch in the frequent domain name.                         |
| FakeLinkInStatusBar           | Boolean   | Checks if there are any fake link in status bar that lures the user towards unsafe websites.            |
| RightClickDisabled            | Boolean   | Check if the right-click option has been disabled in the URL.   |
| PopUpWindow                   | Boolean   | Checks if the URL contains any pop up windows when opened or accessed.                                  |
| SubmitInfoToEmail             | Boolean   | Checks whether a URL requires you to submit your information to email.                                  |
| IframeOrFrame                 | Boolean   | Check if the given URL has used iframes or frames.  |
| MissingTitle                  | Boolean   | Check if there are any missing title.   |
| ImagesOnlyInForm              | Boolean   | Checks if there are only images in the action form.   |

TABLE IV. LIST OF URL FEATURES IN CORRELATED FEATURE GROUP

| Feature                            | Data Type | Description  |
|------------------------------------|-----------|--|
| SubdomainLevelRT                   | -1, 0, 1  | Checks if the subdomain levels are correlated.                         |
| UrlLengthRT                        | -1, 0, 1  | Checks if the URL lengths are correlated.                              |
| PctExtResourceUrlsRT               | -1, 0, 1  | Checks if the percentage of external URL is correlated.                |
| AbnormalExtFormActionR             | -1, 0, 1  | Checks the relationship of different abnormal action forms in the URL. |
| ExtMetaScriptLinkRT                | -1, 0, 1  | Checks the correlation of meta script links                            |
| PxtExtNullSelfRedirectHyperlinksRT | -1, 0, 1  | Checks the correlation of the percentage of self-directed hyperlinks.  |
| Class_label                        | 0, 1      | Identifying the 2 classes of Phishing and Real Website.                |

2) *Random forest*: We used Gini importance to calculate a node's importance for each decision tree. This was based under the assumption that the tree is binary, and so each node has at most two children. For the elimination of branches in the tree, we used the equation (3),

$$n_{ij} = w_j C_j - w_{left(j)} C_{left(j)} - w_{right(j)} C_{right(j)} \quad (3)$$

For calculating the importance of each feature on a decision tree, we used the equation (4),

$$f_i = \frac{\sum_{j:node\ j\ splits\ on\ feature\ i} n_{ij}}{\sum_{k:all\ nodes} n_{ik}} \quad (4)$$

These can be normalized afterward to a value between 0 and 1 by the equation (5),

$$norm f_i = \frac{f_i}{\sum_{j:all\ features} f_{ij}} \quad (5)$$

And the sum of the feature's importance value on each tree is calculated by the equation (6) and divided by the total number of trees.

$$RF f_i = \frac{\sum_{j:all\ trees} norm f_{ij}}{T} \quad (6)$$

A random forest classifier consists of a large number of decision trees that work as an ensemble. At first, it will create a bootstrap dataset of size "N" that will randomly take samples from our dataset. A random forest can then use these bootstrap samples to create a tree. For example, if our training data was [a, b, c, d, e, f], we might give one of our trees the following list [a, b, b, c, f, f]. It should be noticed that both samples are of the same size, and "b" and "f" are repeated in the bootstrap dataset because we sample with replacement. After taking in the samples from the bootstrap dataset, it begins to build trees by first choosing a root node. Random forest differs from decision trees because it uses a method called Feature Randomness. This means that when it comes to choosing a root node for a random tree forest will only allow the trees to choose a root node from a subset of features. The Gini impurity is measured among these subsets of features, and the lowest score will be used as the root node, and the different subsequent nodes are chosen in the same way. After creating the trees, the random forest classifier is ready to make predictions. It will take an unknown sample from our test dataset and run the sample among all of the trees. All of the individual trees give a class prediction, and the class that has the most votes will be the class of the unknown sample. One of the main reasons random forest classifier does well with large

datasets is because it maintains diversity between models by using bootstrap aggregation and feature randomness.

3) *Support vector machines*: We used the equation (7) to calculate the loss function for our support vector machine,

$$\min_w \lambda \left( \|w\|^2 + \sum_{i=1}^n (1 - y_i \langle x_i, w \rangle)_+ \right) \quad (7)$$

For calculating gradients, we used the equation (8),

$$\frac{\partial}{\partial w_k} \lambda \left( \|w\|^2 + \sum_{i=1}^n (1 - y_i \langle x_i, w \rangle)_+ \right) = \begin{cases} 2w_k, & \text{if } y_i \langle x_i, w \rangle \geq 1 \\ -y_i x_{ik}, & \text{else} \end{cases} \quad (8)$$

By using SVM, we plot each data item as a point in  $n$ -dimensional space (where  $n$  is the number of features and in our dataset it is 48) with the value of each feature being a value of a specific coordinate. After that SVM finds a hyperplane or a decision boundary that can properly differentiate between the classes. An optimal hyperplane is one where it has equal and maximum distance between two data points, which are considered as support vectors. SVM is very easy to apply when the data points can be easily divided by a linear line, but it is rare to find such datasets in the real world. This is where the kernel trick of SVM comes to work. One of the reasons why SVM works well with our large dataset is that it can work in infinite dimensions. The best part is that the kernel does not necessarily generate the infinite dimensions but simulates the lower dimension data so as if they are working in infinite dimensions. The kernel is very useful here because it can make a non-separable problem into a separable problem by adding more dimensions to it, and the number of dimensions depends on the number of features each sample has; some of the kernels that we found compelling are Linear Kernel, Polynomial Kernel, and the Radial Basis Function (RBF) kernel.

4) *Logistic regression*: Logistic regression is based on the linear regression, where a line is plotted its axes for a given dataset.

The conditional probability function we used gives a binary output for the variable  $Y$  as a function of  $X$ . Any unknown parameters in the function are estimated by maximum likelihood. The conditional probability is calculated by using equation (9).

$$\Pr(Y = 1 | X = x) = \log \frac{p(x)}{1-p(x)} = \beta_0 + x \cdot \beta \quad (9)$$

We also used equation (10) for the sigmoid function,

$$S(x) = \frac{1}{1+e^{-x}} = \frac{e^x}{e^x+1} \quad (10)$$

Equation (11) is the cost function,

$$-\frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} \log h_0(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_0(x^{(i)})) \right] \quad (11)$$

We calculate the gradient by using the equations (12), (13), (14), and (15).

$$J = -\frac{1}{m} \left[ \sum_{i=1}^m y_i \log h_i + (1 - y_i) \log 1 - h_i \right] \quad (12)$$

$$\frac{\partial J}{\partial \theta_n} = -\frac{1}{m} \cdot \left[ \sum_{i=1}^m \frac{y_i}{h_i} \cdot h_i^2 \cdot x_n \cdot \frac{1-h_i}{h_i} + \frac{1-y_i}{1-h_i} \cdot -h_i^2 \cdot x_n \cdot \frac{1-h_i}{h_i} \right] \quad (13)$$

$$\frac{\partial J}{\partial \theta_n} = -\frac{1}{m} \cdot \left[ \sum_{i=1}^m x_n \cdot (1 - h_i) \cdot y_i + x_n \cdot h_i \cdot (1 - y_i) \right] \quad (14)$$

$$\frac{\partial J}{\partial \theta_n} = \frac{1}{m} \cdot x_i \cdot \left[ \sum_{i=1}^m h_i - y_i \right] \quad (15)$$

## IV. RESULT ANALYSIS

### A. ROC Curve

Now let us look at our ROC curves of different models.

Fig. 1 shows the ROC curve of the support vector machine. The X-axis indicates the false positive rate, and the Y-axis indicates the True positive rate. The AUC value for this is 0.97.

Fig. 2 shows the ROC curve of the non-uniform support vector machine. The X-axis indicates the false positive rate, and the Y-axis indicates the True positive rate. The AUC value for this is 0.96.

Fig. 3 shows are the ROC curve of the linear support vector machine. The X-axis indicates the False Positive rate, and the Y-axis indicates the True positive rate. The AUC value for this is 0.98. This is the highest and best one so far. We can see the steepness in the curve is much closer to the top-left position of the plot.

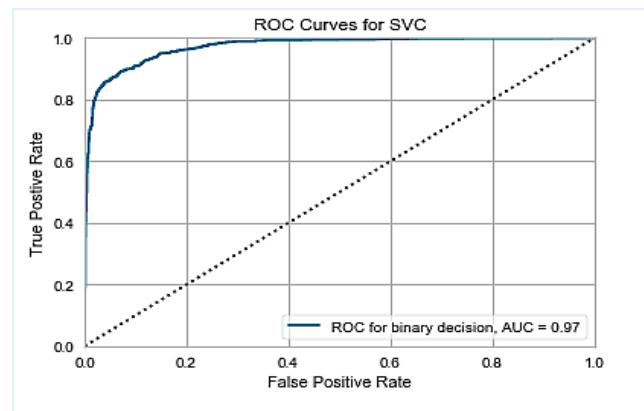


Fig. 1. Curves for SVC.

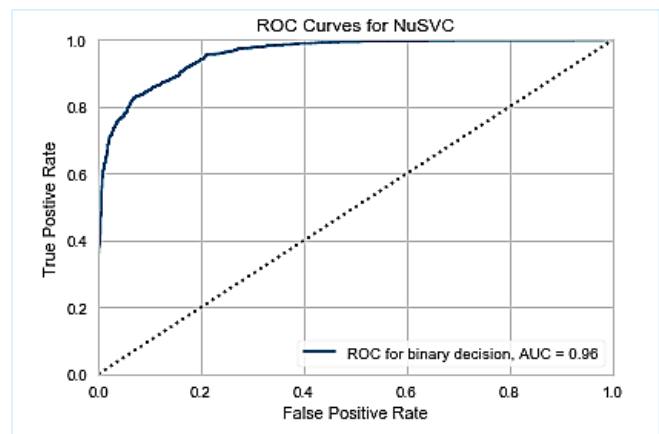


Fig. 2. ROC Curves for NuSVC.

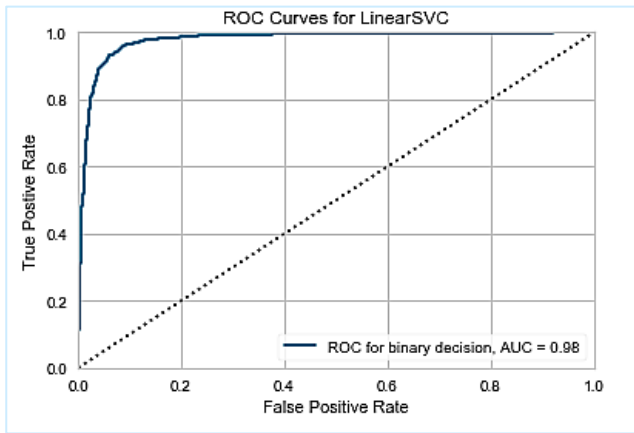


Fig. 3. ROC Curves for LinearSVC.

Fig. 4 shows the ROC curve of KNN. The X-axis indicates the false positive rate, and the Y-axis indicates the True positive rate. Here the AUC of class 0 (phishing website) is 0.94, and class 1 (real website) is 0.94. The AUC of the macro and micro average of the ROC curve is also 0.94.

Fig. 5 shows the ROC curve of Logistic Regression. The X-axis indicates the False Positive rate, and the Y-axis indicates the True positive rate. Here the AUC of class 0 (phishing website) is 0.96, and class 1 (real website) is 0.96. The AUC of the macro and micro average of the ROC curve is also 0.96.

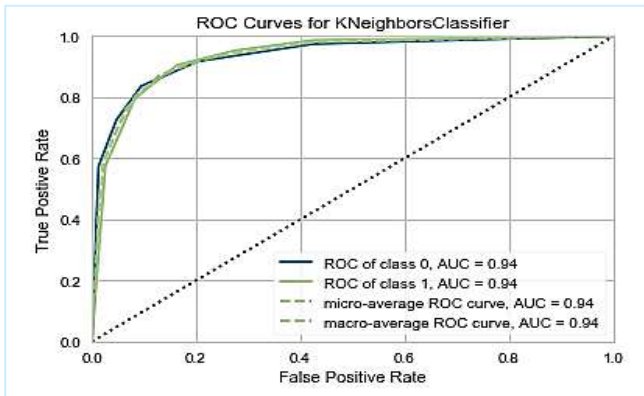


Fig. 4. ROC Curves for KNeighborsClassifier.

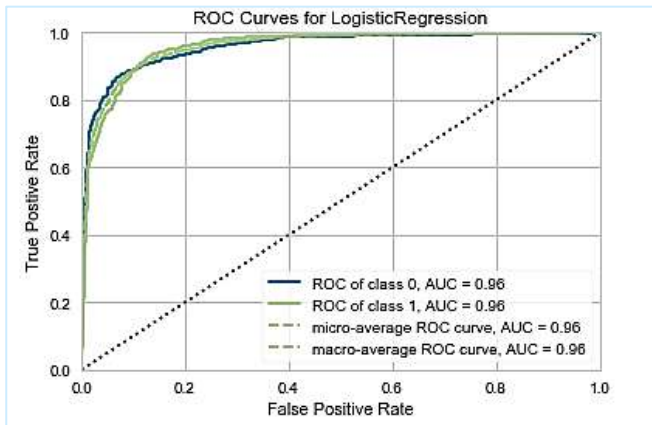


Fig. 5. ROC Curves for LogisticRegression.

Fig. 6 shows the ROC curve of stochastic gradient descent (SGD). The X-axis indicates the false positive rate, and the Y-axis indicates the True positive rate. The AUC value for this is 0.97.

Fig. 7 shows the ROC curve of logistic regressionCV. The X-axis indicates the false positive rate, and the Y-axis indicates the True positive rate. Here the AUC of class 0 (phishing website) is 0.98, and class 1 (real website) is 0.98. The AUC of the macro and micro average of the ROC curve is also 0.98.

Fig. 8 shows the ROC curve of the bagging classifier. The X-axis indicates the false positive rate, and the Y-axis indicates the True positive rate. Here the AUC of class 0 (phishing website) is 0.99, and class 1 (real website) is 0.99. The AUC of the macro and micro average of the ROC curve is also 0.99. This is the best ROC curve so far.

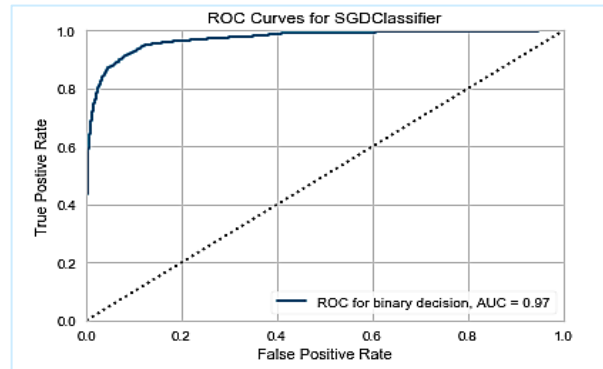


Fig. 6. ROC Curves for SGD Classifier.

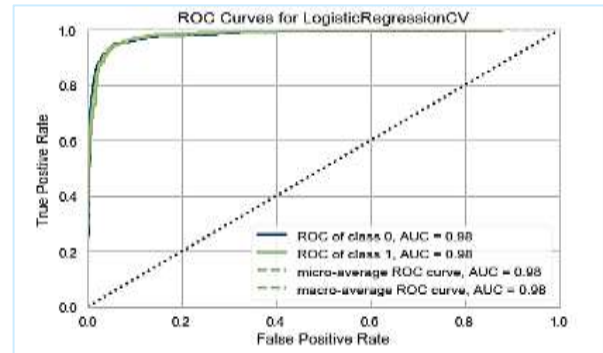


Fig. 7. ROC Curves for Logistic Regression CV.

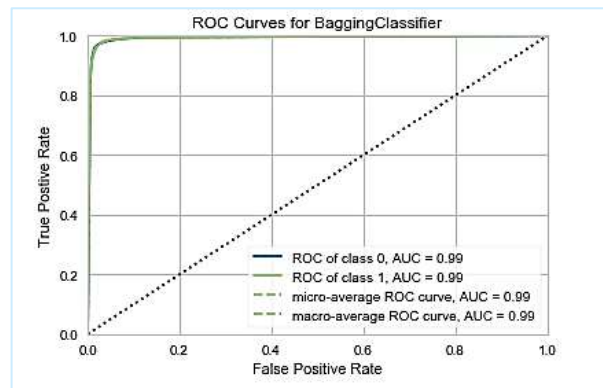


Fig. 8. ROC Curves for Bagging Classifier.

Fig. 9 shows the ROC curve of the extra trees classifier. The X-axis indicates the False Positive rate, and the Y-axis indicates the True positive rate. Here the AUC of class 0 (phishing website) is 1.00, and class 1 (real website) is 1.00. The AUC of the macro and micro average of the ROC curve is also 1.00. This is the best ROC curve so far. We can see that the steepness of the curve is at the most top left corner.

Fig. 10 shows the ROC curve of the random forest classifier. The X-axis indicates the False Positive rate, and the Y-axis indicates the True positive rate. Here the AUC of class 0 (phishing website) is 1.00, and class 1 (real website) is 1.00. The AUC of the macro and micro average of the ROC curve is also 1.00. This is the same as the Extra Trees classifier. We can see that the steepness of the curve is at the most top left corner. Hence it can be said that the extra trees classifier and random trees classifier has the best ROC curve.

**B. Discrimination Threshold**

Let us look at the discrimination threshold of our models.

Fig. 11 shows the threshold plot for the support vector machine. On the X-axis, we have the discrimination threshold, and on the Y-axis, we have the score. Here we see that the discrimination threshold for this is 0.03. For this threshold, we see that the precision, recall, and f1 score are approximately around 0.89.

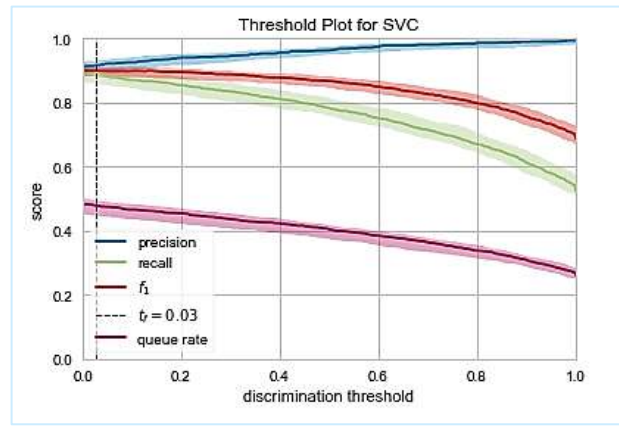


Fig. 11. Threshold Plot for SVC.

Fig. 12 shows the threshold plot for the non-uniform support vector machine. On the X-axis, we have the discrimination threshold, and on the Y-axis, we have the score. Here we see that the discrimination threshold for this is 0.00. For this threshold, we see that the precision, recall, and f1 score are approximately around 0.86.

Fig. 13 shows the threshold plot for the linear support vector machine. On the X-axis, we have the discrimination threshold, and on the Y-axis, we have the score. Here we see that the discrimination threshold for this is 0.05. For this threshold, we see that the precision, recall, and f1 score are approximately around 0.9.

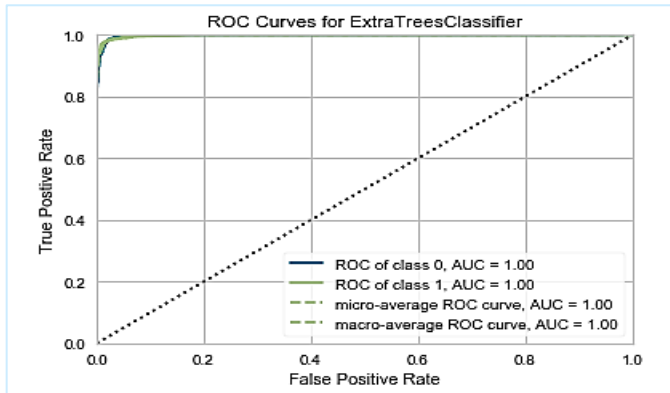


Fig. 9. ROC Curves for ExtraTrees Classifier.

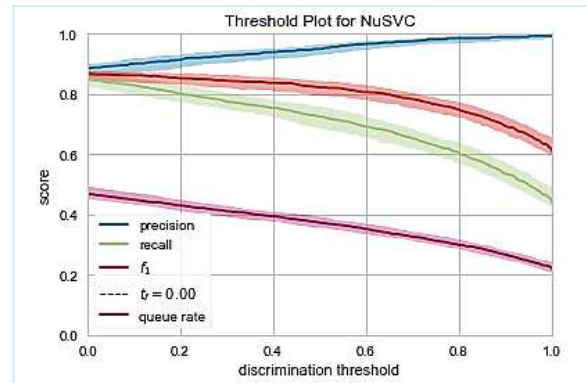


Fig. 12. Threshold Plot for NuSVC.

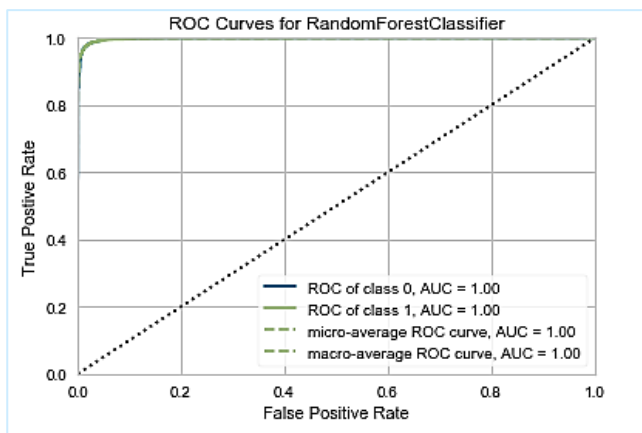


Fig. 10. ROC Curves for Random Forest Classifier.

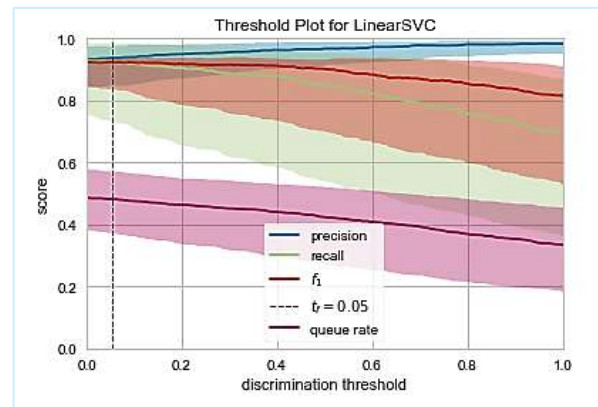


Fig. 13. Threshold Plot for Linear SVC.

Fig. 14 shows the threshold plot for KNN. On the X-axis, we have the discrimination threshold, and on the Y-axis, we have the score. Here we see that the discrimination threshold for this is 0.50. For this threshold, we see that the precision, recall, and f1 score are approximately 0.82 to 0.89.

Fig. 15 shows the threshold plot for logistic regression. On the X-axis, we have the discrimination threshold, and on the Y-axis, we have the score. Here we see that the discrimination threshold for this is 0.46. For this threshold, we see that the precision, recall, and f1 score are approximately 0.85 to 0.9.

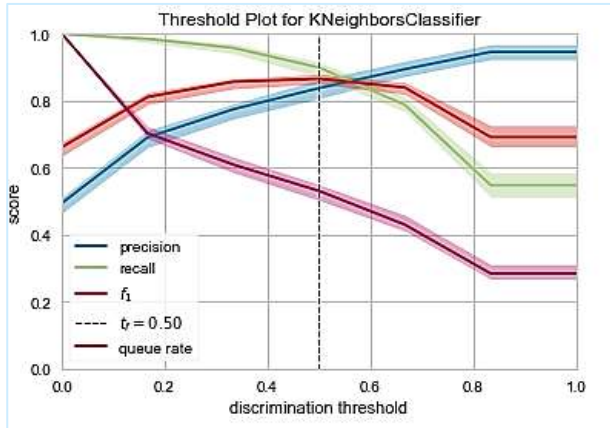


Fig. 14. Threshold Plot for KNeighbors Classifier.

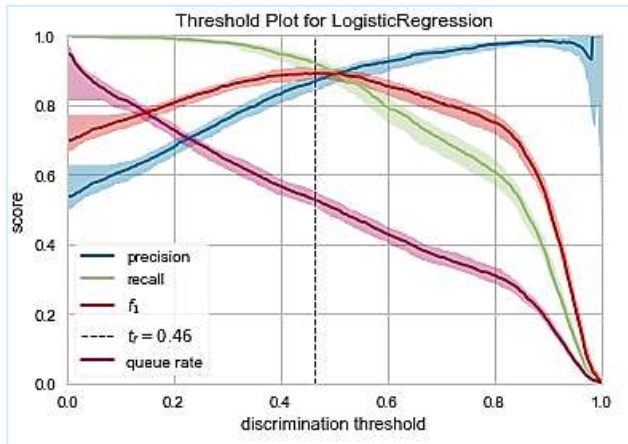


Fig. 15. Threshold Plot for Logistic Regression.

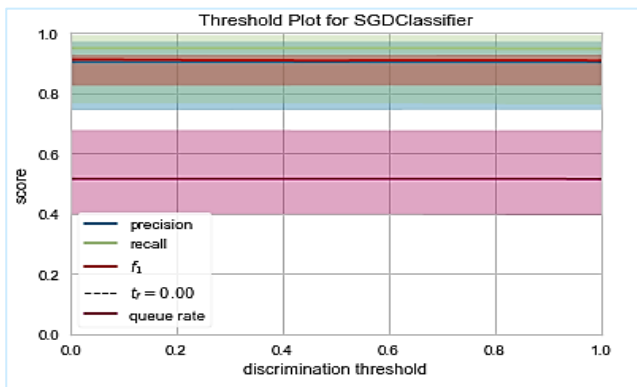


Fig. 16. Threshold Plot for SGD Classifier.

Fig. 16 shows the threshold plot for stochastic gradient descent (SGD). On the X-axis, we have the discrimination threshold, and on the Y-axis, we have the score. Here we see that the discrimination threshold for this is 0.00. For this threshold, we see that the precision, recall, and f1 score are approximately 0.8 to 0.9.

Fig. 17 shows the threshold plot for logistic regressionCV. On the X-axis, we have the discrimination threshold, and on the Y-axis, we have the score. Here we see that the discrimination threshold for this is 0.58. For this threshold, we see that the precision, recall, and f1 score are approximately around 0.95.

Fig. 18 shows the threshold plot for Bagging Classifier. On the X-axis, we have the discrimination threshold, and on the Y-axis, we have the score. Here we see that the discrimination threshold for this is 0.56. For this threshold, we see that the precision, recall, and f1 score are approximately around 0.98.

Fig. 19 shows the threshold plot for random forest classifier. On the X-axis, we have the discrimination threshold, and on the Y-axis, we have the score. Here we see that the discrimination threshold for this is 0.48. For this threshold, we see that the precision, recall, and f1 score are approximately around 0.99.

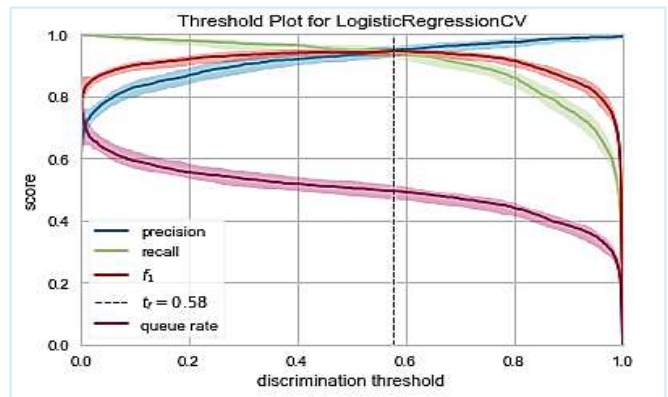


Fig. 17. Threshold Plot for Logistic Regression CV.

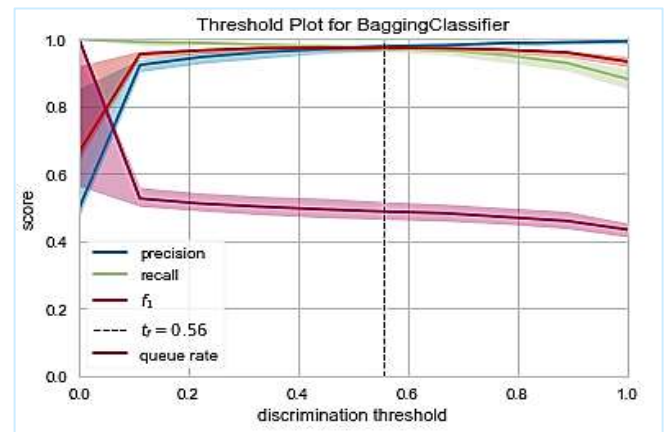


Fig. 18. Threshold Plot for Bagging Classifier.

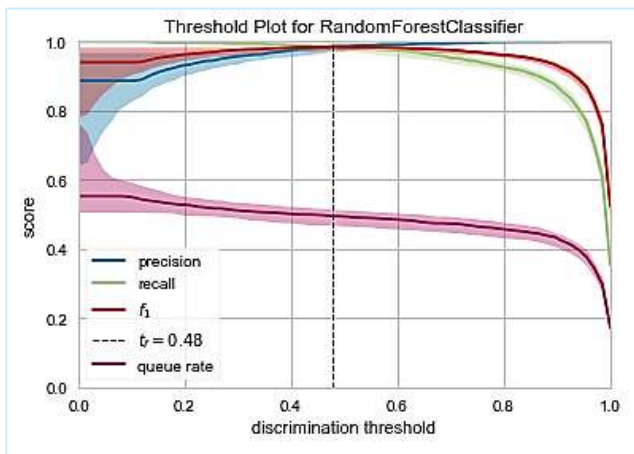


Fig. 19. Threshold Plot for Random Forest Classifier.

## V. CONCLUSION

Our work analyses different machine learning techniques when implemented over a dataset of features regarding websites and their corresponding details that may prove useful to detect a possible phishing website. This document aims to be useful to its readers to provide a conclusive analysis of these methods and to verify our observations regarding the random forest classifier's optimal performance. F1 score for the random forest is 0.99, which indicate that both false positive and false negative rate are in the satisfactory level. The graphs and details we have added to the document aim to help others carry out further experimentation to conclude our work. And we, ourselves, also intend to carry on our work with further modifications to the dataset and applying other machine learning techniques with modified parameters to hopefully open more possibilities in the hopes of improving the world's defenses against the cyber attackers out there. The internet is both fantastic and dangerous. And our work's main objective is to help minimize the danger by addressing a pervasive security issue of the modern world. In this paper, we apply basic machine learning algorithms. In the future, we will deploy deep learning techniques like multilayer perception and artificial neural networks to improve the performance of the detection system.

## REFERENCES

- [1] H. Alqahtani et al., "Cyber Intrusion Detection Using Machine Learning Classification Techniques," in *Computing Science, Communication and Security*, Singapore, 2020, pp. 121-131: Springer Singapore.
- [2] T. Bukth and S. S. Huda, *The soft threat: The story of the Bangladesh bank reserve heist*. SAGE Publications: SAGE Business Cases Originals, 2017.
- [3] P. Black, I. Gondal, R. J. C. Layton, and Security, "A survey of similarities in banking malware behaviours," vol. 77, pp. 756-772, 2018.
- [4] S. Hossain, A. Abtahee, I. Kashem, M. M. Hoque, and I. H. Sarker, "Crime Prediction Using Spatio-Temporal Data," in *Computing Science, Communication and Security*, Singapore, 2020, pp. 277-289: Springer Singapore.
- [5] S. Hossain, et al., "A Belief Rule Based Expert System to Predict Student Performance under Uncertainty," in *2019 22nd International Conference on Computer and Information Technology (ICIT)*, 2019, pp. 1-6.
- [6] S. Hossain, D. Sarma, R. J. Chakma, W. Alam, M. M. Hoque, and I. H. Sarker, "A Rule-Based Expert System to Assess Coronary Artery

- Disease Under Uncertainty," in *Computing Science, Communication and Security*, Singapore, 2020, pp. 143-159: Springer Singapore.
- [7] V. Shreeram, M. Suban, P. Shanthi and K. Manjula, "Anti-phishing detection of phishing attacks using genetic algorithm," *2010 International Conference on Communication Control and Computing Technologies*, Ramanathapuram, 2010, pp. 447-450, doi: 10.1109/ICCCCT.2010.5670593.
- [8] H. Huang, J. Tan and L. Liu, "Countermeasure Techniques for Deceptive Phishing Attack," *2009 International Conference on New Trends in Information and Service Science*, Beijing, 2009, pp. 636-641, doi: 10.1109/NISS.2009.80.
- [9] M. N. Feroz and S. Mengel, "Phishing URL Detection Using URL Ranking," *2015 IEEE International Congress on Big Data*, New York, NY, 2015, pp. 635-638, doi: 10.1109/BigDataCongress.2015.97.
- [10] S. Abu-Nimeh and S. Nair, "Bypassing Security Toolbars and Phishing Filters via DNS Poisoning," *IEEE GLOBECOM 2008 - 2008 IEEE Global Telecommunications Conference*, New Orleans, LO, 2008, pp. 1-6, doi: 10.1109/GLOCOM.2008.ECP.386.
- [11] Erkkila, J. "Why we fall for phishing." *Proceedings of the SIGCHI conference on Human Factors in Computing Systems CHI 2011*. ACM, 2011.
- [12] Khan, Ahmad Alamgir. "Preventing phishing attacks using one time password and user machine identification." *arXiv preprint arXiv:1305.2704* (2013).
- [13] A. Oest, Y. Safaei, A. Doupe, G. Ahn, B. Wardman and K. Tyers, "PhishFarm: A Scalable Framework for Measuring the Effectiveness of Evasion Techniques against Browser Phishing Blacklists," *2019 IEEE Symposium on Security and Privacy (SP)*, San Francisco, CA, USA, 2019, pp. 1344-1361, doi: 10.1109/SP.2019.00049.
- [14] M. Sharifi and S. H. Siadati, "A phishing sites blacklist generator," *2008 IEEE/ACS International Conference on Computer Systems and Applications*, Doha, 2008, pp. 840-843, doi: 10.1109/AICCSA.2008.4493625.
- [15] A. Belabed, E. Aïmeur and A. Chikh, "A Personalized Whitelist Approach for Phishing Webpage Detection," *2012 Seventh International Conference on Availability, Reliability and Security*, Prague, 2012, pp. 249-254, doi: 10.1109/ARES.2012.54.
- [16] Li, Linfeng, Marko Helenius, and Eleni Berki. "A usability test of whitelist and blacklist-based anti-phishing application." *Proceeding of the 16th International Academic MindTrek Conference*. 2012.
- [17] Usuff, Rahamathunnisa & Manikandan, N. & Kumaran, US & Niveditha, C.. (2017). Preventing from phishing attack by implementing url pattern matching technique in web. *International Journal of Civil Engineering and Technology*. 8. 1200-1208.
- [18] Hason N., Dvir A., Hajaj C. (2020) Robust Malicious Domain Detection. In: Dolev S., Kolesnikov V., Lodha S., Weiss G. (eds) *Cyber Security Cryptography and Machine Learning*. CSCML 2020. Lecture Notes in Computer Science, vol 12161. Springer, Cham. [https://doi.org/10.1007/978-3-030-49785-9\\_4](https://doi.org/10.1007/978-3-030-49785-9_4).
- [19] Hossein Shirazi, Bruhadeshwar Bezawada, and Indrakshi Ray. 2018. "Kn0w Thy DomaIn Name": Unbiased Phishing Detection Using Domain Name Based Features. In *Proceedings of the 23rd ACM on Symposium on Access Control Models and Technologies (SACMAT '18)*. Association for Computing Machinery, New York, NY, USA, 69-75. DOI:<https://doi.org/10.1145/3205977.3205992>.
- [20] Lasota K., Kozakiewicz A. (2011) Analysis of the Similarities in Malicious DNS Domain Names. In: Lee C., Seigneur JM., Park J.J., Wagner R.R. (eds) *Secure and Trust Computing, Data Management, and Applications*. STA 2011. Communications in Computer and Information Science, vol 187. Springer, Berlin, Heidelberg. [https://doi.org/10.1007/978-3-642-22365-5\\_1](https://doi.org/10.1007/978-3-642-22365-5_1).
- [21] Akinyelu, Andronicus A., and Aderemi O. Adewumi. "Classification of phishing email using random forest machine learning technique." *Journal of Applied Mathematics* 2014 (2014).
- [22] Gangavarapu, T., Jaidhar, C.D. & Chanduka, B. Applicability of machine learning in spam and phishing email filtering: review and approaches. *Artif Intell Rev* (2020). <https://doi.org/10.1007/s10462-020-09814-9>.

- [23] A. Belabed, E. Aïmeur and A. Chikh, "A Personalized Whitelist Approach for Phishing Webpage Detection," 2012 Seventh International Conference on Availability, Reliability and Security, Prague, 2012, pp. 249-254, doi: 10.1109/ARES.2012.54.
- [24] Alsaleh, M., Alarifi, A., Al-Quayed, F., & Al-Salman, A. (2015). Combating Comment Spam with Machine Learning Approaches. 2015 IEEE 14th International Conference on Machine Learning and Applications (ICMLA). doi:10.1109/icmla.2015.192.
- [25] Cuzzocrea, A., Martinelli, F., & Mercaldo, F. (2018). Applying Machine Learning Techniques to Detect and Analyze Web Phishing Attacks. Proceedings of the 20th International Conference on Information Integration and Web-Based Applications & Services - iiWAS2018. doi:10.1145/3282373.3282422.
- [26] Jeeva, S. Carolin, and Elijah Blessing Rajsingh. "Phishing URL detection-based feature selection to classifiers." International Journal of Electronic Security and Digital Forensics 9.2 (2017): 116-131.
- [27] Kulkarni, Arun D., and Leonard L. Brown III. "Phishing websites detection using machine learning." (2019).
- [28] Mao, Jian, et al. "Phishing page detection via learning classifiers from page layout feature." EURASIP Journal on Wireless Communications and Networking 2019.1 (2019): 43.
- [29] N. Sanglerdsinlapachai and A. Rungsawang, "Using Domain Top-page Similarity Feature in Machine Learning-Based Web Phishing Detection," 2010 Third International Conference on Knowledge Discovery and Data Mining, Phuket, 2010, pp. 187-190, doi: 10.1109/WKDD.2010.108.
- [30] Sahingoz, Ozgur Koray, et al. "Machine learning based phishing detection from URLs." Expert Systems with Applications 117 (2019): 345-357.
- [31] Ke Tian, Steve T. K. Jan, Hang Hu, Danfeng Yao, and Gang Wang. 2018. Needle in a Haystack: Tracking Down Elite Phishing Domains in the Wild. In Proceedings of the Internet Measurement Conference 2018 (IMC '18). Association for Computing Machinery, New York, NY, USA, 429–442. DOI: <https://doi.org/10.1145/3278532.3278569>.
- [32] T. Yue, J. Sun and H. Chen, "Fine-Grained Mining and Classification of Malicious Web Pages," 2013 Fourth International Conference on Digital Manufacturing & Automation, Qingdao, 2013, pp. 616-619, doi: 10.1109/ICDMA.2013.145.

# Identifying Critical Success Factors of Financial ERP System in Higher Education Institution using ADVIAN® Method

Ayogebob Epizitone<sup>1</sup>, Oludayo. O. Olugbara<sup>2</sup>

ICT and Society Research Group, South Africa Luban Workshop  
Durban University of Technology  
Durban, South Africa

**Abstract**—Enterprise resource planning (ERP) has been widely accepted by many organizations as an information technology process to seamlessly integrate, manage, and boost performance in different units of an organization. However, there linger an unpleasant chasm on success and satisfaction rates of ERP system implementation that have limited the effective use of the system. Moreover, the critical success factors (CSFs) of ERP system implementation have not been investigated in the literature for the case of financial functions in higher education institutions. This paper, through the application of advanced impact analysis (ADVIAN®) method exploits the CSFs of ERP system to support financial functions in a higher education institution. The applied ADVIAN® method highlights the CSFs that are measured according to the measures of criticality, integration, and stability. Furthermore, using precarious, driving, and driven measurements for ranking the factors, an effective model of CSFs for a financial ERP system implementation is attained to support financial functions. The study findings provide a comprehensive methodological scheme that can be used as a reference guide and as an orientation point for efficacious planning, implementing, and using ERP systems to support financial functions in higher education institutions.

**Keywords**—Cross impact; enterprise resource; impact analysis; resource planning; success factor

## I. INTRODUCTION

Enterprise resource planning (ERP) system implementation has been contextualized in past research studies as a medium for modernization and reformation of many organizations [1, 2]. The impetus for the adoption of ERP systems has been evident in literature for reasons such as pressure of competitiveness and improvement in operational efficiency [1, 3]. Many definitions exist in literature which describe ERP system as a unified software interface application with one large database, software package to facilitate seamless integration, configurable information system packages and computer-based systems [4, 5]. Al-Hadi and Al-Shaibany [4] consolidated diverse definitions of an ERP system as a software, a concept, a system or a package that integrates multiple modules as a separate functional area that includes a set of business processes with data flowing into a central database that could be uploaded locally or into the cloud. In the context of higher education institutions (HEIs), an ERP system has been defined as an information technology application that integrates automated recruitment, student

admission, financial aid, student records, and many academic and administrative services [2, 6].

Notwithstanding of the diverse definitions, an ERP system can be implemented in any organization regardless of size or nature [4]. Literature highlights the notable increase in an ERP system implementation and utilization from its onset till the present [3, 6]. However, this endeavor has been considered problematic, especially in HEIs that implement ERP systems to enhance quality of academic and administrative services [3, 6, 7]. Previous authors have attributed the causes of these snags to the lack of contextualization and deficiency of specific knowledge required [1, 3, 8]. Furthermore, authors have highlighted the existence of a solemn literature chasm that submerges the implementation of ERP systems in developing countries [9], and insufficient research conducted on successful implementations of ERP systems in HEIs [8]. Hence, with these chasms in mind, this paper seeks to make a significant contribution to the ERP phenomenon using a rigorous scientific method that will benefit practitioners and researchers.

The objective of this study was to investigate the critical success factors (CSFs) of ERP system implementation for the case of financial functions in higher education institutions using advanced impact analysis (ADVIAN®) method. Financial function is the business process of planning, acquiring, controlling, managing, and utilizing funds for the effective operations of an organization. The content of this paper is concisely structured as follows. Section 1 presents the introductory message. Section 2 reviews the related literature while Section 3 describes the ADVIAN® method. Section 4 presents the results of the analysis. Section 5 discusses the results and the paper is succinctly concluded in Section 6.

## II. LITERATURE REVIEW

Large scale fragmentation in HEIs has caused many of the institutions to seek for a unified ERP platform that can seamlessly integrate disparate information systems [1]. However, the implementation of ERP systems has been flouted with many challenges as previously mentioned. In the case of financial systems, there are reports of ERP systems lagging in real-time and not being supportive of the major financial functions in HEIs [10, 11]. The modular design of an ERP system caters for different organizational processes, but



each module that is not financially inclined correlates regularly with the financial module. In a complex interconnected system, a non-financial module provides and concomitantly needs useful information from the financial module for aspects such as reporting, budgeting, salaries, and fees. This concern was echoed by Noaman and Ahmed [11], that highlighted the deficiency of vendor proprietary ERP systems which need to response to the real functionalities of HEIs.

There have been red flags of the ERP systems not meeting the functional needs of HEIs despite that the space is progressing rapidly [11, 12]. Despite being one of the most frequently used systems [6], there is an alarming high rate of failure stemming from the implementation of ERP systems in HEIs [5]. Different authors have posited the implementation of ERP systems in HEIs to be complicated and presenting risks with no factors that can guarantee a successful implementation [4-6]. The plethora of intrinsic challenges of ERP systems posit the inadequacy of quality research that consider the uniqueness of the functional needs of HEIs [3, 5].

Previous studies have attempted to fill these chasms by proposing the identification of CSFs for implementation of ERP systems that targeted principal areas of business operations to ensure success. However, there has been censure for the inadequacy of sound scientific methods that can provide robust empirical evidence to support research findings [5, 8]. The concept of CSFs has been well research in many organizational environments, but little contributions have been experienced in the context of HEIs. Higher education system is attracting the keen interest of software vendors who view the system as a lucrative industry that is worth several hundred of billion dollars in revenue [11]. Several studies have delved into this dimension to derive CSFs for implementation of ERP systems in HEIs [5].

There have been a plethora of extensive studies on ERP implementation over the past decades that have contributed significantly to the understanding of the concept of CSFs [13]. However, within the extant literature, the identification of CSFs still require further investigation with authors calling for the application of more rigorous scientific methods [5, 8]. As a result, many of these authors have responded to the call by adopting CSFs to highlight areas that are critical for the successful implementation of ERP systems [13]. It was alluded by Loonam, Kumar, Mitra and Abd Razak [13] that CSFs are a powerful enabler that are extensively recognized within the ERP literature for identifying germane issues that require organizational attention before and during an ERP system implementation.

Shatat [14] has emphasized the importance of identifying CSFs of implementing ERP systems to facilitate the continuous success of the system and guarantee an improved influence on business performance. Sowan, Tahboub and Khamayseh [6] concentrated on technical success factors of an ERP system implementation. They identified 10 CSFs indicating the importance of factors that could support the structure that needs to be followed during implementation. Soliman and Karia [15] discussed a relatively small number of CSFs for implementing ERP systems in HEIs. They

highlighted the importance of CSFs while providing better understanding of whether their role is limited to influencing results at relevant stages in an innovation process. In addition, authors have highlighted the need for CSFs to maximize the potential outcome of an ERP system implementation with literature evidence to serve as footing to derive the factors using a sound scientific method that satisfies the need for analytically derived factors [16].

### III. METHOD OF ADVIAN®

This study follows a series of stages to identify CSFs that were subjected to expert evaluation which has led to the determination of cross-impact matrix and subsequent application of ADVIAN®, which is discussed briefly in this section. The ADVIAN® [17, 18] has been used in this study to determine the impacts that various CSFs have on a financial ERP system. This will allow cross-impact analysis that offers a provision for organizations to explore the current challenges and adequately prepare the right decisions for future endeavors in a manner that is participatory [5, 8, 17]. Impact analysis has been extensively employed in previous times in scenario analysis and future undertakings [17, 19]. It is currently employed in performance measurement to map tangible and intangible relationships [5, 17, 19].

The cross-impact analysis method utilizes an impact matrix that is filled based on the impact strength of the factors concerned. There are various measures of impact strength that have been used such as a rating score from 0 to 3, where 0 signifies no impact, 1 is a weak impact, 2 is a medium impact and 3 indicates a strong impact [5, 17]. This study rather uses a normalized cross-impact matrix with the impact strength normalized in the interval of 0 to 1 [5]. In this normalized range, 0.00 to 0.25 signifies no impact, 0.26 to 0.50 is a weak impact, 0.51 to 0.75 is a medium impact and 0.76 to 1.00 signifies a strong impact. Early cross-impact analysis methods did not have the ability of analyzing indirect interrelationships that gave rise to the alternate. Matrice d'Impacts Croisés Multiplication Appliquée à un Classement (MICMAC) addressed the inadequacy of direct impact [19]. However, this method along with other cross-impact analysis methods such as Papier computer and Fuzzy approach present intrinsic deficiencies because they consider either direct or indirect interrelations and fail to deal with both aspects concomitantly [5, 17-20].

The ADVIAN® is an improved quantitative impact analysis method that considers the indirect interactions amongst factors in a more reliable manner. The method uses active sum, passive sum and the impact strength of the corresponding factor for determining indirect interdependencies. Base on the method, CSFs can be identified reliably. However, it is limited because it does not measure the conditions of impact factors, but their interactions [18]. Hence, ADVIAN® does not give the status of a system, but rather identifies the important factors necessary for the entire system performance and regulation that justifies its application in this study. Additionally, ADVIAN® is favored over other methods for reasons such as the ability to provide deeper insights where other methods have proven inadequate and to provide an understanding of mutual relationship

between two single resources [5, 17, 18, 20]. Furthermore, it is suitable for explorative modeling as it overcomes the impossibility of privation of theory-based computational models that are because of inadequate theoretical advancement, establishing interrelationships and mutual connections based on expert judgments [19, 21]. The ADVIAN® has the capability of analyzing all interrelationships in diverse cross-impact matrices along with being able to perform supplementary measures such as criticality, stability, integration, driven, driving and precarious [5, 18, 19].

#### IV. RESULTS OF ANALYSIS

The cross-impact matrix data for this study were collected from nine financial system experts using an online survey instrument to elicit their judgements of 20 CSFs described in Table I. These 20 factors were obtained from 205 factors reported in 127 related research papers on ERP system implementation and categorized into four categories of resource (data, valuable, infrastructure, consultant and support in terms of management support, vendor support and training support), culture (communication, commitment, change, participation and values), project (implementation, team, leader and goals) and process (customization, package, plan and evaluation). Frequency distribution based on the number of previous research papers that cited a factor as being critical was constructed and normalized to realize a probability distribution. A mean probability threshold value of 0.88 was calculated to select 20 factors above the threshold value (Epizitone and Olugbara, 2019). Thereafter, these factors were presented to the experts to generate impact scores of factors. An understanding was reached among all participating experts by providing them with the necessary information and framework of the study. Given the nature of the study coupled with timeframe constraint, a minimum of 4 or 5 participants was accepted in the previous studies for impact score evaluation [8, 22, 23]. In this study, nine experts were engaged to provide their opinions on interrelationships of CSFs for implementing an ERP system in relation to financial functions. The number of participating experts is deemed acceptable for the research work of this nature.

During the process of data collection, an individual expert can use the online survey instrument administered to them to create several lists of preference chains of factors as suggested by Thompson, Olugbara and Singh [5]. The preference list approach enables the experts to list factors that are judged to be associated by transitive preference relation. The online survey instrument will then automatically process the numerous preference lists of factors created by experts to hatch a normalized cross-impact matrix (CIM) as shown in Table II.

##### A. Factor Relationships

The application of ADVIAN® method has yielded results that have provided useful insights into the interrelationships of CSFs. The normalized CIM has made it possible to assess interrelationships that revealed a maximum direct active sum of 14.11 and maximum direct passive sum of 7.78 to dictate how the investigated factors impact the implementation of

financial ERP system at varying levels. Table III presents further values of the average by two third of the standard deviation to comprehensively assess the characteristics of CSFs for a more detailed analysis. Emphasis was placed on CSFs whose values exceed the average by two third of the standard deviation to highlight significance of factors for implementing ERP systems [5, 19]. To supplement the assessment of CSFs, an additional analysis has been performed to consider factors that fell below the average not as much as two-third of the standard deviation. The interrelationship between two factors is usually expressed by active sum and passive sum that provide insightful knowledge about impacts of CSFs.

Direct relationship as demonstrated by active sum reveals the degree to which a CSF impacts on the implementation of financial ERP systems. Whereas passive sum demonstrates indirect relationship as the degree to which a financial ERP system impacts on CSFs. Relative scores were obtained by converting active sum and passive sum to percentile scale of 100 from 0 to permit the use of any number of factors in a system. The computational results based on the metrics of ADVIAN® are shown in Table III, where ellipses of DAS, RDAS, DPS, RDPS, RIAS, RIPS, stand for direct active sum, relative direct active sum, direct passive sum, relative direct passive sum, relative indirect active sum, and relative indirect passive sum respectively. Direct relationships produce high relative scores of 100 for active sum in the case of factor F1. It has a maximum impact on the financial system based on its impact on other factors in the system. Factor F17 gave a score of 48.82 relative direct active sum while factors F3 and F4 followed with the same value of 46.45. For the passive sum that divulges the factors that are directly influenced and affected by other factors in the system, F14 has the highest relative direct passive sum of 100 to indicate that it is the most affected and influenced by other factors. F12 has a relative direct passive sum of 88.57 while three factors F6, F11 and F13 have relative direct passive sum of 84.28, 82.85 and 81.43 respectively.

The average relative scores for the first order are 36.06 and 65.42 for relative active sum and relative passive sum. The factors that exceeded the average by two-third of the standard deviation for either passive sum or active sum have an adverse effect on the other factors. F1 with maximum relative active sum has the lowest relative passive sum. The relative active sums for the factors F11, F12, F14 and F18 fall below the average and less the two third of the standard deviation (23.33). However, they have high values that exceeded the 79.31 average by two-third of the standard deviation for relative passive sum to reveal the degree of relationships that exist in the ERP system. This result indicates that a factor with high degree of direct impact on system has less influence on other factors. F1 has the strongest impact exerted on other factors but is less influenced by other factors with a value of 100 to 5.86. The impact of F1 can be clearly identified as shown in Table III to be less significant but has the strongest impact on the system. The factors F4, F8, F9 and F13 are highly influenced by F1. The most significant impact on factor F14 emanated from factors F1, F3, F4, F5, F7, F10 and F19.

TABLE I. CATEGORY AND DESCRIPTION OF FACTORS

| Factor | Category | Description  |
|--------|----------|--|
| F1     | Resource | Top management support and commitment  |
| F2     | Culture  | Interdepartmental communication and cooperation throughout the institution   |
| F3     | Culture  | Commitment to business process reengineering to do away with redundant processes   |
| F4     | Project  | Implementation of project management from initiation to closing  |
| F5     | Culture  | Change management program to ensure awareness for any changes that may happen  |
| F6     | Project  | Project team competence, formulation, composition, and involvement   |
| F7     | Resource | Education and training for stakeholders, including end users, technical and IT staff   |
| F8     | Project  | Project champion presence to lead the implementation, authorized to use internal and external resources to complete implementation |
| F9     | Project  | Project mission and goals for the system with clear objective agreed upon  |
| F10    | Resource | ERP expert consultant use to guide the implementation process  |
| F11    | Process  | Minimum level of customization to utilize ERP functionalities to a maximum   |
| F12    | Process  | Package selection, carefully and professionally selected   |
| F13    | Culture  | Understanding the institutional culture, norms, values, and beliefs  |
| F14    | Culture  | User involvement and participation throughout implementation   |
| F15    | Resource | ERP vendor support and partnership   |
| F16    | Process  | Business vision and plan   |
| F17    | Resource | Adequate IT infrastructure   |
| F18    | Process  | Monitoring management especially evaluation of performance metrics   |
| F19    | Resource | Allocating and dedicating valuable resources   |
| F20    | Resource | Data management plan that ensures that data are accurately and efficient migrated to a new system and analysed properly            |

TABLE II. NORMALIZED CROSS-IMPACT MATRIX OF 20 FACTORS

| Factor | F1  | F2  | F3  | F4  | F5  | F6  | F7  | F8  | F9  | F10 | F11 | F12 | F13 | F14 | F15 | F16 | F17 | F18 | F19 | F20 |
|--------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| F1     | 0.0 | 1.0 | 1.0 | 0.7 | 0.8 | 0.9 | 0.9 | 0.7 | 0.7 | 0.9 | 0.4 | 0.6 | 0.9 | 0.7 | 0.8 | 0.7 | 1.0 | 0.3 | 0.7 | 0.7 |
| F2     | 0.0 | 0.0 | 0.4 | 0.2 | 0.2 | 0.4 | 0.3 | 0.4 | 0.2 | 0.4 | 0.3 | 0.3 | 0.4 | 0.3 | 0.3 | 0.2 | 0.2 | 0.2 | 0.1 | 0.2 |
| F3     | 0.0 | 0.3 | 0.0 | 0.2 | 0.4 | 0.3 | 0.6 | 0.3 | 0.3 | 0.3 | 0.3 | 0.6 | 0.4 | 0.6 | 0.3 | 0.3 | 0.3 | 0.2 | 0.2 | 0.3 |
| F4     | 0.1 | 0.4 | 0.3 | 0.0 | 0.4 | 0.6 | 0.3 | 0.4 | 0.1 | 0.3 | 0.4 | 0.3 | 0.2 | 0.6 | 0.4 | 0.2 | 0.2 | 0.3 | 0.2 | 0.4 |
| F5     | 0.0 | 0.3 | 0.4 | 0.1 | 0.0 | 0.3 | 0.3 | 0.2 | 0.2 | 0.4 | 0.3 | 0.4 | 0.1 | 0.7 | 0.2 | 0.3 | 0.2 | 0.2 | 0.1 | 0.4 |
| F6     | 0.0 | 0.1 | 0.2 | 0.0 | 0.2 | 0.0 | 0.0 | 0.6 | 0.2 | 0.2 | 0.3 | 0.2 | 0.2 | 0.3 | 0.1 | 0.1 | 0.0 | 0.2 | 0.0 | 0.2 |
| F7     | 0.0 | 0.4 | 0.3 | 0.1 | 0.3 | 0.6 | 0.0 | 0.2 | 0.2 | 0.4 | 0.6 | 0.4 | 0.2 | 0.6 | 0.3 | 0.3 | 0.3 | 0.3 | 0.2 | 0.3 |
| F8     | 0.1 | 0.1 | 0.2 | 0.1 | 0.2 | 0.1 | 0.1 | 0.0 | 0.4 | 0.3 | 0.2 | 0.3 | 0.3 | 0.3 | 0.2 | 0.1 | 0.1 | 0.2 | 0.2 | 0.2 |
| F9     | 0.1 | 0.2 | 0.2 | 0.3 | 0.2 | 0.3 | 0.2 | 0.2 | 0.0 | 0.4 | 0.2 | 0.4 | 0.3 | 0.3 | 0.2 | 0.1 | 0.3 | 0.3 | 0.3 | 0.2 |
| F10    | 0.0 | 0.2 | 0.3 | 0.1 | 0.1 | 0.2 | 0.1 | 0.1 | 0.1 | 0.0 | 0.3 | 0.6 | 0.4 | 0.6 | 0.3 | 0.3 | 0.3 | 0.2 | 0.3 | 0.2 |
| F11    | 0.0 | 0.1 | 0.1 | 0.0 | 0.0 | 0.1 | 0.1 | 0.1 | 0.1 | 0.0 | 0.0 | 0.3 | 0.1 | 0.2 | 0.0 | 0.0 | 0.0 | 0.2 | 0.0 | 0.1 |
| F12    | 0.0 | 0.1 | 0.1 | 0.0 | 0.0 | 0.1 | 0.1 | 0.1 | 0.1 | 0.0 | 0.2 | 0.0 | 0.1 | 0.1 | 0.1 | 0.1 | 0.2 | 0.2 | 0.2 | 0.2 |
| F13    | 0.1 | 0.2 | 0.2 | 0.1 | 0.3 | 0.3 | 0.3 | 0.2 | 0.2 | 0.3 | 0.2 | 0.3 | 0.0 | 0.3 | 0.2 | 0.1 | 0.1 | 0.3 | 0.4 | 0.3 |
| F14    | 0.0 | 0.2 | 0.0 | 0.0 | 0.0 | 0.2 | 0.0 | 0.2 | 0.1 | 0.0 | 0.2 | 0.2 | 0.3 | 0.0 | 0.1 | 0.1 | 0.0 | 0.2 | 0.0 | 0.3 |
| F15    | 0.0 | 0.3 | 0.3 | 0.0 | 0.3 | 0.4 | 0.1 | 0.3 | 0.2 | 0.3 | 0.4 | 0.4 | 0.6 | 0.4 | 0.0 | 0.4 | 0.1 | 0.3 | 0.4 | 0.3 |
| F16    | 0.0 | 0.3 | 0.3 | 0.1 | 0.2 | 0.3 | 0.2 | 0.3 | 0.3 | 0.2 | 0.4 | 0.4 | 0.6 | 0.3 | 0.3 | 0.0 | 0.2 | 0.3 | 0.4 | 0.3 |
| F17    | 0.0 | 0.4 | 0.3 | 0.1 | 0.3 | 0.6 | 0.4 | 0.3 | 0.2 | 0.3 | 0.4 | 0.3 | 0.6 | 0.3 | 0.6 | 0.4 | 0.0 | 0.4 | 0.4 | 0.2 |
| F18    | 0.0 | 0.1 | 0.1 | 0.0 | 0.1 | 0.1 | 0.0 | 0.1 | 0.0 | 0.1 | 0.1 | 0.1 | 0.0 | 0.1 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| F19    | 0.0 | 0.3 | 0.2 | 0.1 | 0.2 | 0.3 | 0.1 | 0.2 | 0.1 | 0.3 | 0.3 | 0.2 | 0.3 | 0.6 | 0.2 | 0.1 | 0.1 | 0.3 | 0.0 | 0.3 |
| F20    | 0.0 | 0.4 | 0.2 | 0.1 | 0.1 | 0.2 | 0.3 | 0.2 | 0.1 | 0.2 | 0.4 | 0.2 | 0.1 | 0.4 | 0.1 | 0.0 | 0.2 | 0.3 | 0.1 | 0.0 |

TABLE III. DIRECT AND INDIRECT RELATIONSHIPS

| Factor         | DAS   | RDAS   | Ranking by RDAS | DPS  | RDPS   | Ranking by RDPS | RIAS   | Ranking by RIAS | RIPS   | Ranking by RIPS |
|----------------|-------|--------|-----------------|------|--------|-----------------|--------|-----------------|--------|-----------------|
| F1             | 14.11 | 100.00 | 1               | 0.44 | 5.71   | 20              | 100.00 | 1               | 5.86   | 20              |
| F2             | 5.56  | 39.37  | 9               | 5.89 | 75.71  | 6               | 36.01  | 7               | 69.06  | 9               |
| F3             | 6.56  | 46.45  | 3               | 5.56 | 71.42  | 8               | 41.81  | 4               | 64.65  | 11              |
| F4             | 6.56  | 46.45  | 4               | 2.44 | 31.43  | 19              | 41.86  | 3               | 25.91  | 19              |
| F5             | 5.56  | 39.37  | 8               | 4.67 | 60.00  | 14              | 34.15  | 10              | 53.47  | 15              |
| F6             | 3.33  | 23.62  | 16              | 6.56 | 84.28  | 3               | 19.16  | 16              | 78.44  | 5               |
| F7             | 6.33  | 44.88  | 5               | 4.67 | 60.00  | 13              | 37.89  | 5               | 53.64  | 14              |
| F8             | 4.11  | 29.13  | 14              | 5.44 | 70.00  | 10              | 27.67  | 13              | 69.87  | 8               |
| F9             | 5.22  | 37.00  | 10              | 4.11 | 52.85  | 16              | 34.78  | 9               | 52.30  | 16              |
| F10            | 5.00  | 35.43  | 11              | 5.78 | 74.28  | 7               | 30.48  | 12              | 67.77  | 10              |
| F11            | 1.67  | 11.81  | 19              | 6.44 | 82.85  | 4               | 8.54   | 19              | 85.56  | 3               |
| F12            | 2.22  | 15.75  | 18              | 6.89 | 88.57  | 2               | 13.67  | 17              | 92.00  | 2               |
| F13            | 4.89  | 34.64  | 12              | 6.33 | 81.43  | 5               | 31.73  | 11              | 79.04  | 4               |
| F14            | 2.33  | 16.53  | 17              | 7.78 | 100.00 | 1               | 12.76  | 18              | 100.00 | 1               |
| F15            | 6.00  | 42.52  | 6               | 5.00 | 64.28  | 12              | 35.52  | 8               | 57.91  | 12              |
| F16            | 5.89  | 41.73  | 7               | 4.11 | 52.85  | 17              | 36.12  | 6               | 48.61  | 17              |
| F17            | 6.89  | 48.82  | 2               | 4.11 | 52.85  | 18              | 43.28  | 2               | 45.2   | 18              |
| F18            | 1.00  | 7.09   | 20              | 5.44 | 69.99  | 11              | 5.79   | 20              | 74.65  | 6               |
| F19            | 4.56  | 32.28  | 13              | 4.56 | 58.57  | 15              | 26.58  | 14              | 54.86  | 13              |
| F20            | 4.00  | 28.34  | 15              | 5.56 | 71.42  | 9               | 23.68  | 15              | 71.36  | 7               |
| Average        | 5.09  | 36.06  |                 | 5.09 | 65.42  |                 | 32.07  |                 | 62.51  |                 |
| STD DEV        | 2.74  | 19.39  |                 | 1.62 | 20.83  |                 | 19.55  |                 | 21.85  |                 |
| AVG+2/3std Dev | 6.91  | 48.99  |                 | 6.17 | 79.31  |                 | 45.11  |                 | 77.08  |                 |
| AVG-2/3std Dev | 3.26  | 23.13  |                 | 4.01 | 51.54  |                 | 19.04  |                 | 47.94  |                 |

### B. Classification of Factors

The classification of CSFs can be done using the measures of criticality, integration, and stability for conditional state of system of factors. The state of a system can be significantly altered because of changes in any factors deem critical in the system. Table IV shows the computed values for criticality, integration and stability of factors based on the calculations of ADVIAN®.

There are significant changes in the system when there are changes to any of the CSFs based on their criticality. A high level of criticality was obtained for factors F2, F3 and F13. These factors have a high level of criticality, when looking at the average by the standard deviation. F3 has the highest criticality score of 51.99, followed by F13 with 50.08 and 49.87 for F2. These factors necessitate an early update for corrective measures to be taken should they change. A low level of criticality which is below 34.24 was obtained for factors F1, F4, F11 and F18, hence changes to these factors render minimal impacts on the system.

Fig. 1 shows the contour lines of the criticality corresponding to the system stability. However, the active sum and passive sum present the reverse dependency. This implies

that factors with high criticality will have a low stability as in the case of F2, F3 and F

The connection of factors in the system can be determined by the level of integration. A high-level integration (55.39) for a factor indicates a strong connection with the rest of the system. A high value of integration was obtained for factor F13 while factor F14 gave a value of 56.38. The other factors with high integration are F1 (52.93), F2 (52.53), F3 (53.23) and F12 (52.83) with values above the 51.09 threshold. This presents the existence of a mutual connection and likely feedback loops among these factors. These feedback loops strongly reinforce each other mutually in indirect connections at different levels to confirm the presence of mutual connections and feedback loops among the factors as can be seen by the contour line of Fig. 2. In the integration system grid, high integration factors of F14, F13 and F1 have high passive sum, high active sum and low integration score with low passive sum and active sum.

The determination of system stability was based on the distribution of factors toward active sum and passive sum axes [5, 17-20]. Factors that are controlled by the system are aligned close to the axis of the passive sum with low active sum, while factors that control the systems are closed to the

active sum axis with low passive sum (Linss and Fried, 2010). Stability level attenuates feedback loops to ensure the absence of uncontrolled feedback loops [5, 17, 19]. A high stability of 73.08 of the average by two-third of the standard deviation was obtained. Factors F1, F11, F12, F14 and F18 contribute heavily to the system stability. F18 (89.26), F1 (88.92) and F11(84.47) are factors with high stability values, followed by F12 (76.20) and F14 (77.36). The combinations of these factors with high passive sum and low active sum coupled with their different integration levels and high stability indicate that they are independent factors within the system which can hardly alter these factors. Fig. 3 shows the contour line for the system stability which position factors F4 and F6 above the system stability of 64.97.

### C. Ranking of Factors

Precarious, driving and driven are essential measures for ranking of CSFs. The precarious measure is obtained for a factor by calculating the geometric mean of the active indirect sum and criticality measurement. The driving measure of a factor considers the geometric mean of indirect sum without a complete percentage of the criticality measure. The driven measure for each factor is an inverse of the driving measure which substitutes the indirect active sum with passive sum. Table V shows the scores computed for precarious, driving and driven based on the ADVIAN® method. Precarious CSFs exert utmost influence on the system and are affected by peripheral forces. The factors F1 (49.21), F2 (42.37), F3

(46.63) and F17 (43.75) present high precarious scores above the average by two third of the standard deviation of 41.40. These factors have strong influence on the system and are invulnerable to external forces because they are not ideal to warrant intervening activities.

Fig. 4 presents the contour line for the determination of the most precarious factors and lowest precarious factors of F18, F11 and F14. High driving ranking for CSFs is essential for implementation success because these factors present high influence on other factors in the system and do not cause strong feedback. Factors of F1 (87.05), F4 (52.98) and F17 (49.13) are drivers with outstanding characteristics to drive a successful financial ERP system implementation. These factors are non-critical with high active sum, but they are good beginning point for intervention activities because of their suitability for external actions that dependent on the ability to influence other factors without causing strong feedback.

Fig. 5 shows the contours for driving factors. Driven impact factors are non-critical with high passive sums. They are more reactive in nature and are not reasonably altered by intervening activities. Nonetheless, they indicate the success of external actions taken on driving factors and not reasonably affected by external changes made to the system. Factors of F6, F11, F12, F14 and F18 are good indicators of success of external intervention on financial ERP system implementation by driving factors. Fig. 6 presents the contour lines for the driven ranking impact factors.

TABLE IV. CLASSIFICATION OF FACTORS BY MEASURES OF CRITICALITY, INTEGRATION, AND STABILITY

| Factor         | Critical | Integration | Stability |
|----------------|----------|-------------|-----------|
| F1             | 24.22    | 52.93       | 88.92     |
| F2             | 49.87    | 52.53       | 52.67     |
| F3             | 51.99    | 53.23       | 49.22     |
| F4             | 32.93    | 33.88       | 68        |
| F5             | 42.73    | 43.81       | 58.32     |
| F6             | 38.77    | 48.8        | 69.2      |
| F7             | 45.08    | 45.76       | 55.59     |
| F8             | 43.97    | 48.77       | 60.36     |
| F9             | 42.65    | 43.54       | 58.23     |
| F10            | 45.45    | 49.12       | 57.95     |
| F11            | 27.03    | 47.05       | 84.47     |
| F12            | 35.46    | 52.83       | 76.2      |
| F13            | 50.08    | 55.39       | 54.71     |
| F14            | 35.73    | 56.38       | 77.36     |
| F15            | 45.35    | 46.71       | 55.97     |
| F16            | 41.9     | 42.37       | 58.55     |
| F17            | 44.23    | 44.24       | 55.78     |
| F18            | 20.79    | 40.22       | 89.25     |
| F19            | 38.19    | 40.72       | 64.19     |
| F20            | 41.11    | 47.52       | 64.44     |
| Average        | 39.88    | 47.29       | 64.97     |
| STD DEV        | 8.46     | 5.7         | 12.16     |
| AVG+2/3std Dev | 45.52    | 51.09       | 73.08     |
| AVG-2/3std Dev | 34.24    | 43.49       | 56.86     |

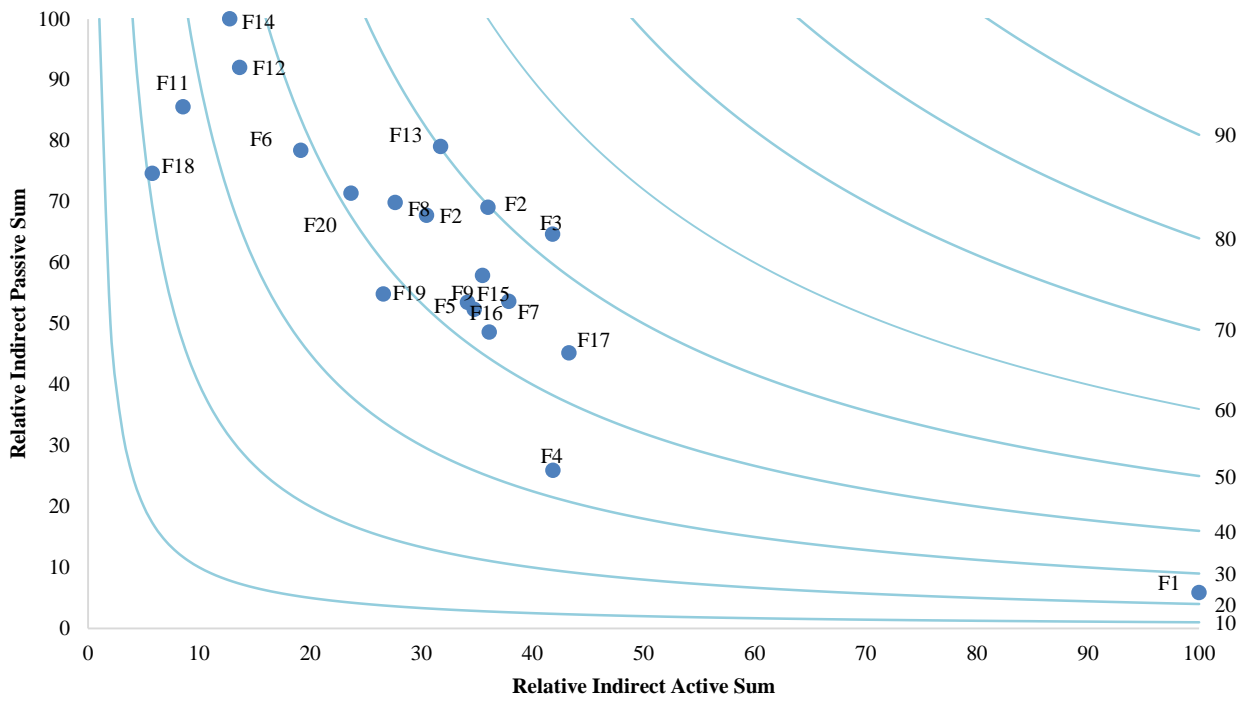


Fig. 1. Classification of Factors by a Measure of Criticality.

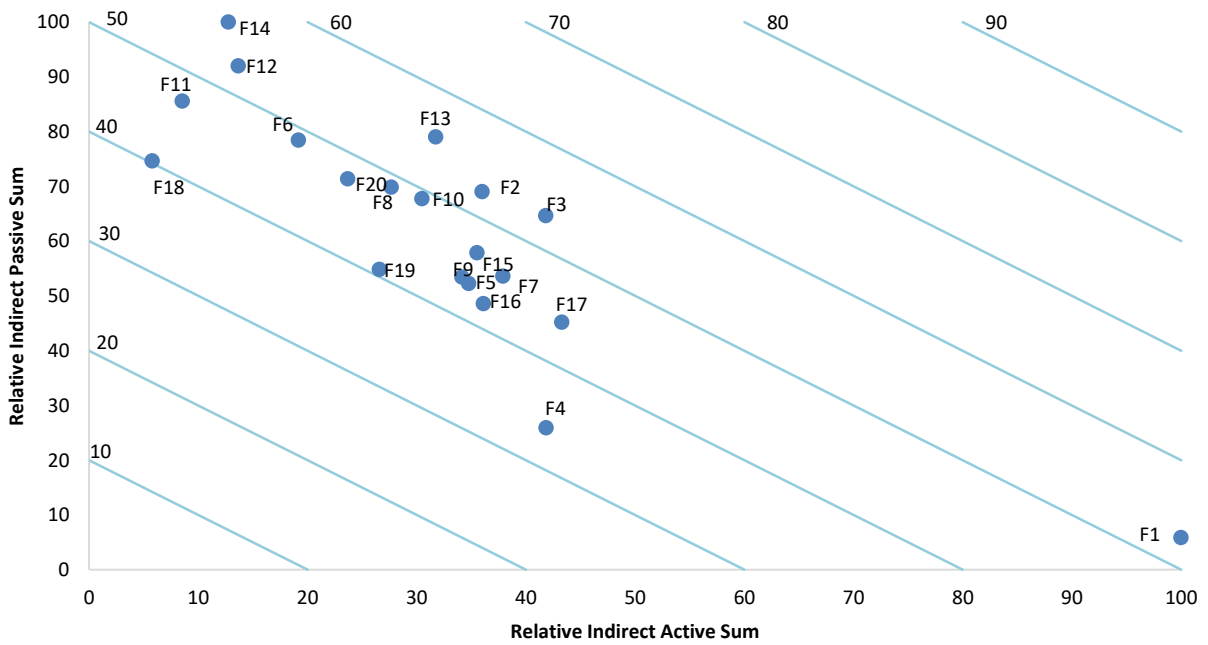


Fig. 2. Classification of Factors by a Measure of Integration.

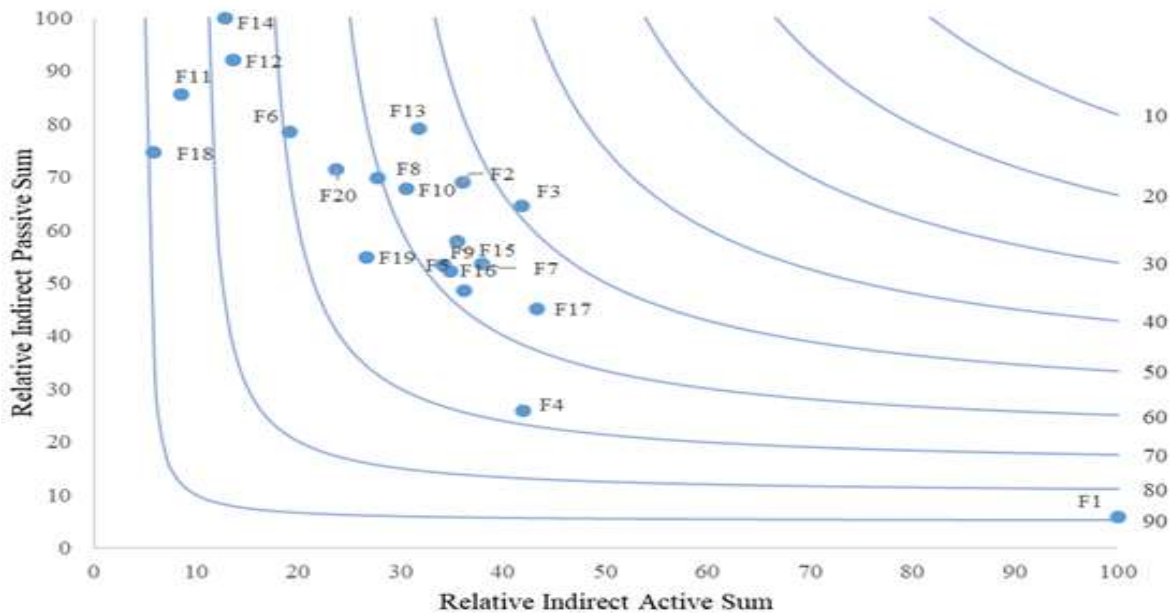


Fig. 3. Classification of Factors by a Measure of Stability.

TABLE V. RANKING OF FACTORS BY MEASURES OF PRECARIOUS, DRIVING, AND DRIVEN

| Factor         | Precarious | Precarious Ranking | Driving | Driving Ranking | Driven | Driven Ranking |
|----------------|------------|--------------------|---------|-----------------|--------|----------------|
| F1             | 49.21      | 1                  | 87.05   | 1               | 21.08  | 20             |
| F2             | 42.37      | 4                  | 42.49   | 10              | 58.84  | 10             |
| F3             | 46.63      | 2                  | 44.80   | 6               | 55.71  | 13             |
| F4             | 37.13      | 12                 | 52.98   | 2               | 41.68  | 19             |
| F5             | 38.20      | 10                 | 44.22   | 8               | 55.34  | 14             |
| F6             | 27.26      | 16                 | 34.25   | 16              | 69.30  | 5              |
| F7             | 41.33      | 5                  | 45.62   | 5               | 54.27  | 16             |
| F8             | 34.88      | 13                 | 39.37   | 14              | 62.57  | 8              |
| F9             | 38.51      | 9                  | 44.66   | 7               | 54.77  | 15             |
| F10            | 37.22      | 11                 | 40.78   | 11              | 60.80  | 9              |
| F11            | 15.20      | 19                 | 24.97   | 19              | 79.01  | 2              |
| F12            | 22.01      | 17                 | 29.70   | 17              | 77.06  | 3              |
| F13            | 39.87      | 7                  | 39.80   | 13              | 62.81  | 7              |
| F14            | 21.35      | 18                 | 28.64   | 18              | 80.17  | 1              |
| F15            | 40.14      | 6                  | 44.06   | 9               | 56.26  | 12             |
| F16            | 38.90      | 8                  | 45.81   | 4               | 53.14  | 17             |
| F17            | 43.75      | 3                  | 49.13   | 3               | 50.21  | 18             |
| F18            | 10.97      | 20                 | 21.42   | 20              | 76.89  | 4              |
| F19            | 31.86      | 14                 | 40.54   | 12              | 58.23  | 11             |
| F20            | 31.20      | 15                 | 37.34   | 15              | 64.83  | 6              |
| Average        | 34.40      |                    | 41.88   |                 | 59.65  |                |
| STD DEV        | 10.28      |                    | 13.36   |                 | 13.72  |                |
| AVG+2/3std Dev | 41.25      |                    | 50.79   |                 | 68.80  |                |
| AVG-2/3std Dev | 27.55      |                    | 32.97   |                 | 50.50  |                |

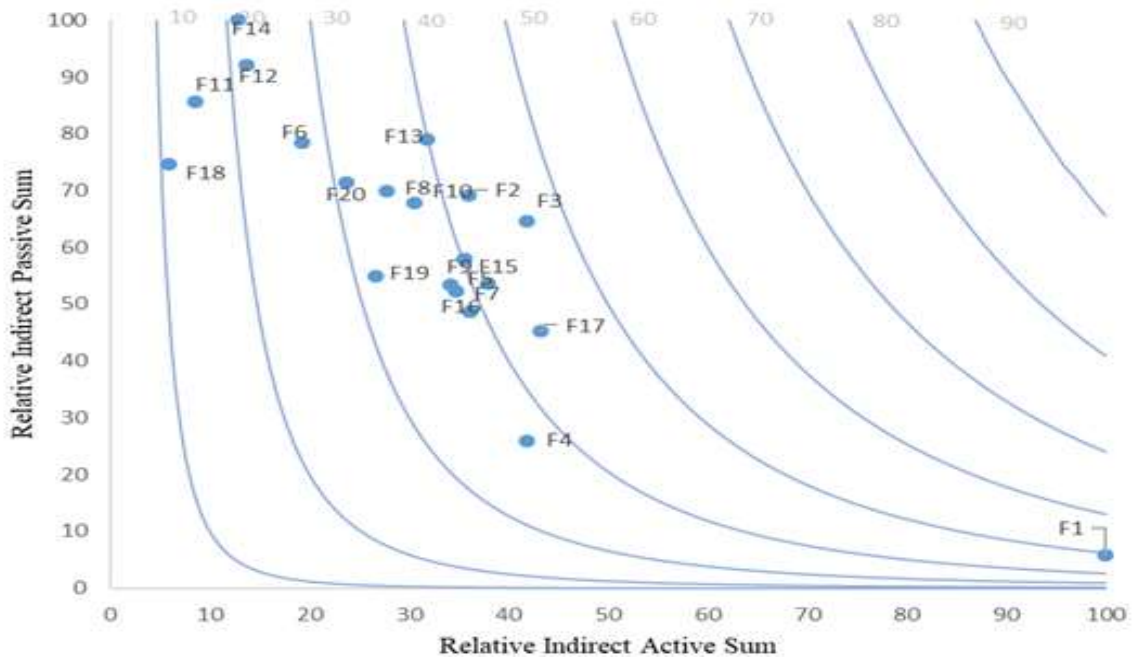


Fig. 4. Ranking of Factors by a Measure of Precarious.

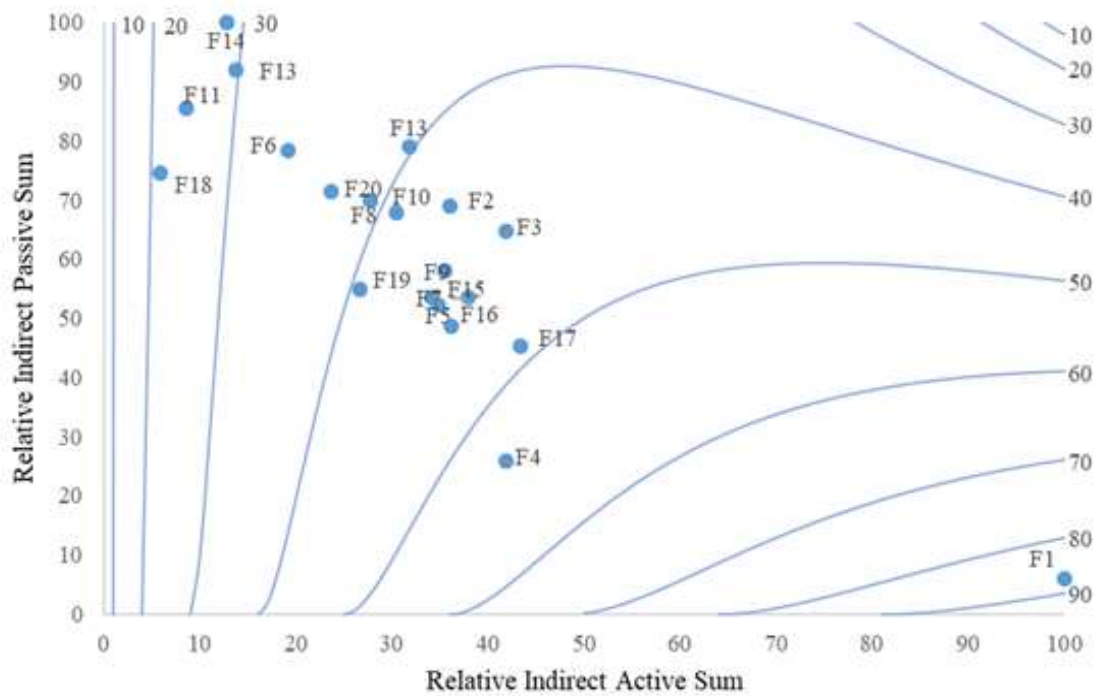


Fig. 5. Ranking of Factors by a Measure of Driving.



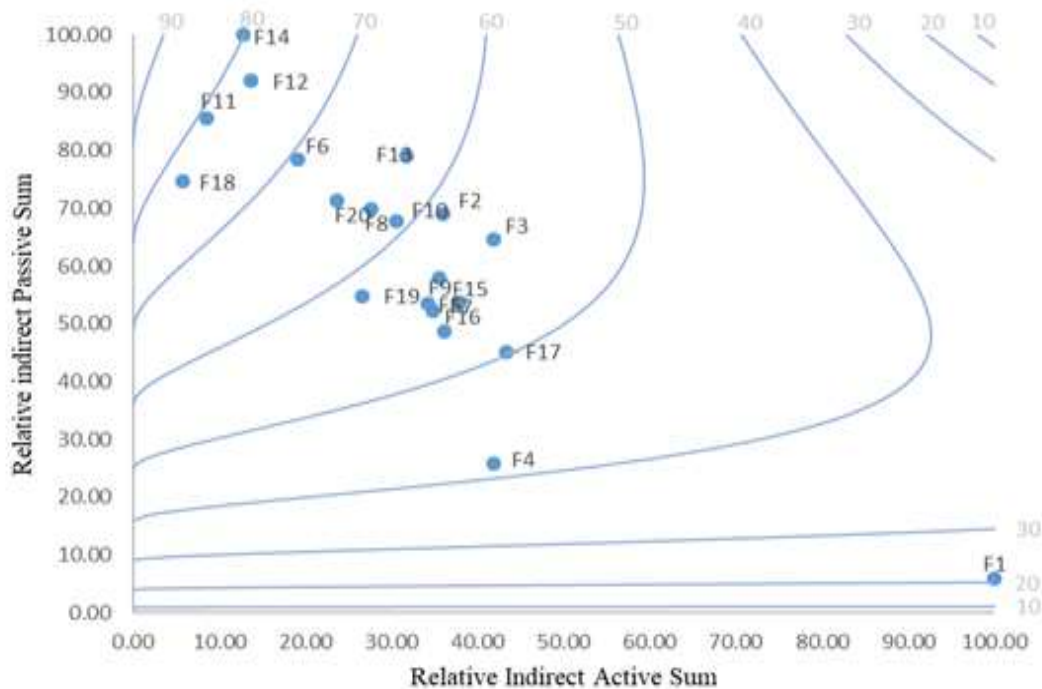


Fig. 6. Ranking of Factors by a Measure of Driven.

## V. DISCUSSION OF RESULTS

The objective of this study was to use the ADVIAN® method to determine factors that are critical for implementing ERP systems to support financial functions. Consequently, 20 top CSFs aggregated from the literature were subjected to ADVIAN®. The exploration of these factors has yielded valuable insights to the successful ERP system implementation that would support financial functions in a HEI. In this study, top management support and commitment (F1) is highly relevant for implementing financial ERP systems. Given that this factor has been cited by several authors in the general case of ERP system implementation is not a coincidence. The factor means authorizing, commissioning, and making available resources that are needed for the implementation of a financial ERP system. A successfully implemented financial ERP system is an asset to top management who rely heavily on outputs of the system to support strategic decision making in an organization. Despite this result, literature has reported the dysfunctionality of ERP system in financial reporting, which requires that attention be given to this factor [24]. The factor presents a high active sum of 100, stability of 88.92 and driving score of 87.05. Despite the low criticality of this factor, any changes presented to it would need intervening activities because it has great impact on other factors in the system. Hence, mirroring the previous and current studies for successful implementation of ERP systems [25-27].

Interdepartmental communication and cooperation (F2), is highly critical for consideration in the implementation of a financial ERP system. The need for open communication and cooperation among different stakeholders cannot be overlooked. Failure to properly communicate and cooperate is likely to cause system failure ([10, 28, 29]. This factor

presents a high passive sum of 75.71 and active sum of 39.37 in explaining the high criticality of 49.87 and precarious ranking of 42.37. This result makes it of great influence because of its dependency on other factors and is strongly influenced by other factors of the system. A high driving value of 42.49 is a clear indication that is it critical. Commitment to business process reengineering to do away with redundant processes (F3) has been found to be critical to financial ERP system implementation where functions have been highly automated. Business process reengineering refers to the restructuring of business processes in an organization. The factor presents a high passive sum of 71.42, which indicates the impact that other factors exert on it. In addition, the factor reflected a high criticality score of 51.99 and integration of 53.23. These values indicate the contributions of the factor to successful implementation of ERP system to support financial functions. This study converges with previous studies that cited this factor as critical, given its classification as an important integrator in the system with high precarious score of 46.63 is an indication of strong influence on the system. Moreover, a driven score of 44.80 above the average makes it critical as one prime organizational resource for improvement that warrants an intervention activity should there be any change to it [5, 30].

The implementation of project management from the initiation phase to the closing (F4) was identified in this study to be a factor with high stability score of 68.00, which implies its significant contribution to the stability of an ERP system implementation. The driving value of 52.98 was obtained for the factor that alludes to its criticality in a previous study and highlights its significance in the success of financial ERP system because of being active ([1]. Project management entails a broader scope of implementing ERP systems to support financial functions. For the financial system to be

successful, a good project management plan must be instituted. To promote successful financial system implementation, it is important to build team, deal with conflict and execute objectives of the implementation [31]. It is important to focus attention on the implementation of a financial system given that ERP project can tend to be complex [32, 33].

Change management program to build awareness and ensure readiness for changes that may happen (F5) was identified in this study to be critical. Presenting high value of 39.37 for active sum and high value of 44.22 for driving than the average. It presents itself to be active in the system implementation and impacting other factors in the system. Change spans a great length ranging from cultural, organizational, and structural changes and therefore, a stable and successful setting is required obviously for successful implementation of ERP system [34-37]. Project team competence of formulation, composition and involvement (F6) has emerged in this study to concretely align with previous studies as a factor of great influence on the system because of other factors such as F1, F4, F7 and F17 that significant impact on it. It has low precarious value of 27.26 indicating that it can be used for intervention given that it may be suitably influenced by external actions. Project team plays an important role in the system implementation indicating that formation, participation, skills and involvement is highly relevant for the success of implementing an ERP system [34-38]. Rightly deploying a project team with competency to execute specific tasks can mediate a successful financial ERP system, given that it has a high driven score of 69.30. Hence, it is a good indicator of successful intervention that requires attention of management.

Education and training for stakeholders including end users, technical and information technology staff (F7) is deemed to be critical in literature on ERP systems. Lack of education can be directly attributed to many implementation challenges such as system failure, employee uneasiness with using the system and training overhead that have been reported in literature [30, 34, 35, 39, 40]. The precarious value of 41.33 implies that this factor is a precarious impact factor with high activities that require actions to be taken during the implementation should there be any issue affecting it. Project champion to lead the financial system implementation, authorized to use internal and external resources to complete financial ERP system implementation (F8) is important. The significance of ERP project champion is recognized in terms of a precarious value of 34.88 that indicates its influence on the system. This study deemed this factor significant because of the driven nature with a value of 62.57, which makes it a good indicator of success and integration score of 48.77 requires that it be monitored.

Project mission and goals for the system with clear objective agreed upon (F9) is rank 9th place to be precarious (38.51) and 7th to be driving (44.66). This supports the listing of the factor as critical to financial ERP system implementation in accordance with previous literature. The factor will require corrective actions and interventions if there is any deviation from the goals and mission of the project that will adversely affect implementation. Consequently, it alludes

to the need for a clear future plan for implementing financial ERP systems as echoed in the previous studies [2, 4, 5]. ERP expert consultant use to guide the implementation process (F10) of financial ERP systems is critical with a value of 45.45 above the average and minimum threshold of 34.21. Driven value of 60.80 computed for this factor coupled with a high passive sum of 74.28 support the previous studies that have highlighted it as a success indicator. The presence of this factor in a feedback loop is highly probable as previous studies link its inclusion in implementation to avoid failure and transfer knowledge to the implementation team [1, 41].

Minimum level of customization to utilize ERP functionalities to a maximum (F11) is required for a successful implementation of financial systems. A high direct passive sum of 82.85 and high indirect passive sum of 85.56 signify high influence of the factor on implementation of financial ERP systems. This is supported by high stability value of 84.47 to indicate that it contributes highly to system stability which is also an indication of the strong impact that other factors in the system have on it. A driven score of 79.01 indicates further that the factor is strongly impacted on by other factors in the financial system which is an indication of success of intervention activities on the system. In addition, a very low precarious score of 15.20 attests to its stability in the system. Previous studies have posited that minimal level of customization should be instituted to prevent delays and dissatisfaction of ERP systems which indirectly implies the need for a vanilla financial ERP system [1, 42-44].

Package selection carefully and professionally selected (F12) for financial ERP systems is a highly significant factor with a driven score of 77.06. This score makes it a success gage of measure of intervention because it has most influenced on the system and is guided by internal impacts on the system. It contributes significantly to the stability of a financial ERP system implementation with a value of 76.20. Previous studies have revealed the focus on selecting the right package to be pivotal to successful implementation of ERP systems [1, 32, 38, 41]. The integration score of 52.83 substantiates the previous studies to postulate the right package selection for financial ERP system to be critical for implementation success.

Understanding the institutional culture which include the norms, values, and beliefs (F13) is regarded as important in this study in terms of criticality (50.08) and integration (55.39). The high passive sum of 62.81 for this factor leads to it been ranked as driven factor which aligns with the preceding studies that have considered it to be highly critical for successful implementation of ERP systems [1, 4, 5, 8, 13, 45, 46]. Consequently, in the context of financial ERP systems, the factor requires a monitoring intervention to be deployed. User involvement and participation throughout the implementation (F14) process of financial ERP systems has been established with a 100 percent passive sum, irrespective of low active sum. This factor presents high stability score of 77.36 and driven score of 80.17 which make it a major success factor for the implementation of financial ERP systems. The integration score of 56.38 supports the stability and driven measures of this factor with a low precarious impact score of 21.35. In this study, we have identified the factor to be highly

influential and driven in nature for monitoring the implementation of financial ERP systems because of its dependency on other factors [5].

ERP vendor support and partnership (F15) has emerged for a financial system implementation with a criticality score of 45.35. However, the low precarious score indicates that it is relevant for intervention because of its high reactive and driving score of 44.06. Previous authors have supported this factor to be important for the successful implementation of ERP systems in form of technical assistance, system maintenance, system update and user training [5, 42]. Business vision and plan (F16) for implementation of a financial ERP system presents low ranking and classification scores below the average values. However, the factor is reactive with an active sum of 41.73 and indirect passive sum of 48.61 above the average less the two-third standard deviation. This indicates that it impacts other factors in the system, it should be monitored and can be used for interventions because of its driving nature (45.81). In addition, a criticality score of 41.90 is observed in the system. A plan with clear vision and agreed upon objectives is reported by previous authors to be important for the successful implementation of ERP systems [9, 23, 34, 37, 47-53].

Adequate IT infrastructure (F17) for implementing financial ERP systems has emerged to be critical with a high driving score of 49.13, criticality score of 44.23 and precarious score of 43.75 over the average. It is an indication that it is reactive in nature and strongly influencing the system which make it suitable for interventions. Prior studies have ranked the adequacy of IT infrastructure high across implementations to indicate its significance [23, 50]. Monitoring management especially evaluation of performance

metrics (F18) is established in this study as driven in nature with a score of 76.89. Its high passivity indicates dependency on other factors in the system and provides an enabler of high stability (89.25). The low precarious score of 10.97 implies it could be used for intervention. This study upholds the prior studies that have elevated this factor to be critical for preventing adverse consequences for the implementation of ERP systems [4, 23, 50].

Allocating and dedicating valuable resources (F19) for implementing financial ERP systems is presented in this study to contradict the results of previous studies that have judged it as critical [54, 55]. Both low values for passive sum and active sum indicate that the role of this factor is neutral with no significant changes. Hence, its presence could likely contribute to the system stability (64.19). It is also an indication that this factor is torpid in nature. Data management plan that ensures that data are accurately and efficiently migrated to a new system and analyzed properly (F20) is presented in this study to be dependent on other factors in the system because of its high passive scores of 71.42 and 71.36. The integration score of 47.52 and critically score of 41.11 have established this factor to be important for implementing financial ERP systems. It is being driven in nature (64.83), it can indicate success, and offer stability (64.44) because of its low precarious score of 37.34.

In this study, CSFs for implementing financial ERP systems have been ranked and classified to structure a comprehensive model (Table VI) that can help organization to mediate a successful system implementation that takes into account the identified factors for supporting financial functions.

TABLE VI. ADVIAN® ANALYSIS OF FACTORS, HIGHLY CRITICAL (√√), ABOVE AVERAGE (√)

| Factor | Interaction   |                | Classification |             |           | Ranking    |         |        |
|--------|---------------|----------------|----------------|-------------|-----------|------------|---------|--------|
|        | Highly Active | Highly Passive | Criticality    | Integration | Stability | Precarious | Driving | Driven |
| F1     | √√            |                |                | √√          | √√        | √√         | √√      |        |
| F2     | √             | √              | √√             | √√          |           | √√         | √       |        |
| F3     | √             | √              | √√             | √√          |           | √√         | √       |        |
| F4     | √             |                |                |             | √         | √          | √√      |        |
| F5     | √             |                | √              |             |           | √          | √       |        |
| F6     |               | √√             |                | √           | √         |            |         | √√     |
| F7     | √             |                | √              |             |           | √√         | √       |        |
| F8     |               | √              | √              | √           |           | √          |         | √      |
| F9     | √             |                | √              |             |           | √          | √       |        |
| F10    |               | √              | √              | √           |           |            |         | √      |
| F11    |               | √√             |                |             | √√        |            |         | √√     |
| F12    |               | √√             |                | √√          | √√        |            |         | √√     |
| F13    |               | √√             | √√             | √√          |           | √          |         | √      |
| F14    |               | √√             |                | √√          | √√        |            |         | √√     |
| F15    | √             |                | √              |             |           | √          | √       |        |
| F16    | √             |                | √              |             |           | √          | √       |        |
| F17    | √             |                | √              |             |           | √√         | √       | √√     |
| F18    |               | √              |                |             | √√        |            |         |        |
| F19    |               |                |                |             |           |            |         |        |
| F20    |               | √              | √              | √           |           |            |         | √      |

The ranking has provided suitable factors that should be selected for system improvement and control of improvement success. The model has considered factors with the highest driving scores to be the best impact factors for intervention changes. Furthermore, selecting a smaller number of low precarious factors along with high driven impact factors as success indicators of intervention activities. The classification of CSFs for implementing financial ERP systems has been achieved and determined by their integration into the system as well as contribution to the system stability while proving the identified critical factors. The classification of factors has served as a good basis for observation with care of implementing financial ERP systems. The most significant measures of successful implementation of ERP systems are driven and driving CSFs. The driving CSFs exert impacts on the system which implies they should be considered as a good starting point for intervening activities. Hence, they are prime resources in an organization for improvement. However, a small number of the selected CSFs can be made based on the exclusion of high precarious CSFs when there are too high driving factors for suitable improvement. Furthermore, driven CSFs are the most influential in the system which makes them good indicators of success for intervening activities taken on driving resources. These factor characteristics make conditions of driven resources a good measurement of success for intervention activities and state of driving resources can be improved by intervention activities.

## VI. CONCLUSION

This study identifies the CSFs for successful implementation of an ERP system to support financial functions. This study has investigated the CSFs of ERP system implementation to support financial functions in the context of HEI using the ADVIAN® method. Our study delivers stimulating insights into the mindset regarding CSFs and implementation of financial ERP systems in the HEIs. This study has employed a rigorous scientific method to generally contribute to information system research and practice as most of the 20 factors identified can be applicable to any information system project. This paper contributes to the implementation of ERP system, financial system, and CSFs research by identifying 20 CSFs derived from the literature. The significance of each identified factor for implementing financial ERP systems has been investigated in this study. In consonance with the objective of this study, it is apparent that factors presented are highly relevant to the implementation of financial ERP systems. The determination of CSFs based on the calculations of ADVIAN® has revealed the contributions of the identified factors for successful implementation financial ERP systems. This study provides the management of finance with useful knowledge that can enable effective implementation strategies for supporting financial functions such as incorporating these factors into the planning, implementation and use of an ERP system that supports financial functions.

The limitations of this study are circumscribed by the availability of domain experts and underserved niche of ERP system implementation research that have delved deeply into the financial domain. However, most of the enterprise wide implementations of ERP systems are customized to the

specific settings that make them potentially viable for studies on larger enterprises to present different findings. In addition, the results of expert evaluation, ranking and classification of factors may turn out differently. Nevertheless, the use of a scientific rigorous approach is highly desirable for this research because only a small amount of scientific research has been conducted in this subject area. Consequently, the chosen approach represents an appropriate method for obtaining a preliminary overview and substantiation of CSFs for implementing ERP systems to support financial functions in the higher education.

The direct implication of this study indicates a clear chasm in the literature on implementation of financial ERP systems. It is in the interest of management to observe very clearly the identified the challenges in the process of implementing ERP systems and to employ holistic strategies to mitigate the challenges. There exists a chasm between what has been investigated within the field of ERP system implementation, CSFs, and financial sector. The chasm can be exploited further by researchers and practitioners to better understand what makes a successful implementation of ERP systems for supporting financial functions. Furthermore, regardless of the study grounded on the settings of HEIs, the outcome can very well be generalized to other similar organizations whose financial management are in dare need of reformation and is of importance. Since many HEIs in developing countries such as South Africa have implemented financial ERP systems it could be contended that the outcome of this study can be generalized geographically. Hence, a possible action for future research endeavor in the domain of ERP system.

## ACKNOWLEDGMENT

The first author would like to thank her fellow researchers for their inspirations and the Durban University of Technology for the support provided during the study.

## REFERENCES

- [1] A. Rabaa'i, "Identifying critical success factors of ERP Systems at the higher education sector," 2009.
- [2] A. I. ALdayel, M. S. Aldayel, and A. S. Al-Mudimigh, "The critical success factors of ERP implementation in higher education in Saudi Arabia: a case study," *Journal of Information Technology and Economic Development*, vol. 2, no. 2, pp. 1, 2011.
- [3] G. Tortorella, and C. E. Fries, "Reasons for adopting an ERP system in a public university in Southern Brazil."
- [4] M. A. Al-Hadi, and N. A. Al-Shaibany, "Critical success factors (CSFs) of ERP in higher education institutions," *International Journal*, vol. 7, no. 4, pp. 92-95, 2017.
- [5] R. C. Thompson, O. O. Olugbara, and A. Singh, "Deriving critical success factors for implementation of enterprise resource planning systems in higher education institution," *African Journal of Information Systems*, vol. 10, no. 1, 2018.
- [6] I. K. Sowan, R. Tahboub, and F. Khamayseh, "University ERP Preparation Analysis: A PPU Case Study," *International Journal of Advanced Computer Science And Applications*, vol. 8, no. 11, pp. 345-352, 2017.
- [7] J. Ram, D. Corkindale, and M.-L. Wu, "Implementation critical success factors (CSFs) for ERP: Do they contribute to implementation success and post-implementation performance?," *International Journal of Production Economics*, vol. 144, no. 1, pp. 157-174, 2013.
- [8] O. O. Olugbara, B. M. Kalema, and R. M. Kekwaletswe, "Identifying critical success factors: the case of ERP systems in higher education," 2014.

- [9] M. Y. M. Al-Sabaawi, "Critical success factors for enterprise resource planning implementation success," *International Journal of Advances in Engineering & Technology*, vol. 8, no. 4, pp. 496, 2015.
- [10] A. Trigo, F. Belfo, and R. P. Estébanez, "Accounting information systems: The challenge of the real-time reporting," *Procedia Technology*, vol. 16, pp. 118-127, 2014.
- [11] A. Y. Noaman, and F. F. Ahmed, "ERP systems functionalities in higher education," *Procedia Computer Science*, vol. 65, pp. 385-395, 2015.
- [12] A. Abugabah, and L. Sanzogni, "Enterprise resource planning (ERP) system in higher education: A literature review and implications," *International Journal of Human and Social Sciences*, vol. 5, no. 6, pp. 395-399, 2010.
- [13] J. Loonam, V. Kumar, A. Mitra, and A. Abd Razak, "Critical success factors for the implementation of enterprise systems: A literature review," *Strategic Change*, vol. 27, no. 3, pp. 185-194, 2018.
- [14] A. S. Shatat, "Critical success factors in enterprise resource planning (ERP) system implementation: An exploratory study in Oman," *Electronic Journal of Information Systems Evaluation*, vol. 18, no. 1, pp. 36-45, 2015.
- [15] M. Soliman, and N. Karia, "Antecedents for the Success of the Adoption of Organizational ERP Among Higher Education Institutions and Competitive Advantage in Egypt," *Engineering, Technology & Applied Science Research*, vol. 7, no. 3, pp. 1719-1724, 2017.
- [16] A. Epizitone, and O. O. Olugbara, "Critical success factors for ERP system implementation to support financial functions.," *Academy of Accounting and Financial Studies Journal* vol. 23, Issue 6, 2019.
- [17] V. Linss, and A. Fried, "Advanced Impact Analysis: the ADVIAN® method—an enhanced approach for the analysis of impact strengths with the consideration of indirect relations," 2009.
- [18] V. Linss, and A. Fried, "The ADVIAN® classification—A new classification approach for the rating of impact factors," *Technological Forecasting and Social Change*, vol. 77, no. 1, pp. 110-119, 2010.
- [19] B. Guertler, and S. Spinler, "Supply risk interrelationships and the derivation of key supply risk indicators," *Technological Forecasting and Social Change*, vol. 92, pp. 224-236, 2015.
- [20] A. Fried, "Performance measurement systems and their relation to strategic learning: A case study in a software-developing organization," *Critical Perspectives on Accounting*, vol. 21, no. 2, pp. 118-133, 2010.
- [21] V. A. Bañuls, M. Turoff, and S. R. Hiltz, "Collaborative scenario modeling in emergency management through cross-impact," *Technological Forecasting and Social Change*, vol. 80, no. 9, pp. 1756-1774, 2013.
- [22] J. L. Worrell, P. M. Di Gangi, and A. A. Bush, "Exploring the use of the Delphi method in accounting information systems research," *International Journal of Accounting Information Systems*, vol. 14, no. 3, pp. 193-208, 2013.
- [23] F. Campos Fernandes Leandro, M. P. Mexas, and G. M. Drumond, "Identifying critical success factors for the implementation of enterprise resource planning systems in public educational institutions," *Brazilian Journal of Operations & Production Management*, vol. 14, no. 4, pp. 529-541, 12, 2017.
- [24] V. Arnold, "The changing technological environment and the future of behavioural research in accounting," *Accounting & Finance*, vol. 58, no. 2, pp. 315-339, 2018.
- [25] H. Abdelghaffar, and R. H. A. Azim, "Significant factors influencing ERP implementation in large organizations: Evidence from Egypt." pp. 1-16.
- [26] C. Leyh, "Critical success factors for ERP projects in small and medium-sized enterprises-The perspective of selected German SMEs." pp. 1181-1190.
- [27] C. Leyh, "Critical success factors for ERP projects in small and medium-sized enterprises—the perspective of selected ERP system vendors," *Multidimensional Views on Enterprise Information Systems*, pp. 7-22: Springer, 2016.
- [28] S. Dezdard, and S. Ainin, "The influence of organizational factors on successful ERP implementation," *Management Decision*, vol. 49, no. 6, pp. 911-926, 2011.
- [29] F. Belfo, and A. Trigo, "Accounting information systems: Tradition and future directions," *Procedia Technology*, vol. 9, pp. 536-546, 2013.
- [30] M. A. AlSudairi, "Analysis and exploration of critical success factors of ERP implementation: a brief review," *International Journal of Computer Applications*, vol. 69, no. 8, 2013.
- [31] T. A. Sykes, V. Venkatesh, and J. L. Johnson, "Enterprise system implementation and employee job performance: Understanding the role of advice networks," *Mis Quarterly*, vol. 38, no. 1, 2014.
- [32] T. M. Somers, and K. G. Nelson, "A taxonomy of players and activities across the ERP project life cycle," *Information & Management*, vol. 41, no. 3, pp. 257-278, 2004.
- [33] F. Fui-Hoon Nah, J. Lee-Shang Lau, and J. Kuang, "Critical factors for successful implementation of enterprise systems," *Business process management journal*, vol. 7, no. 3, pp. 285-296, 2001.
- [34] S. Finney, and M. Corbett, "ERP implementation: a compilation and analysis of critical success factors," *Business Process Management Journal*, vol. 13, no. 3, pp. 329-347, 2007.
- [35] Y. B. Moon, "Enterprise Resource Planning (ERP): a review of the literature," *International journal of management and enterprise development*, vol. 4, no. 3, pp. 235-264, 2007.
- [36] S. Dezdard, and A. Sulaiman, "Successful enterprise resource planning implementation: taxonomy of critical factors," *Industrial Management & Data Systems*, vol. 109, no. 8, pp. 1037-1052, 2009.
- [37] L. Shaul, and D. Tauber, "Critical success factors in enterprise resource planning systems: Review of the last decade," *ACM Computing Surveys (CSUR)*, vol. 45, no. 4, pp. 55, 2013.
- [38] R. G. Saade, and H. Nijher, "Critical success factors in enterprise resource planning implementation: a review of case studies," *Journal of Enterprise Information Management*, vol. 29, no. 1, pp. 72-96, 2016.
- [39] K. Al-Fawaz, T. Eldabi, and A. Naseer, "Challenges and influential factors in ERP adoption and implementation," 2010.
- [40] Z. Alias, E. Zawawi, K. Yusof, and N. Aris, "Determining critical success factors of project management practice: A conceptual framework," *Procedia-Social and Behavioral Sciences*, vol. 153, pp. 61-69, 2014.
- [41] M. Al-Mashari, A. Al-Mudimigh, and M. Zairi, "Enterprise resource planning: A taxonomy of critical factors," *European Journal of Operational Research*, vol. 146, pp. 352-364, 2003.
- [42] E. Reitsma, and P. Hilletoft, "Critical success factors for ERP system implementation: a user perspective," *European Business Review*, vol. 30, no. 3, pp. 285-310, 2018.
- [43] E. Ziemba, and I. Oblak, "Critical success factors for ERP systems implementation in public administration." pp. 1-19.
- [44] H. Akkermans, and K. van Helden, "Vicious and virtuous cycles in ERP implementation: a case study of interrelations between critical success factors," *European journal of information systems*, vol. 11, no. 1, pp. 35-46, 2002.
- [45] S. Venkatraman, and K. Fahd, "Challenges and success factors of ERP systems in Australian SMEs," *Systems*, vol. 4, no. 2, pp. 20, 2016.
- [46] A. Y. Aremu, A. Shahzad, and S. Hassan, "Determinants of Enterprise Resource Planning Adoption on Organizations' Performance Among Medium Enterprises," *LogForum*, vol. 14, no. 2, 2018.
- [47] N. Ahmad, A. Haleem, and A. A. Syed, "Compilation of critical success factors in implementation of enterprise systems: a study on Indian organisations," *Global journal of flexible systems management*, vol. 13, no. 4, pp. 217-232, 2012.
- [48] M. Ashja, A. H. Moghadam, and H. Bidram, "Comparative study of large information systems' CSFs during their life cycle," *Information Systems Frontiers*, vol. 17, no. 3, pp. 619-628, 2015.
- [49] O. François, M. Bourgault, and R. Pellerin, "ERP implementation through critical success factors' management," *Business Process Management Journal*, vol. 15, no. 3, pp. 371-394, 2009.
- [50] P. Hanafizadeh, R. Gholami, S. Dadbin, and N. Standage, "The core critical success factors in implementation of enterprise resource planning systems," *International Journal of Enterprise Information Systems (IJEIS)*, vol. 6, no. 2, pp. 82-111, 2010.

- [51] E. W. Ngai, C. C. Law, and F. K. Wat, "Examining the critical success factors in the adoption of enterprise resource planning," *Computers in industry*, vol. 59, no. 6, pp. 548-564, 2008.
- [52] J. Ram, M.-L. Wu, and R. Tagg, "Competitive advantage from ERP projects: Examining the role of key implementation drivers," *International Journal of Project Management*, vol. 32, no. 4, pp. 663-675, 2014.
- [53] J. M. Denolf, J. H. Trienekens, P. N. Wognum, J. G. van der Vorst, and S. O. Omta, "Towards a framework of critical success factors for implementing supply chain information systems," *Computers in industry*, vol. 68, pp. 16-26, 2015.
- [54] P. Saa, O. Moscoso-Zea, A. C. Costales, and S. Luján-Mora, "Data security issues in cloud-based Software-as-a-Service ERP." pp. 1-7.
- [55] D. Saxena, and J. McDonagh, "Yet Another 'List' of Critical Success 'Factors' for Enterprise Systems: Review of Empirical Evidence and Suggested Research Directions."

# Machine Learning based Analysis on Human Aggressiveness and Reactions towards Uncertain Decisions

Sohaib Latif<sup>1</sup>, Abdul Kadir Abdullahi Hasan<sup>2</sup>, Abdaziz Omar Hassan<sup>3</sup>

School of Mathematics and Big Data  
Anhui University of Science and Technology  
Huainan, China

**Abstract**—Tweet data can be processed as a useful information. Social media sites like Twitter, Facebook, Google+ are rapidly growing popularity. These social media sites provide a platform for people to share and express their views about daily routine life, have to discuss on particular topics, have discussion with different communities, or connect with globe by posting messages. Tweets posted on twitter are expressed as opinions. These opinions can be used for different purposes such as to take public views on uncertain decisions such as Muslim ban in America, War in Syria, American Soldiers in Afghanistan etc. These decisions have direct impact on user's life such as violations & aggressiveness are common causes. For this purpose, we will collect opinions on some popular decision taken in past decade from twitter. We will divide the sentiments into two classes that is anger (hatred) and positive. We will propose a hypothesis model for such data which will be used in future. We will use Support Vector Machine (SVM), Naive Bayes (NB), and Logistic Regression (LR) classifier for text classification task. Further-more, we will also compare SVM results with NB, LR. Research will help us to predict early behaviors & reactions of people before the big consequences of such decisions.

**Keywords**—Opinion mining; Naïve Bayes; linear regression; support vector machine

## I. INTRODUCTION

Internet is providing all the services a normal user looking for. Starting from the health, education, government and business, all categories of modern life have been covered in the shape of internet. Internet provides connectivity [1,2] between people and information publicly shared globally. Similarly, social media such as Facebook, Twitter, YouTube are platform to remain updated with current news and a airs. Through social media people [3,4,5] can share news, share their opinions and participate in activities being held online. Social Networking Sites (SNS) such as Twitter and Facebook have a beneficial effect on our way of life. SNS has been used for expressing opinions on different issues. In this work, we propose a sentiment based method for the predication of aggressive estimation.

In the age of technology [6,7] millions of people are using social media sites like Facebook, Twitter, Google Plus, etc. to share and express their views, emotions, and opinion about their daily lives. Through the online communities, we get an interactive media where consumers inform and influence

others through forums. Social media are now become rich of data in the form of tweets, status updates, posts, blog, comments, reviews, etc. [8, 9]. These social sites are not just using for personal use, but now it become a fastest tool to reach the people. It provides an opportunity for businesses by giving a platform to connect with their customers for advertising. Mostly people rely on user generated content or reviews to a great extent for decision making. The online content generated by users is too rich to analyze by normal user. The thing is to automate the process to take the views of user's as opinion. The online contents are mainly consider as opinions, sentiments, attitudes, and emotions [10].

## II. LITERATURE REVIEW

Machine Learning, Data mining and Natural Language processing all used together for the classifications of text documents widely. These three techniques also used to discover patterns from the electronic documents. Text mining is used to discover hidden useful information from the documents and deals with the operations like, retrieval, classification (supervised, unsupervised and semi supervised) and summarization [11].

There have been many e orts regarding text classifications in the past. Krishna and Gonghzu [12] have analyzed large data from clinics and try to find the clinical disorders. Sonia and Shruti [15] have used Machine Learning techniques for analysis of social network E-Health data. Roshan and Rio D Souza [13] have analyzed product value using sentimental analysis publicly given on Twitter. Both have [16] worked to solve the problem of reading millions of reviews by a single user for a particular product, they have developed a model using reviews posted which gives product classification in term of positive, negative and neutral reviews. In the same context, Barnaghi et al. [14] used Twitter sentiments to predict event winner. They used Bayesian Logistic Regression (BLR). They manually labelled tweets into two categories positive and negative. A model proposed [17,19] by them can be used to predict winner of any event using sentiments. In our research we will propose a methodology to analyses the pattern of human behaviors towards uncertain decisions. Our proposed methodology saves time and cost for such a huge public review posted daily on social networks. Nirbhay Kashyap et al. [20] have worked on music lyrics to categorize the mood of individuals. They have used different text mining and data

mining approaches to deal with such a problem. They have considered music associations, melody choice and music proposal as a feature to demonstrate the data. It is beneficial for predicting more accurate understanding of the music mood in the mood mapping process. Similarly, many studies have been found to investigate the online business trends using social data. Online business and larger companies' world-wide used user feedback which has been given on social sites for the improvement of product and business need with the passage of time. The amount of text and information shared on twitter in the form of tweets have valid information and it can be used to track the progress of product. They have categorized the data into different categories such as against, positive and negative and used machine learning clustering algorithms to do so. They have found that the data available online can be used for the process of information extraction and it is beneficial for the companies to track the progress of their product and handy for future considerations [21].

Santoshi et al. [22] have used twitter data differently. They have tried to figure out the user behavior towards political parties. They have captured twitter data before the election and categorized the raw data into 5 different categories such as positive, negative, happy, sad and neutral. This type of information is very handy for political parties before the election. It is also effective to solve the real problems of people so that you can change the thinking of users. They have considered BJP and INC for their purpose. These are the biggest political parties in India. Using text mining and unsupervised lexical method classified tweets related to these parties to identify people emotions for the parties.

Xin Li [23] have adopted the same platform for his studies with his group mates. They have used different Natural language processing techniques [25] for the awareness of social issues human facing. Social awareness information is analyzed by applying text mining and social network analysis.

AK Rathore et al. [24] has collected twitter data for the prediction of Pizza success after its launch. It is very handy information they have worked. This type of methodologies can be used to predict the behavior of any user for a particular product. Rathore and his company has used R and NodeXL for analyzing tweets collected from twitter. Furthermore, they have used different text mining, Natural Language Processing and Network Analysis techniques to predict user behavior. Any company or food delivering company can use this sort of information for the purpose of success and failure of product. Nobody has worked to analyze the behavior of certain decision and their impact of human life before. In our research we will propose a methodology to analyse the pattern of human behaviors towards uncertain decisions. Our proposed methodology saves time and cost for such a huge public review posted daily on social networks.

### III. PROPOSED METHOD

The solution we suggest involves Twitter data. Tweets collected with Twitter Search API [18]. Our methodology consists of two steps: training and testing phases. Feature representation and tweets collection and classifier training comes in training phase, while the testing phase have four phases: tweets collection for testing, feature representation,

hypothesis prediction and evaluation. The first two tasks (i.e. tweets collection and feature representation) are shared between training and testing phase. Some popular classifiers such as SVM, NB and LR used in training and hypothesis. We have used WEKA tool for training and testing of our proposed methodology. Firstly, we divided the data sets into two parts, training data and secondly testing data.

#### A. Preprocessing

Preprocessing reshape the data into desired form. The data collected is not purified for the process of classification, for this we have applied data Processing methodologies to transform the data into meaningful features.

Fig. 1 is showing the training of dataset. This involves mainly tokenization (or featurizing), feature weighting and data cleaning (removal of irrelevant features). Once the data is collected, URLs from the tweets and replies were removed. Data only with image or with a link but there was no textual information was also removed. Stop words also do not give any information about topic and just create noise in the data so using stop word-list they were also removed from the data. Pre-processing is the key process in data classification tasks. It also improves the effectiveness of proposed classifier. When data is pre-processed it helps in saving classifier time while classifying. Collected tweets are further pre-processed with following steps.

1) *Tokenization*: Tokenization deals with breaking of long text strings into substrings which may include phrases and words collectively known as tokens. Among two ways of tokenization (phrase and word tokenization), word-level tokenization is considered as more effective due to statistical significance. In this process, the sentence for instance "Trump is mentally disturbed person" was broken into tokens "Trump", is, mentally, disturbed, person. The algorithms which are used to tokenize a sentence separate the tokens with whitespace and some are based on built-in dictionary. Text can be tokenized in two ways, by words (often called bag of words) or phrases.

2) *Feature Weighting*: A standard function to compute the weights is TF-IDF. TF-IDF scheme is based on two parts: TF and IDF. TF stands for term frequency which is used to count the represented terms/tokens in a document. It can give a complete measure of term occurrence. IDF stands for inverse document frequency of a term in a collection of documents.

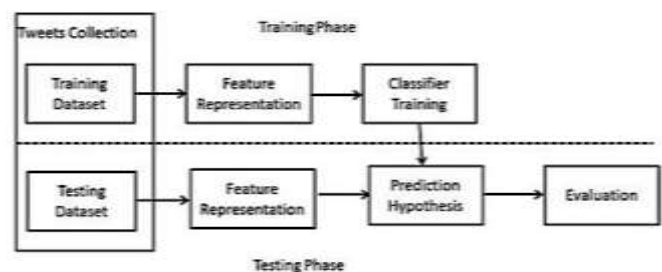


Fig. 1. Blocked Diagram of Proposed Methodology.



B. Sentiment Classification

Once we applied the pre-processing, we have data in a suitable format to apply classification algorithm on it. We have categorized the data into two formats. A data with false words labeled as Negative and data with positive words labeled as Positive. A sample of tweets rows which we have labeled. Different algorithms are available in this domain that can be used to train the classification task. Different experimental studies have been directed to analyze these methods for text categorization.

Once we applied the pre-processing, we have data in a suitable format to apply classification algorithm on it. We have categorized the data into two formats. A data with false words labeled as Negative and data with positive words labeled as Positive. A sample of tweets rows which we have labeled. Different algorithms are available in this domain that can be used to train the classification task. Different experimental studies have been directed to analyze these methods for text categorization.

IV. CLASSIFICATION

Supervised classification is a machine learning approach in which training data are used to construct the model and test data are used to evaluate the constructed model on unseen data to measure the performance of algorithm. There are a number of classifiers that exist to classify data, and below in Table I we will discuss the classifiers which we have explored in this work.

TABLE I. TWEETS ROWS WITH LABELED POSITIVE AND NEGATIVE SAMPLES

| Sr. No | Positive Samples  | Negative Samples   |
|--------|---|--|
| 1      | RT @joshua_landis: Major U.S. polycystatement on #Syria by Gen. Mattis: "US tofight IS in Syria until IS declares thatthey're done." Als...[War in Syria] | RT @Palespanish: Saudi Arabia: *Played a huge part in destabilizing Syria [War inSyria]  |
| 2      | catalonia election spain s king felipe warnsseparatists many truly seek god s mess[Catalonia]   | not only that mexican people areindigenous to north america any nativeshould be against the wall [NoBanNoWall]                             |
| 3      | israeli terrorist politician harasses palesoren hazan rightwing israeli known forpublicity stunts was f**k [Jerusalem]                                    | so in your opinion jerusalem is in whichcountry [Jerusalem]  |
| 4      | #Modi RafaleScam msg to cntry..thinkbig..do big.. forget ppl #NoteBan#notebandi #amitshahkiloat [Notebandi]   | Bjp is likely to show sunny leone ji CD sothat people will forget about GST & NOTEBANDI [Notebandi]  |
| 5      | RT @newsbusters: Says it all! The liberalmedia are STILL obsessed with trashinghttps://t.co/Y1HWxcpesP[Trump Victory]                                     | RT @leedsgarcia: The administration let DACA renewals sit inmailboxes — and then rejected them for being "late" https://t.co/4LYHOVictory] |

SVM provides better results than other Machine Learning algorithms in sense larger boundary distributions. SVM also supports high dimensional data. SVM is suitable for millions of features at the same time. SVM also supports optimization problems. Software libraries present for the implementation of SVM are lib-linear, libsvm. In logistic regression function, we have the hypothesis below, and sigmoid activation function.

Nave Bayes is probabilistic classifier which strongly based on Bayes Theorem. Simple Bayes, Independence Bayes are common names which are used. It is mostly used in classifying text information into their respective categories. There are some other example which are associated with the classifier such as to check either email is spam or not, either emails is related to sports or not.

V. EXPERIMENTAL SETUP

A. Evaluation Measures

We used various evaluation measures to assess the results, and these measures are described below in Table II.

The results of sentiment classification using Logistic classification are given in Table III. Precision, recall, and f-measure are approximately 83%, 84%, and 84%, respectively.

Here we have given the results of sentiment classification. The results of sentiment classification using Support Vector Machine (SVM) classification are given in Table IV. Precision, recall, and f-measure are approximately 92%, 85%, and 88%, respectively.

The results of sentiment classification using Naïve Bayes (NB) classification are given in Table V. Precision, recall, and f-measure are approximately 85%, 86%, and 85%, respectively.

TABLE II. BAG OF WORDS USED FOR CLASSIFICATIONS

| Sr. No | Hash Tag            | Words  |
|--------|---------------------|--|
| 1      | #Trump              | shithole, criminal trump, trump shutdown, turned sour, reject, failed socialist, evils choice, Moscow's victory, what nonsense, shocking crimes, disgrace, slap, Pathetic  |
| 2      | #SaudiWomenCanDrive | condemned, cheesy, car accident, lack of intellect, Protests ,condemnation, break,license, disasters   |
| 3      | #PanamaVerdict      | squeezed, shameless, Patwari, Haram Family, corruption, trashed, bark, barking, gangster, anti, wolf   |
| 4      | #NoteBandi          | waste of time, economy mess up, destroyed, black money, laundering, su er, cor-ruption, e ect, common man, impact, a ected, self-destructive, slap, idiot, com-plaining, ruthless, corrupted gov, history, disasters, stunts, nashbandi, Nobandi |
| 5      | #NoBanNoWall        | attacks, hurt, worse, wasteful, monument, hurdle, damage, environment, hate, break the wall, harmed, wreck less, darkness, must resist, rapists, hurt brownpeople, discriminatory ban stupid wall, resist  |

TABLE III. RESULTS OBTAINED USING LOGISTIC REGRESSION

| Sr. No | Class              | Precision | Recall | F-Measures |
|--------|--------------------|-----------|--------|------------|
| 1      | Brexit             | 0.814     | 0.763  | 0.784      |
| 2      | Catalonia          | 0.879     | 0.885  | 0.882      |
| 3      | PanamaVerdict      | 0.881     | 0.872  | 0.876      |
| 4      | NoteBandi          | 0.86      | 0.822  | 0.844      |
| 5      | Jerusalem          | 0.86      | 0.817  | 0.834      |
| 6      | SaudiWomenCanDrive | 0.857     | 0.85   | 0.855      |
| 7      | Trump              | 0.856     | 0.864  | 0.86       |
| 8      | MuslimBan          | 0.855     | 0.854  | 0.856      |
| 9      | NoBanNoWall        | 0.89      | 0.891  | 0.891      |
| 10     | SyriaWar           | 0.863     | 0.879  | 0.869      |

TABLE IV. RESULTS OBTAINED USING SVM CLASSIFIER

| Sr. No | Class              | Precision | Recall | F-Measures |
|--------|--------------------|-----------|--------|------------|
| 1      | Brexit             | 0.865     | 0.883  | 0.863      |
| 2      | Catalonia          | 0.841     | 0.908  | 0.873      |
| 3      | PanamaVerdict      | 0.922     | 0.925  | 0.906      |
| 4      | NoteBandi          | 0.900     | 0.899  | 0.889      |
| 5      | Jerusalem          | 0.926     | 0.929  | 0.921      |
| 6      | SaudiWomenCanDrive | 0.911     | 0.900  | 0.901      |
| 7      | Trump              | 0.860     | 0.817  | 0.834      |
| 8      | MuslimBan          | 0.850     | 0.873  | 0.843      |
| 9      | NoBanNoWall        | 0.916     | 0.924  | 0.909      |
| 10     | SyriaWar           | 0.922     | 0.921  | 0.905      |

TABLE V. RESULTS OBTAINED USING NAÏVE BAYES

| Sr. No | Class              | Precision | Recall | F-Measures |
|--------|--------------------|-----------|--------|------------|
| 1      | Brexit             | 0.817     | 0.821  | 0.819      |
| 2      | Catalonia          | 0.864     | 0.846  | 0.855      |
| 3      | PanamaVerdict      | 0.866     | 0.856  | 0.862      |
| 4      | NoteBandi          | 0.899     | 0.889  | 0.879      |
| 5      | Jerusalem          | 0.872     | 0.886  | 0.875      |
| 6      | SaudiWomenCanDrive | 0.866     | 0.867  | 0.865      |
| 7      | Trump              | 0.819     | 0.815  | 0.817      |
| 8      | MuslimBan          | 0.798     | 0.824  | 0.809      |
| 9      | NoBanNoWall        | 0.875     | 0.880  | 0.878      |
| 10     | SyriaWar           | 0.864     | 0.885  | 0.869      |

Precision (Positive Predictive value) can be defined as relevant instances from the retrieved instances. The concept is used for binary classifications. Whereas recall is the number of relevant instances from total number of relevancy. This is also known as sensitivity.

To get good performance of classifier precision and recall are often used together [28]. F-Measure can be defined as harmonic mean of precision and recall.

### B. Tools for Evaluation

To perform desired task, we used WEKA. WEKA is open source free software which has been used for various machine learning problems using data. It contains tools which can be used for classifications, pre-processing, clustering, visualization, association rules etc. Machine Learning is nothing without giving an artificial intelligence to your data. Machine learning methods are very similar to data mining algorithms. WEKA has collection of Machine Learning (ML) algorithms which are applied on data to extract desired results from it.

### C. Comparative Analysis

A comparison analysis of classifiers for sentiment classification is given in Table VI. We can see that SVM provides best results and it gives approximately 88% F-measure which is much better than from NB and LR results.

TABLE VI. COMPARATIVE ANALYSIS OF RESULTS OBTAINED USING ALL THREE CLASSIFIERS

| Class | Precision | Recall | F-Measures |
|-------|-----------|--------|------------|
| SVM   | 0.92      | 0.901  | 0.899      |
| NB    | 0.807     | 0.811  | 0.801      |
| LR    | 0.799     | 0.789  | 0.799      |

## VI. DISCUSSIONS

In last chapter we have described tools, data source, and different technologies that we have used in our approach. In this chapter we will present the obtained results. Fig. 2 is showing the work flow of our experimentation. Three classifiers Support Vector Machine, Naive Bayes and Logistic Regression are used in our experiment and to measure the effectiveness of each classifier we have used three measurements i.e. recall, precision, and f-measure by applying standard 10-folded cross-validation.

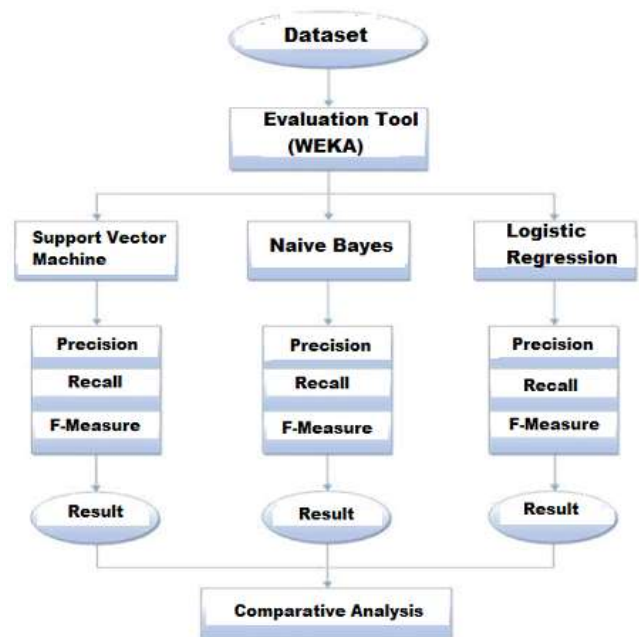


Fig. 2. Workflow of Experimentation.

## VII. CONCLUSION

Twitter is one of the most important social sharing platform for useful information. Tweets posted on twitter are expressed as opinions. These opinions can be used for different purposes such as to take public views on uncertain decisions such as Muslim ban in America, War in Syria, American Soldiers in Afghanistan, etc. These decisions have direct impact in users life such as violations & aggressiveness are common causes. We have collected tweets of such decisions and labeled the tweets into two categories such as anger (hatred) and positive. We have used classifier algorithms such as Support Vector Machine (SVM), Naive

Bayes (NB), and Logistic Regression (LR) for building models. We have also compared SVM results with NB, LR. This research is useful for predicting early behaviors & reactions of people before the big consequences of such decisions.

In the future we interested to build a tool which can work as a recommender system to classify tweets automatically into two categories such as Anger and Positive.

$$a^2 + b^2 = c^2$$

#### ACKNOWLEDGMENT

This work is partially supported by the National Natural Science Foundation of China (No.61572035, No.61902002, No.61402011 and No.61272153), the Natural Science Foundation of Educational Government of Anhui Province of China(No.KJ2016A208), Anhui Provincial Big Data Foundation(No. 2017032), the Foundation for top-notch academic disciplines of Anhui Province(No.gxbjZD11),the Academic and Technology Leader Foundation of Anhui Province(No.2019H239), the Open Project Program of Key Laboratory of Embedded System and Service Computing of Ministry of Education(No.ESSCKF2018-04). We also gratefully acknowledge the helpful comments and suggestions of the reviewers, which have improved the presentations.

In the end, researcher would like to pay thanks to my father and mother, especially because both of them suffer a lot during my studies.

#### REFERENCES

- [1] Liu.B, Sentiment analysis and opinion mining Synthesis lectures on human language technologies, 5(1), 1-167, 2012.
- [2] Khan.F.H, Bashir.S. And Qamar.U. TOM: Twitter opinion mining framework using hybrid classification scheme Decision Support Systems, 57, 245-257, 2014.
- [3] Lin.C, And He.Y. Joint sentiment/topic model for sentiment analysis. In Proceedings of the 18th ACM conference on Information and knowledge management, pp. 375-384, ACM, 2009.
- [4] Jiang. L, Yu. M, Zhou. M, Liu.X and Zhao.T .Target-dependent twitter sentiment classification. In Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies- Association for Computational Linguistics .Volume 1, pp. 151-160, 2011.
- [5] Barbosa.L, and Feng. J. Robust sentiment detection on twitter from biased and noisy data. In Proceedings of the 23rd International Conference on Computational Linguistics Association for Computational Linguistics Posters, pp. 36-44. 2010.
- [6] Torunoglu.D, Telseren.G, Sagturk.O and Ganiz.M.C. Wikipedia based semantic smoothing for twitter sentiment classification. In Innovations in Intelligent Systems and Applications (INISTA), 2013 IEEE International Symposium on pp. 1-5, IEEE. 2013.
- [7] Go. A, Bhayani.R and Huang.L. Twitter sentiment classification using distant supervision, CS224N Project Report, Stanford, 1(2009), 12. 2009.
- [8] Liu. K.L, Li.W.J and Guo.M. Emoticon smoothed language models for twitter sentiment analysis In AAAI.publications, 2012.
- [9] Bravo Marquez.F, Mendoza.M and Poblete.B Combining strengths, emotions and polarities for boosting Twitter sentiment analysis. In Proceedings of the Second International Workshop on Issues of Sentiment Discovery and Opinion Mining (p. 2) ACM, 2013.
- [10] Goncalves.P, Araujo.M, Benevenuto.F and Cha.M. Comparing and combining sentiment analysis methods. In Proceedings of the first ACM conference on Online social networks pp. 27-38, ACM, 2013, October.
- [11] Khan. A, Baharudin.B, Lee.L. H and Khan, K. A review of machine learning algorithms for text-documents classification. Journal of advances in information technology 1(1), 4-20, 2010.
- [12] Chodey.K. P and Hu.G. Clinical text analysis using machine learning methods. In Computer and Information Science (ICIS), 2016 IEEE/ACIS 15th International Conference on pp. 1-6. IEEE. 2016.
- [13] Fernandes.R and D'Souza. R Analysis of product Twitter data through opinion mining. In India Conference (INDICON), 2016 IEEE Annual, pp. 1-5 .IEEE, 2016.
- [14] Barnaghi.P, Ghaffari.P and Breslin.J. G. Opinion mining and sentiment polarity on twitter and correlation between events and sentiment. In Big Data Computing Service and Applications (Big Data Service), 2016 IEEE Second International Conference on pp. 52-57, IEEE, 2016.
- [15] Saini.S and Kohli.S. Machine learning techniques for effective text analysis of social network E-health data. In Computing for Sustainable Global Development (INDIA.Com), 2016 3rd International Conference on pp. 3783-3788, IEEE, 2016.
- [16] Nawaz.M.S, Bilal.M, Lali.M.I. Ul Mustafa.R, Aslam.W and Jajja.S. Effectiveness of Social Media Data in Healthcare Communication. Journal of Medical Imaging and Health Informatics, 7(6), 1365-1371, 2017.
- [17] Liu.B and Zhang.LA survey of opinion mining and sentiment analysis. In mining text data Springer US pp. 415-463, 2012.
- [18] Twitter Search API. Retrieved from website: <https://dev.twitter.com/rest/public/search>, 2015.
- [19] Hall.M, Frank.E, Holmes.G, Pfahringer.B, Reutemann.P and Witten. I. H. The WEKA data mining software an update ACM SIGKDD explorations newsletter, 11(1), 10-18, 2009.
- [20] Kashyap, Nirbhay et al. Mood Based Classification of Music by Analyzing Lyrical Data Using Text Mining Micro-Electronics and Telecommunication Engineering (ICMETE), 2016 International Conference on. IEEE, 2016.
- [21] Zunic, Emir, AlmirDjedovic, and DzenanaDonko. Application of Big Data and text mining methods and technologies in modern business analyzing social networks data about traffic tracking. Telecommunications (BIHTEL), 2016 XI International Symposium on IEEE, 2016.
- [22] Kuamri, Santoshi and NarendraBabu.C. Real time analysis of social media data to understand people emotions towards national parties, 2017 8th International Conference on Computing, Communication and Networking Technologies (ICCCNT).IEEE, 2017.
- [23] Vespignani, Alessandro. Predicting the behavior of techno-social systems. Science 325.5939, 425-428, 2009.
- [24] Rathore, Ashish Kumar and VigneswaraIlavarasan.P. Social media analytics for new product development Case of a pizza. Advances in Mechanical, Industrial Automation and Management Systems (AMIAMS), 2017 International Conference on IEEE, 2017.
- [25] Ricardo. A, Calix, Automated semantic understanding of human emotions in writing and speech, 2011.
- [26] Pang.B , Lee.L, Vaithyanathan.S, Thumbs up sentiment classification using machine learning techniques .in Proceedings of the ACL-02 conference on Empirical methods in natural language processing- Volume 10, pp. 9-86, 2002.
- [27] Janyce.M.Wiebe, Bruce.Rebecca.F, O'Hara and Thomas P. Development and use of a gold-standard data set for subjectivity classifications in Proceedings of the 37th annual meeting of the Association for Computational Linguistics on Computational Linguistics, pp. 246—253, 1999.
- [28] Full list of bad words and top swear banned by Google Available at: <https://www.freewebeheaders.com/full-list-of-bad-words-banned-by-google/>Lastaccessed ,2019.

# A Cluster-Based Mitigation Strategy Against Security Attacks in Wireless Sensor Networks

Jahangir Khan<sup>1</sup>

Community College  
Department of Information Systems  
King Khalid University, Tehama Branch  
Kingdom of Saudi Arabia

Ansar Munir Shah<sup>2</sup>

Department of Computer Science and IT  
Institute of Southern Punjab (ISP)  
Multan, Pakistan

Babar Nawaz<sup>3</sup>

Department of Computer Science  
IIC University of Technology  
Phnom Penh  
Cambodia

Khalid Mahmood<sup>4</sup>

College of Science and Arts  
Department of Information Systems  
King Khalid University Tehama Branch  
Kingdom of Saudi Arabia

Muhammad Kashif Saeed<sup>5</sup>

Community College  
Department of Information Systems  
King Khalid University, Tehama Branch  
Kingdom of Saudi Arabia

Mehmood ul Hassan<sup>6</sup>

Department of Computer Skills  
Najran University  
Kingdom of Saudi Arabia

**Abstract**—Wireless Sensor Networks (WSNs) applications range across distinct application comprising of event detection at real-time. WSNs can be deployed for not only mobile nodes but also for static sensor nodes (SNs) for various applications which may include health care system, smart parking, environmental monitoring etc. Sensor nodes in WSN are constrained in terms of energy contents of each node and can be accessible by other nodes in a wireless medium are more likely to be susceptible to various categories of attacks. Wireless Network are more likely prone to various kinds of security attacks, one such type of attack caused by a malicious attacker, which can result to decay in the lifetime of the network and an adverse scenario can even lead to congestion in the entire network. This paper presents the overview of various attacks and their consequences on different layers and evaluates defense strategy used to mitigate the various categories of attacks on Wireless Sensor Networks. This study proposes a cluster-based approach for each node of a WSN where the nodes of network constrained by energy can organize and perform network duties as per the network performance for this one node performs the role of cluster head (CH) which is elected on the basis of the "Reputation" of a node which is an indicator of nodes individual behavior in the network and "Net\_Credit\_Score" which determines the cooperating behavior of sensor node in the cluster. Further, this study highlights few parameters which can be implemented to further enhance the defense strategy by taking into account the factors such as Cluster count, Stability factor of both the Cluster and Cluster Head and Intra-Cluster topology which can be crucial. This will result in formulating a road map for designing a secure and resistant reputation-based system for WSN to overcome the various security related attacks.

**Keywords**—Wireless sensor network; security attacks; security issues; clusters

## I. INTRODUCTION

Wireless Sensor Networks (WSN) can be defined as a versatile communications system that makes use of the wireless medium (radio frequency) in order to transmit and receive data, therefore reducing their dependency on wired connections. It can be termed as a group of spatially dispersed and dedicated sensors deployed to monitor and collect the details such as the physical conditions existing in a region such as the collection of data, scientific examination, military applications etc. But the sensor nodes are constrained due to various factors such as security issues and limited resource energy; they are more likely to be vulnerable to security attack. Security threats can easily affect the WSNs. Sensor network attackers do not need to be constrained by the characteristic of sensor nodes since they are able to use costly transceivers and power supply nodes, making it possible for this type of network to be affected. The local storage ability of sensor network is very minimal [1], so in order to run very complex cryptology protocols, security mechanisms for sensor networks can not enable each sensor node to store long-sized keys. Sensor nodes have low power consumption, so power preservation must be the priority of sensor network protocols. Usually, sensor networks comprise of a huge number of communication nodes, do not have a global identification number and might face simple node failure.

The purpose of the attacker is to render target domains inaccessible to legitimate users. In several areas a sensor network without adequate security against attacks will not be deployable. Wireless networks are more vulnerable to security attacks due to the transmission medium being broadcasting in nature. The security attacks [2] are the situation of violation of

system security, carried out by an intelligent entity and is a deliberate act of threat on an organization(system) resource or service. We need to design a security algorithm for secure working condion.

This paper considers cluster-based strategy for mitigation against security attacks carried out on WSNs and proposes a cluster-based approach which is Fault-tolerant and by categorizing the sensor nodes into cluster also considers the balancing of network load.

The paper is organized as follows. In Section II we present related work to various security attacks in detail. The Security requirements of WSNs is explained in Section III. Section IV gives a review of attacks on WSN. A proposed cluster-based mitigation strategy is described in Section V. Section VI presents conclusions and future enhancements related to security attacks on WSN have been considered.

## II. RELATED WORK

The paper [3] examines and evaluate the various kinds of attacks carried out on WSN. The main focus of this study is to examine how such attacks can be prevented for WSNs by creating a sound understanding of various kinds of attacks in WSNs. In [4] the authors have conducted a review on DDoS attack to present its impact on networks and to present various defensive, detection and preventive measures which can be adopted in order to mitigate attacks on WSNs. Various parameters related to methods used for selection of clusters [5] need of re-clustering [6] and study of the QoS parameters such as performance [7] of nodes in WSN. The approach used in [8] determines that the cluster head is selected on the basis of a threshold value "T" which can be calculated using the remaining energy and relative position of the node in the network. In the study [9] CH is elected based on assigning weight factor to the nodes such as Reputation-based system such as RFSN [10] and DRBTS [11], energy, mobility and distance between the nodes and based on the weighted value of these three parameters Cluster head can determine.

1) *Physical layer*: The attack primarily focused on this layer that may affect (leading to starvation) or may not affect (resulting in sniffing) the physical environment needed to send the data.

2) *Data link layer*: DDoS attacks (active as well as passive) can be carried out resulting in increase in the packet drop or in adverse situation may even lead to decrease in the lifetime of the network.

3) *Network layer*: Sniffing attack and intelligently carried out DoS attacks (that allow the traffic to pass through it) and then ultimately slowly increasing the magnitude to block the route(congestion)on increase the magnitude of the attack.

4) *Transport layer*: Denial of service attack at this layer is aimed to make use of the information of the network resources(machines) working, the main aim of the attack is to cause adverse impact leading to halt in the working(congestion) of the entire network. Both online, as well as offline services, are likely to be affected through this attack.

5) *Session and presentation layers*: Till date, any attack mainly targeting these layers have not been discovered.

6) *Application layer*–This layer disclosed to both active as well as passive attacks. Distributed denial of service is common at this layer.

Table I presents various types of attacks at the multiple layers. The consequences of these attacks depend on the impact caused by the outcomes of the resources affected by these attacks.

Table II presents various protocols, one of the strategy is to monitor the malicious characteristic of nodes on the basis of Non-Cooperative nodes. So in order to mitigate the attacks on the wireless network, we have taken into account the malicious activity shown in terms of the non Cooperating behavior characteristics exhibited by the nodes of the network. Algorithm 2 describes assigning of Reputation value to each sensor node of the network. One method to identify non-cooperative nodes is to assign a "Node\_Reputation ( )" value to each of the node cooperating in the transmission process. Since each node in Mobile Adhoc networks and WSNs have no other way of collecting the information about the nodes located outside their range, and therefore there is a greater chance of uncertainty in the communication information related to them. So in order to enhance trust and reputation - based system for MANETs and WSNs in particular is a challenging issue. MANETs are assumed to be self-configuring collection of nodes mobile in nature connected by wireless links. These nodes are exhibit random movement and is the primarily the reason for rapidly changing topology of the network.

Load balancing [12] in WSN is critical to classify the sensor node into equal size groups so as to ensure that expected network performance is achieved for each node, Fault tolerance [13] is the feature of a network which ensures reliability and trust aspect of dependency of each sensor node on other nodes of the network.

TABLE I. DDoS ATTACKS AT VARIOUS LAYERS IN WIRELESS SENSOR NETWORK

| Layers                   | Types of Attacks           | Consequences of Attacks            | Impact of Attack |
|--------------------------|----------------------------|------------------------------------|------------------|
| Physical                 | Sniffing                   | Nodes exhibiting malicious pattern | Mild             |
|                          | DDoS                       | Starvation, Depletion of Resources | Severe           |
|                          | Tampering (node capturing) |                                    |                  |
| Data Link                | Active DDos Attacks        | Increase in packet Drop Ratio      | Mild             |
|                          | Passive DDos Attacks       | Decrease in lifetime of network    | Severe           |
| Network                  | Sniffing                   | Nodes exhibiting malicious pattern | Mild             |
|                          | DOS                        | Congestion                         | Severe           |
| Transport                | DoS                        | Congestion                         | Severe           |
| Session and Presentation | No attack noticed yet      | -----                              | -----            |
| Application              | Sniffing                   | Loss of Packests                   | Mild             |
|                          | Spoofing (IP, ARP, DNS)    | Delay in the packets transmitted   | Severe           |
|                          | DDoS                       | Lifetime Decay                     | Severe           |

TABLE II. DEFENSE SCHEME USED AT THE VARIOUS PROTOCOL LAYERS

| Protocols Layers            | Types of Attacks      | Defense Used           |
|-----------------------------|-----------------------|------------------------|
| Physical                    | Node Destruction      | Hide Nodes             |
| MAC (Medium Access Control) | Denial of sleep       | Sleep                  |
| Network                     | Spoofing              | Authentication         |
|                             | Hello Floods          | Geographic Routing     |
|                             | Homing                | Header Encryption      |
| Transport                   | SYN flood             | SYN Cookies            |
|                             | Desynchronization     | Path Authentication    |
| Application                 | Path Based DoS        | Anti-replay protection |
|                             | Reprogramming Attacks |                        |

In these networks, each node plays the dual role of being the end-system as well as the task of relaying the packets to the other nodes. Since the nodes in MANET are autonomous without any common interest, so there is a greater tendency for a node to not participate in a cooperative manner with other nodes of the network. This Non-cooperative behavior exhibited by the node explained in Fig. 1 may lead to malicious activities such as leading to DoS (Denial of service) attacks and various other deviation from ideal expected behavior by the sensor node.

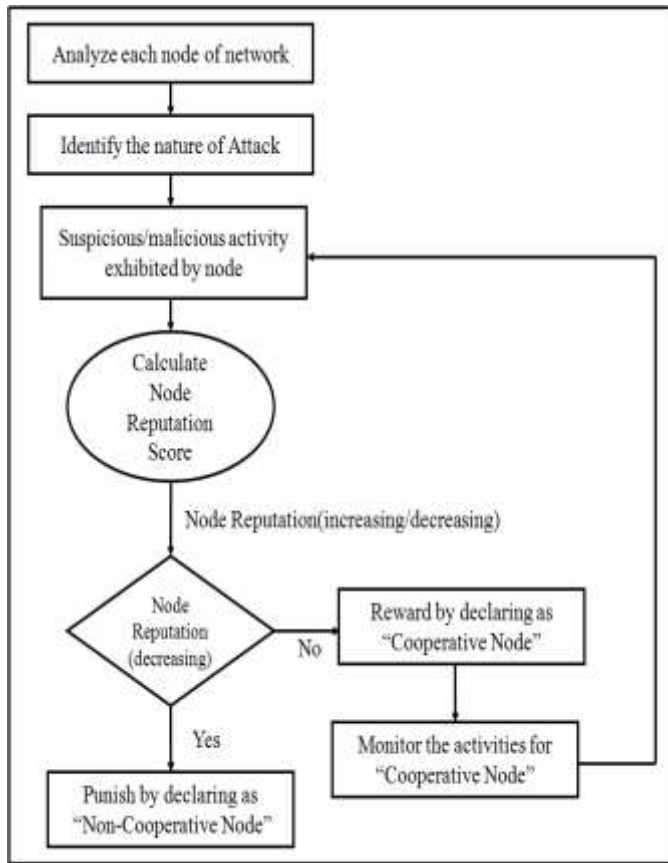


Fig. 1. Identification of Node Behavior in Wireless Sensor Network.

Unlike MANET, in case of WSNs all the sensors nodes are confined to be the part of a single group and they operate to attain same goal. So there arise a need to classify the nodes of WSNs into groups called "Clusters".

### III. SECURITY CONCERNS IN WSNs

WSNs have been emerging as the most widely deployed networks in various application areas. The Security concern for WSNs are as follows [14][15]:

#### A. Confidentiality

It is the measure which assures that sensor nodes control or influence what information related to them may be collected and stored and by whom (sensor nodes) and to whom that piece of data or information may be disclosed.

#### B. Integrity

It is the measure which ensures that the data or information received by a sensor node must not be altered maliciously by the member nodes of WSNs.

#### C. Authentication

It is the measure which ensures that the entity (sensor node) being a genuine member of the network which can be trusted and verified against the data sent and received being the legitimate sender and receiver of the data or information.

#### D. Authorization

The authorization is used to ensure and assure that only the authorized (legitimate) sensor nodes are allowed to perform the required operations in WSNs.

#### E. Availability

It is the measure which ensures that information access is on a timely basis and reliably that is WSN services must be available whenever the WSN users need them.

#### F. Secrecy (Forward and Backward)

Forward secrecy is deployed in WSN in order to disallow a sensor node that has left a Wireless Sensor Network from accessing (read) any future data, whereas Backward secrecy means preventing a new incoming sensor node to a Sensor Network from reading any previous data.

### IV. ATTACK ON WIRELESS SENSOR NETWORKS

There are wide ranges of attacks. Fig. 2 presents security attacks that are classified as passive attacks and active attacks.

#### A. Passive Attacks

In Passive attacks, the main goal of the intruder (opponent) is to monitor (examine) and obtain the information that is being transmitted between the sender and the receiver. These are the attacks against the privacy of wireless sensor network. Some of the passive attacks are the release of message contents, eavesdropping and traffic analysis, etc. In addition to these various attacks aimed to obtain various critical information such as decoding the poorly enciphered traffic, and observing the important information such as secret message and identification. The consequence of these attacks is the exposure of information or any feasible source of data to an attacker.

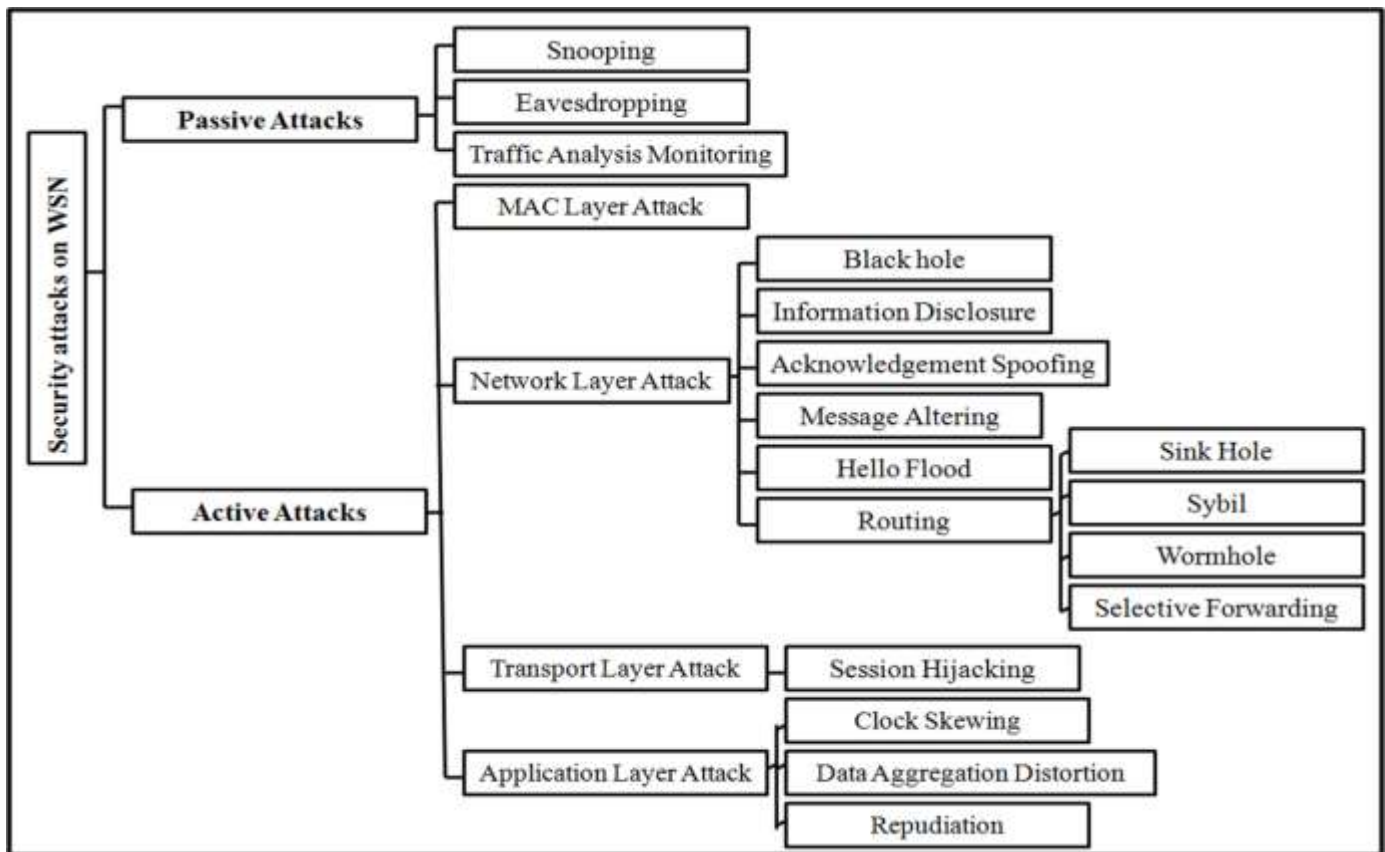


Fig. 2. Classification of Security Attacks for Wireless Sensor Networks.

### B. Active Attacks

Active attacks are mainly targeted with primary aim leading to the modification of the data stream or the creation of a false stream in order to disturb their operations. The attacker alters the data stream to masquerade one entity as some other. As a result of Active attack may be the exposure of data files and their amendment or in worst scenarios may even lead to denial of service (DoS). Various detection and prevention methods can be used to avoid multiple DoS attacks. A DDoS attack is initiated by flooding a massive number of data packets or bogus requests to a victim's network that leads to increase the bandwidth requirements. Therefore, exceeding beyond the capacity of handling the application by the victim(server), so that the processing node is flooded with undue requests that prevent legitimate users from receiving the service and hence leads to congestion or starvation [16], [17].

Such types of attacks target both the service provider and user in addition the main resources for attack can be aimed to disrupt the processing unit or the memory, to drain the energy of the sensor node(battery power), and the bandwidth of the

wireless network. These attacks also affect connectivity and reduce the throughput and quality of service (QoS). As the sensor nodes of WSNs continuously monitor the dynamically changing parameters in the network. Therefore any issue such as packet drop etc. is shared with the neighbouring nodes, based on such information possibility of a kind of denial of service attack can be identified and due to preventive measures can be adopted by the nodes of WSNs.

### V. PROPOSED WORK

Many studies over recent years have summarized the role of various parameters related to methods used for selection of clusters need of re-clustering and study of the QoS parameters such as performance of nodes in WSN. Some approaches uses the concept that the cluster head is selected on the basis of defining a threshold value which can be calculated based on critical parameters such as residual energy and relative position of the node in the network. In this study CH is elected (given in Algorithm 1) on the basis of selecting the maximum value of Reputation and Score of an individual node and Cluster head, respectively.

**Algorithm 1: SELECTION OF CLUSTER HEAD IN WIRELESS SENSOR NETWORK**

```

Score-Scr
Reputation-Rep
Node-Nod
Sm- variable assigning maximum value of scr
Rm- variable assigning maximum value of rep
Value-Val
ClusterHeadSelection ( )
{
Scr = NodScrVal ( )
Rep=NodRep ( )
Sm=Max(Scr ( ) )
Rm=Max(Rep ( ) )
For each cluster Cj ∈ WSN
{
For (each node Ni && Ni = 1 to Ni = n)
{
If ( Sm || Rm )
{
Set Nod Ni as (CH)j
// ith nod is assigned to the jth cluster as CH
}
else
{
Add node jth Cluster i.e. Cj
}
}
}
}

```

Identifying the nodes of WSN into clusters can lead to attain the objectives, namely, as:

**A. Load Balancing**

Load balancing of nodes in a WSN is a measure of the distribution of the various overheads related to data processing or various other intra-cluster management task confined to the cluster head (CH) node within the network. So there arise a need to maintain a balance of the load among the nodes of a WSN, so that every node can meet its expected performance goals. Specifically for WSN where the CH are to be selected among the sensor nodes of the network. Therefore formulation of "Cluster" is crucial in order to extend the lifetime of the network and meet the expected performance criteria.

**B. Fault-Tolerane**

As WSNs are expected to operational in extreme and adverse working situations such as military applications such as battlefield surveillance, border surveillance, Disaster management, security surveillance etc. therefore these networks are likely to suffer from physical damage and malfunction etc. Failure of a node of WSN can have a significant impact on the network and this situation can worsen

if the affected node is a cluster head, as the loss or failure of a CH means loss of certain critical sensor data. So we need an intuitive way to overcome the overcome the failure of a Cluster Head.

**Algorithm 2. REPUTATION-BASED CREDIT SCORE OF A NODE OF A CLUSTER IN A WIRELESS SENSOR NETWORK**

```

Nod_Rep ( )
{
For each nod Ni ∈ WSN
Each node maintains ( Nod_id, Net_Rep_Scr )
//Where Net_Rep_Scr = S_credit-U_credit
If ( S_Credit >= U_Credit )
{
Net_Rep_Scr val Increases
}
Else
Net_Rep_Scr val Decreases
}
where
Nod_id – Nod id in the cluster
S_credit – score corresponding to correctly Packet forwarding capability of a node
U_credit – is score corresponding to Un- correct forwarding capability of a node

S_credit(A,B) = +ive value // Successful_Credit_Score
U_credit(A,B) = -ive value // Unsuccessful_Credit_Score
//Net_Credit_Scr which represents its reputation and is calculate as
Net_Credit_Scr = ∑S_Credit(i , j) + ∑U_Credit(i , j)
// where nod "i" is request for service from "j"

If ( Net_Credit_Scr >= 0 )
{
Nod has + rep.
}
else
{
Nod is consider malicious nod
}
rtn(Net_Credit_Scr);

```

In this study, we have considered incentive (based on *Reputation\_Score()* of a node of the network) based approach can be used to enhance trust for nodes in WSNs to behave in a cooperative manner. Many Reputation and trust-based systems based system has been successfully modelled for WSNs. WSN is an autonomous collection of mobile nodes driven by constrained resources such as energy content of node so in order to enhance the network lifetime is a major concern. In order to address issues such as scalability, energy of a node, the nodes are often grouped into disjoint clusters. Each cluster is monitored by a node referred as cluster head (CH).The selection of Cluster Head is based on calculation



"Node\_Reputation( )" which is the characteristics of individual sensor node of WSN where as *NodeScoreValue( )* which is the characteristic of node behavior in the cluster.

## VI. CONCLUSION

In this work, we have proposed cluster-based mitigation technique basis of *Net\_Credit\_Score* assigned to the nodes of a wireless sensor network. A positive "*Net\_Credit\_Score*" increases the trust & Reputation of a sensor node among the nodes of a WSN, whereas a negative value is an indicator of nodes exhibiting malicious or suspicious behavior. Some of the critical factors which can be considered for future work can be the constrained energy of sensor node of WSN, as the power of each sensor node is limited the network lifespan of WSN is critical issue to consider. Similarly, Cluster count (i.e. size of cluster), Stability of Cluster and Cluster Head and Intra-Cluster topology can also be some critical parameters to consider in devising strategies for mitigating against the security attacks carried out on WSNs.

## ACKNOWLEDGMENT

The authors would like to express their gratitude to King Khalid University, Saudi Arabia for providing administrative support.

## REFERENCES

- [1] I. F. Akyldiz, W. Su, Y. Sankarasubramaniam, E. Cayirci, "Wireless sensor networks: a survey," *Computer Networks*, vol. 38, 2002, pp:393-422.
- [2] Khan, Rizwan, and A. K. Vatsa. "Detection and control of DDOS attacks over reputation and score based MANET." *Journal of Emerging Trends in Computing and Information Sciences* 2.11 (2011): 646-655.
- [3] Upavi .E.Vijay1, Nikhil Sameul2, "Study of Various Kinds of Attacks and Prevention Measures in WSN", *International Journal of Advanced Research Trends in Engineering and Technology (IJARTET)*, Vol. II, Special Issue X, March 2015.
- [4] Sonali Swetapadma Sahu et.a. "Distributed Denial of Service Attacks: A Review", *IJ Modern Education and Computer Science*, 2014, 1, 65-71 Published Online January 2014 in MECS.
- [5] Bhaskar P. Deosarkar1, Narendra Singh Yada and R.P. Yadav, 2008. Clusterhead Selection in Clustering Algorithms for Wireless Sensor Networks: A Survey, *Proceedings of the 2008 International Conference on Computing, Communication and Networking (ICCCN 2008)*, 2008 IEEE.
- [6] Ramesh, K. and Dr. K. Somasundaram, 2011. A Comparative Study of Clusterhead Selection Algorithms in Wireless Sensor Networks, *International Journal of Computer Science & Engineering Survey (IJCSSES)* 2(4).
- [7] Dechene, D.J., A. El Jardali, M. Luccini and A. Sauer 2010. A Survey of Clustering Algorithms for Wireless Sensor Networks, *Network-Based Information Systems (NBIS)*, 2010 13th International Conference on Networking & Broadcasting, pp: 14-16.
- [8] Hyung Su Lee, Kyung Tae Kim and Hee Yong Youn, XXXX. A New Cluster Head Selection Scheme for Long Lifetime of Wireless Sensor Networks, *School of Information and Communications Engineering, Sungkyunkwan University, Korea*.
- [9] Saaty, T.L., 2000. *Fundamentals of Decision Making and Priority Theory with the Analytic Hierarchy Process*, RWS Publications, USA.
- [10] A. Srinivasan, J. Teitelbaum, and J.Wu, DRBTS: Distributed reputation-based Beacon trust system, in *Proceedings of the 2nd IEEE International Symposium on Dependable, Autonomic and Secure Computing*
- [11] (DASC'06), Indianapolis, USA, pp. 277–283, 2006.
- [12] S. Ganeriwal and M. Srivastava, Reputation-based framework for high integrity sensor networks, in *Proceedings of the 2nd ACM Workshop on Security of Ad Hoc and Sensor Networks (SASN '04)*, New York, USA, October 2004, pp. 66–77.
- [13] Alam, M., & Varshney, A. K. (2016). A new approach of dynamic load balancing scheduling algorithm for homogeneous multiprocessor system. *International Journal of Applied Evolutionary Computation (IJAE)*, 7(2), 61-75.
- [14] Gholamreza Kakamanshadi, Savita Gupta, and Sukhwinder Singh. 2015. A survey on fault tolerance techniques in Wireless Sensor Networks. In *Proceedings of the 2015 International Conference on Green Computing and Internet of Things (ICGCIoT) (ICGCIOT '15)*. IEEE Computer Society, USA, 168–173.
- [15] Reddy Y.B. A game theory approach to detect malicious nodes in wireless sensor networks; *Proceedings of the Third International Conference on Sensor Technologies and Applications (SENSORCOMM)*; Athens, Greece. 18–23 June 2009; pp. 462–468.
- [16] Agah A., Asadi M., Das S.K. Prevention of DoS Attack in Sensor Networks using Repeated Game Theory; *Proceedings of the ICWN*; Las Vegas, NV, USA. 26–29 June 2006; pp. 29–36.
- [17] K. Giotis, M. Apostolaki, and V. Maglaris, "A reputation-based collaborative schema for the mitigation of distributed attacks in SDN domains," *Proc. NOMS 2016 - 2016 IEEE/IFIP Netw. Oper. Manag. Symp.*, no. Noms, pp. 495–501, 2016.
- [18] R. Wang, Z. Jia, and L. Ju, "An entropy-based distributed DDoS detection mechanism in software-defined networking," *Proc. - 14th IEEE Int. Conf. Trust. Secur. Priv. Comput. Commun. Trust.* 2015, vol. 1, pp. 310–317, 2015.

# A Predictive Model for the Determination of Academic Performance in Private Higher Education Institutions

Francis Makombe<sup>1</sup>, Manoj Lall<sup>2</sup>  
Department of Computer Science  
Tshwane University of Technology  
Pretoria, South Africa

**Abstract**—The growth and development of predictive models in the current world has influenced considerable changes. Today, predictive modelling of academic performance has transformed more than a few institutions by improving their students' academic performance. This paper presents a computational predictive model using artificial neural networks to predict whether a student will pass or fail. The model is unique in the current literature as it is specifically designed to evaluate the effectiveness of the predictive strategies on neural networks as well as on five additional algorithms. The analysis of the experimental results shows that Artificial Neural Networks outperformed the eXtremeGBoost, Linear Regression, Support Vector Machine, Naive Bayes, and Random Forest algorithms for academic performance prediction.

**Keywords**—Classification modelling; data mining; higher education institutions; accuracy; academic performance

## I. INTRODUCTION

Public higher education providers are institutions that have been established and funded by the state through the Department of Higher Education and Training (DHET). Public providers include universities, universities of technology, and comprehensive universities. Private providers are owned by private organizations or individuals. Higher education institutions (HEIs) operate in an increasingly complex and challenging environment. Competition has increased, and previously anticipated government funding has become scarce [1]. In such circumstances, HEIs must succeed in a financial sense or else they will go out of business [2]. In their quest for survival, common practices adopted by HEIs are to increase the intake of students and try to improve on their success rates. Since, many government and private funds depends on the throughput rates of institutions, being able to predict the chances of any new student's success is very important. This study aims to improve the pass rates of students' in a particular private academic institution by providing a classification model to assist in identifying student at risk of failing a program. Being able to identify such students, the educational institutes can provide a targeted support mechanisms to the needy students. The author in [3] mention that the reasons for the identification of a student at risk of dropouts or attrition early enough are to be able to provide necessary support and interventions for the student with the

goal of reducing dropouts, increasing retention, performance and graduation rate.

Application of the appropriate data mining technique that suits the current scenario is important in order to identify useful patterns. In this article, factors that have an impact on the pass rates of students are identified and used in the classification model. The following algorithms are applied in the construction of the classification model-Artificial Neural Networks, Logic Regression, eXtremeGBoost, SVM, Naive Bayes, and Random Forest algorithms.

The rest of this article is structured as follows: the literature review is presented in Section II. The description of the data and the methodology used are presented in Sections III and IV. The results and its discussion are presented in Section V. In Section VI, conclusions and recommendations are presented.

## II. LITERATURE REVIEW

In a research conducted by [4], the researchers attempted to explore the applicability of Fuzzy C-Means clustering technique for academic performance of students. They found that fuzzy C-Means clustering algorithm serves as a good benchmark to monitor the progression of students modelling in educational domain. The author in [5] also recommended a fuzzy logic-based expert system that periodically evaluates student performance and supplies students with feedback on progress within data grid environment. The system made use of the fuzzy logic theory and develop the decision making process based on fuzzy rules to assess whether a student gets very poor, poor, good, average or excellent performance.

In an attempt to identify the main attributes that may affect the performance of students in engineering, [6] applied data mining concepts such as k-Means clustering and Decision tree Techniques. They used records of 1500 students enrolled for various subjects in engineering. The author in [7] investigated the impact of classroom attendance and gender on academic performance of university students in an Organic Chemistry course. Data was collected through survey involving real time documentation of attendance for each student at each class lesson over a three month period. Their findings show that attendance had a significant impact on the performance. In another study, [8] analysed the impact of class attendance, practical work and assignments in a course on the success rate.

They found that the number of given assignment has a negative impact on the academic performance. They used C4.5 as the classification algorithm for their work. Several other studies conducted have shown that class attendance is an important predictor of academic outcomes which conclude that students who attend more classes generally earn higher final grades [9].

In a study by [10], one of the factors that influences a student's ability to succeed is the socioeconomic conditions. This fact is supported by [11] who state that Student poverty and the lack of sufficient funding have consistently been cited as key reasons for student academic failure and progression difficulties. In the study by [12], they used marks of four academic batches of Computer Science & Information Technology (CS&IT) students for predicting performance. In their study, they collected records of 347 undergraduate students have been mined with classifiers such as Decision tree, Neural Networks and Naive Bayes.

In another study, [13] applied Naïve Bayes for the classification of student evaluation. Their dataset consisted of the following parameter-age, place of birth, gender, high school status (public or private), department in high school, organization activeness, age at the start of high school level, and progress GPA score.

Discriminate analysis was done by [14] to predict the success and failure of students in a specific physics course. Discriminate analysis is a similar technique to multiple regression except that it is used for categorized data. They used this technique to provide a function that contains the variables that should be used for predicting the success of a student. They collected the data for 1622 students who enrolled in Electricity & Magnetism course, which had a high rate of failure. At first they identified many possible predictors such as, SAT grade, MATH GPA, Overall GPA. In another study [15], applied predictive modelling techniques to identify students at risk of dropping out of their registered qualification. They used Support Vector Machine, Naïve Bayes, Decision tree, K-nearest neighbors and Random Forest on 1156 students.

### III. DATA DESCRIPTION

This research followed a quantitative approach. Questionnaires were administered to private academic institutions in an anonymously manner to enhance the privacy and anonymity of the participants. The questionnaires in this study were distributed in two ways: manually handed out and also using the online survey tool survey monkey. The dataset consisted of the following attributes:

- Study hours per week.
- Bursary - whether a student has a bursary or not.
- Class Attendance.
- Student workload (number of modules registered).
- Fulltime study or attending through part-time classes.
- English language proficiency marks.

- Number of employed parents or guardians.
- Group Assignment marks.
- Test marks.
- Individual Assignment marks.

The scatterplot (Fig. 1) shows the distribution of individual test marks in relation to the individual assignment marks. In analysis of this scatterplot, most of the students perform well in both tests and individual assignments. There are a few outliers who perform very well in individual assignments but poorly in tests. According to this scatterplot, the approximate range for tests with most students' marks is 40 to 80, and that for the individual assignments is 50 to 90. This shows that students are generally performing better in individual assignments than in tests.

The scatterplot (Fig. 2) for Test and Group assignment marks shows that a greater proportion of students perform very well in group assignments, where they take part in research activities. By comparison, a lot of students fail the tests as shown by the large concentration of test marks below the mark of 50, compared to the test mark greater than 50. This could provide a basis for intervention by the private institution in efforts to assist the students prepare better for tests.

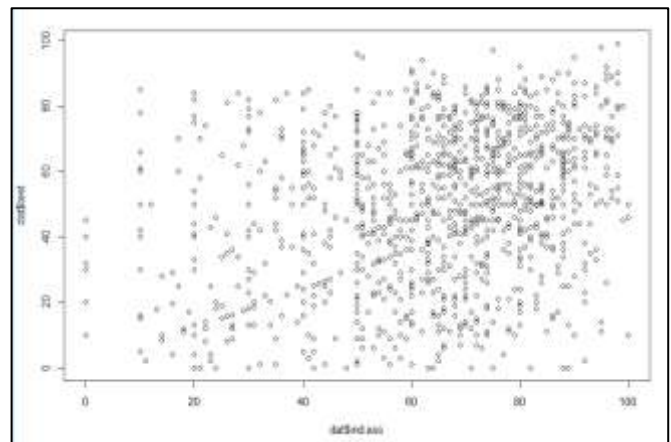


Fig. 1. Scatterplot of Test and Individual Assignment marks.

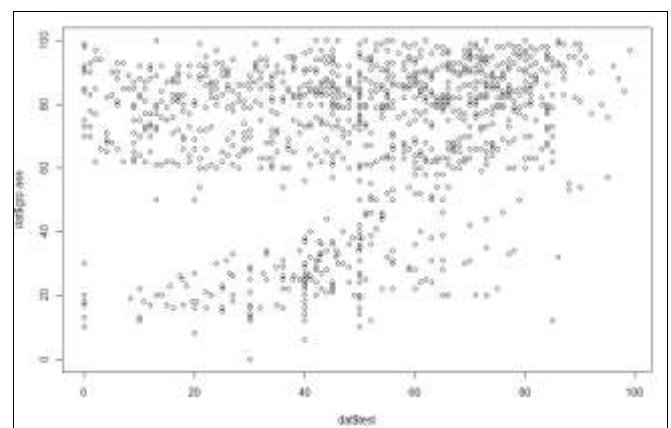


Fig. 2. Scatterplot of Test and Group Assignment Marks.

#### IV. METHODOLOGY

In order to assess the effect of data quality, data attribute significance and class number in the academic performance prediction in this study, six classification algorithms have been selected and implemented in R programming language. These algorithms were chosen because they cover the different approaches used by classifiers for learning and they are state of the art algorithms that are often used in data mining applications [16].

##### A. Random Forests

Instead of building a single tree for classification, the Random Forests constructs a set of trees, and uses them all to classify or to predict. Random Forests were developed by [17], and they create a (forest) collection of decision trees by the method of bagging. Random Forests are sets of learning models where the unknown input is listed according to the majority vote of decision-making bodies. This means that the class predicted by most of the trees would be the last class in the set. Random Forests, increase the classification performance, avoiding overfitting and are robust to outliers and noise [17].

##### B. Neural Network- Multilayer Perceptron (MLP)

This refers to an artificial neural feed network class in which at least three layers of nodes are present: one input layer, one hidden layer and one output layer. Every node is a non-linear activation neuron except for the input nodes. MLP uses a supervised learning method called training backpropagation [18]. Every node layer is fully connected to the next layer, which generates a finite acyclic graph (DAG). Except the input node, each node is a processing node that is used to calculate the output based on an input using a non-linear activation function. Each link of two nodes has a change in weight depending on the training data set. The weight adjustments are based on the error of the measured output difference and the predicted output. The weights are adjusted to reduce the error by using a gradient descent.

##### C. Support Vector Machines (SVMs)

SVMs for binary classification were developed by [19]. This is an approach that is used to solve classification problems using linear methods for both datasets having linearly and nonlinearly separable classes [20].

##### D. Linear Regression

Linear regression helps to predict the value of the Y outcomes variable on the basis of one or more X variables (Equation 1). The objective is to create a linear relationship (a mathematical formula) between a predictor variable(s) as well as the response variable, such that the value of the Y answer is determined by using this formula only when the values (Xs) of the predictors are known. In general, the formula for linear regression is provided as follows:

$$Y = \beta_1 + \beta_2 X + \epsilon \quad (1)$$

where,  $\beta_1$  is the intercept and  $\beta_2$  is the slope. These are called regression coefficients, and  $\epsilon$  is the error term, which refers to the area of Y, that the regression model cannot be able to explain.

##### E. Naïve Bayes Classifiers

These refer to a collection of "probabilistic classifiers" which are based on the application of Bayes' theorem with strict (Naïve) independence assumptions amongst features. Naïve Bayes classifiers are very scalable.

##### F. eXtreme Gradient Boosting

eXtreme Gradient Boosting (XGBoost) is a versatile and enhanced gradient algorithm booster variant designed for efficiency, machine speed and performance of the model. It is an ensemble learning technique that combines multiple machine learning algorithms to lessen errors and increase prediction accuracies.

#### V. RESULTS AND DISCUSSION

The following chart demonstrates the different accuracy, sensitivity, and F-measure values obtained (Fig. 3). Inaccuracies are also shown for each of the six algorithms used in this research. Fig. 1 shows that neural networks algorithm had the best accuracies which also had the least inaccuracies. It also had high precision, and F-measure values where a good classifier has an F-measure value of close to 1, whilst the worst classifier has an F-measure close to 0.

##### A. Receiver Operating Characteristic (ROC) Curve

The purpose of the Receiver Operating Characteristic (ROC) curve is to primarily assess the accuracy of a continuous measurement that is performing a binary outcome prediction. The best classifier has an area under the curve (AUC) value close to 1. Fig. 4 shown below are the AUC values for three classifiers, two with the best performance and one with a poor performance.

The following values were obtained for the AUC. This was done for three classifiers, which were the two best classifiers, and the worst classifier Table I.

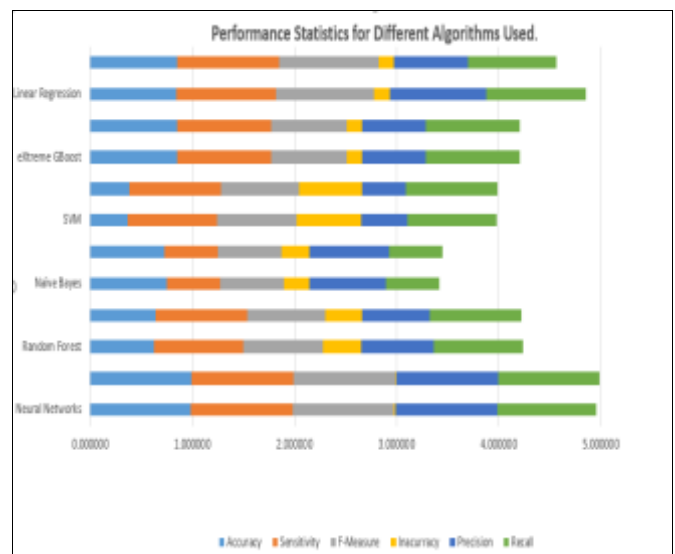


Fig. 3. Performance Statistics for different Algorithms used.

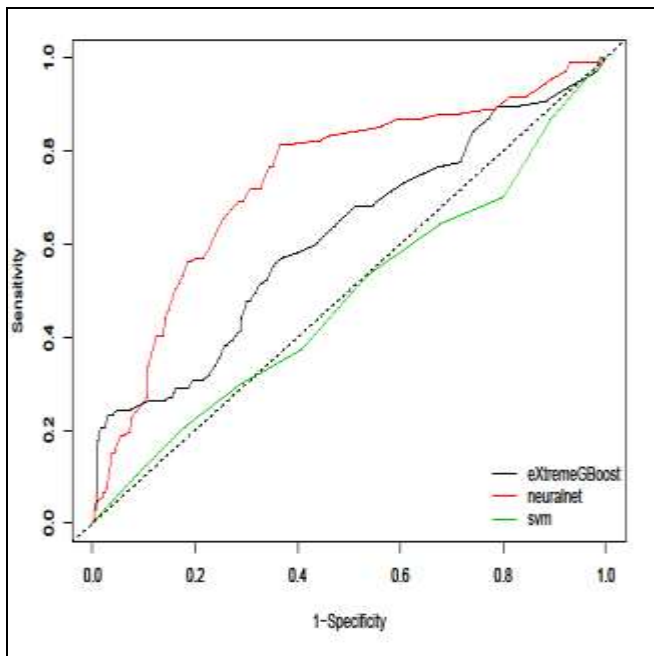


Fig. 4. ROC for Two Most Accurate and the Least Accurate Classifier.

TABLE I. AUC VALUES FOR TWO MOST ACCURATE CLASSIFIERS AND THE LEAST ACCURATE CLASSIFIER

| Algorithm       | AUC  |
|-----------------|------|
| eXtremeGBoost   | 0.62 |
| Neural networks | 0.86 |
| SVM             | 0.43 |

By making use of the AUC and accuracy values obtained in this experiment the neural networks algorithm was selected to be the most suitable algorithm for the prediction of academic performance for this study. The performance of ANN was followed by eXtremeGBoost and then SVM. It can therefore be concluded that the neural net algorithm outperformed the other five algorithms for academic performance classification.

**B. Confusion Matrix Results**

Table II below summarizes the experimental results obtained for both the training and testing dataset, and it also demonstrates the accuracies and misclassification errors obtained using a neural network defined with the simple learning rate algorithm.

TABLE II. NEURAL NETWORK ALGORITHM WITH SIMPLE LEARNING RATE CONFUSION MATRIX RESULTS

| CONFUSION MATRIX |               |                |                |               |                         |
|------------------|---------------|----------------|----------------|---------------|-------------------------|
| Dataset          | True Positive | False Positive | False Negative | True Negative | Misclassification Error |
| Training data    | 5478          | 88             | 86             | 3636          | 0.019                   |
| Test data        | 1223          | 17             | 23             | 1059          | 0.017                   |

Fig. 5 shows the predictions of a sample of six students using the neural networks. These are the computed values which show the predicted value of whether a student will pass or fail a module. The simple learning rate algorithm was used for these predictions. The value of 0.4566725 for the first student in the dataset means that the student is more likely to fail this module. Similarly, the value of 0.6010540 (which is greater than 0.5) for the second student would mean that this student is more likely to pass this module.

```
> head(output$net.result)
      [,1]
[1,] 0.4566725
[2,] 0.6010540
[3,] 0.5961648
[4,] 0.4566725
[5,] 0.4566725
[6,] 0.5130902
> |
```

Fig. 5. Output of Neural Networks.

**VI. CONCLUSIONS AND RECOMMENDATIONS**

In the study the researcher shows the degree of accuracy of the six algorithms used in the study, and their related misclassification errors. It was observed that ANN performed better than Logic Regression, eXtremeGBoost, SVM, Naive Bayes, and Random Forest algorithms. It was observed that bursary and group assignments had a positive correlation with the pass rate. The recommendations, based on the results obtained, are: (1) Group assignments have a positive correlation concerning whether a student will pass or fail as they have a direct effect. Hence it is recommended that students should be encouraged to take a more active role in group assignments. (2) Bursaries have a positive correlation with academic performance; therefore, it is recommended for the private institute to provide bursaries to successful applicants. (3) There should be provision made for booster or support classes meant for students predicted to fail. To have a more accurate assessment of a student's academic performance, data from other domains of higher education value chain such as psychosocial domain, cognitive domain, institutional domain, personality domain, and demographic domain should be considered as future work.

**REFERENCES**

- [1] E. J. Dumond and T. W. Johnson, "Managing university business educational quality: ISO or AACSB?," *Quality Assurance in Education*, 2013.
- [2] H. J. Juhl and M. Christensen, "Quality management in a Danish business school—A head of department perspective," *Total Quality Management*, vol. 19, no. 7-8, pp. 719-732, 2008.
- [3] O. W. Adejo and T. Connolly, "Predicting student academic performance using multi-model heterogeneous ensemble approach," *Journal of Applied Research in Higher Education*, 2018.
- [4] R. S. Yadav and V. P. Singh, "Modeling academic performance evaluation using fuzzy c-means clustering techniques," *International Journal of Computer Applications*, vol. 60, no. 8, 2012.
- [5] S. Patel, P. Sajja, and A. Patel, "Fuzzy logic based expert system for students performance evaluation in data grid environment," *International Journal of Scientific & Engineering Research*, vol. 5, no. 1, 2014.

- [6] V. Sreenivasarao and C. G. Yohannes, "Improving academic performance of students of defence university based on data warehousing and data mining," *Global Journal of computer science and technology*, 2012.
- [7] O. D. Ayodele, "Class attendance and academic performance of second year university students in an organic chemistry course," *African Journal of Chemical Education*, vol. 7, no. 1, pp. 63-75, 2017.
- [8] N. A. Yassein, R. G. M. Helali, and S. B. Mohomad, "Predicting student academic performance in KSA using data mining techniques," *Journal of Information Technology and Software Engineering*, vol. 7, no. 5, pp. 1-5, 2017.
- [9] A. Kirby and B. McElroy, "The effect of attendance on grade for first year economics students in University College Cork," *Vol. XX, No. XX, Issue, Year, 2003*.
- [10] D. E. Roby, "Research on school attendance and student achievement: A study of Ohio schools," *Educational Research Quarterly*, vol. 28, no. 1, pp. 3-16, 2004.
- [11] S. Mngomezulu, R. Dhunpath, and N. Munro, "Does financial assistance undermine academic success? Experiences of at risk students in a South African university," *Journal of Education (University of KwaZulu-Natal)*, no. 68, pp. 131-148, 2017.
- [12] R. Asif, A. Merceron, and M. K. Pathan, "Predicting student academic performance at degree level: a case study," *International Journal of Intelligent Systems and Applications*, vol. 7, no. 1, p. 49, 2014.
- [13] N. Dengen, E. Budiman, M. Wati, and U. Hairah, "Student Academic Evaluation using Naïve Bayes Classifier Algorithm," in *2018 2nd East Indonesia Conference on Computer and Information Technology (EIConCIT)*, 2018: IEEE, pp. 104-107.
- [14] E. W. Thomas, M. J. Marr, A. Thomas, R. M. Hume, and N. Walker, "Using discriminant analysis to identify students at risk," in *Technology-Based Re-Engineering Engineering Education Proceedings of Frontiers in Education FIE'96 26th Annual Conference*, 1996, vol. 1: IEEE, pp. 185-188.
- [15] R. Lottering, R. Hans, and M. Lall, "A model for the identification of students at risk of dropout at a university of technology," in *2020 International Conference on Artificial Intelligence, Big Data, Computing and Data Communication Systems (icABCD)*, 2020: IEEE, pp. 1-8.
- [16] I. H. Witten and E. Frank, "Data mining: practical machine learning tools and techniques with Java implementations," *Acm Sigmod Record*, vol. 31, no. 1, pp. 76-77, 2002.
- [17] L. Breiman, "Random forests," *Machine learning*, vol. 45, no. 1, pp. 5-32, 2001.
- [18] J. Gao, X. He, and L. Deng, "Deep learning for web search and natural language processing," 2015.
- [19] C. Cortes and V. Vapnik, "Support-vector networks," *Machine learning*, vol. 20, no. 3, pp. 273-297, 1995.
- [20] D.-M. Tsai and C.-C. Lin, "Fuzzy C-means based clustering for linearly and nonlinearly separable data," *Pattern recognition*, vol. 44, no. 8, pp. 1750-1760, 2011.

# The Effect of Requirements Quality and Requirements Volatility on the Success of Information Systems Projects

Eman Osama<sup>1</sup>, Mohamed Abdelsalam<sup>3</sup>

Business Information Systems department  
Faculty of Commerce and Business Administration  
Helwan University  
Cairo, Egypt

Ayman Khedr<sup>2</sup>

Information Systems Department  
Faculty of Computers and Information Technology  
Future University in Egypt (FUE)  
Cairo, Egypt

**Abstract**—This study aims to identify the effect of poorly written requirements specifications of software development and its continuous changes; on information systems' projects success and its influence on time and cost overrun of the project based on empirical understanding in practice. As the world is moving towards the internet of things and due to the dramatic increase in demand on complex information systems projects, the development of information systems became more difficult and handling the customers' requirements became very challenging. This research follows a conclusive design, Using a descriptive research design was held first to reveal and discover the characteristics of a good requirement, and then a quantitative method was used through conducting questionnaire and distributing to more than 400 participants in the software industry in Egypt, to understand the relationship between variables and how to improve the quality of data based on real world observations or experiment. The data collected was analyzed using python and R analysis techniques. The results indicates that, the organizations with the highest quality of requirements and less requirement volatility, have higher software success rates in terms of Project's efficiency as well as Business and direct organizational success, while the requirements volume doesn't have significant effect on success rates. From this analysis we developed an initial model.

**Keywords**—Requirement engineering; Software Requirements Specification (SRS); requirements quality; requirements volatility; project success factor

## I. INTRODUCTION

Within the huge and speedy development in computers and information systems, along with the covid-19 pandemic is forcing governments all over the world towards technology to deal with the huge amount of data and to provide reliable information to undertake right decision [1]. Also many governments designed apps to assist in COVID-19 battle [2].

In addition to the increase in number of enterprises that use computers and systems and as the worldwide IT Expenditure in 2019 was almost USD 3.7 trillion [3]. According to the Standish Group Chaos Report, it stated that more than 31% of the projects were failed before they completed, and that only 16.2% failed the project schedule and budget [4]. And since the software products development is very complex and mainly based on customer requirements. Thus the research

main objective is to provide deeper interpretation of requirements specification phase; reduce the requirements vagueness, and identify how they affect the project success, dictates the good requirements specifications measurements then how to enhance requirements quality based on empirical study in the real-world practice within the Egyptian IT industry. This research is one of the few research applied on using empirical study in Egypt industry. As the faster-growth of Egyptian IT investment over the medium term, according To Egyptian ICT, growth rate of ICT GDP has increased from 14.1% in 2017/2018 to 16% in 2018/2019; ranking the sector among the highest growing sectors in the economy [5]. The research objective is not only concerns with defining the impact of requirements' quality and requirements volatility on the success of information systems' projects but also to define an equation that calculates the estimated percentage of project success given the percentage of requirements volatility and quality.

## II. LITERATURE REVIEW

### A. The Antecedents of a Successful Project

As "project" is authoritative well-defined by "BS 6079-2:2000 Project Management Vocabulary" as number of activities to be completed within a specific interval of time and budget to build certain product originally intended human resources [6-7]. According this definition, the success of projects can be stated that it depends on three main factors: Time, Cost and functionalities that meets the stakeholders' needs. The ultimate focus in this research is whether the quality of requirements, requirements volume and the rate of requirements volatility affects the success of software development (SD) projects or not. There are different definitions and criteria for measuring the project success [8-10]. Ramos & Mota, 2016, stated that there is a different point of view on project performance for every project stakeholder. E.g., the project manager and the team members of a project may consider it successful while it may not be successful for the CEO [11]. Hence, the success parameters vary depends on each stakeholder [12-13].

Common success factors were noticed that can be measured to identify whether the project is successful or not. It is categorized into two different categories; the first

category is the project efficiency which measures whether “the project has been completed within its budget and on schedule” [14-17], and whether the product met the customers’ requirements and the end users are satisfied by the end product or not [18-20]. The other category is the Business organizational success; which measures if the project increased the profitability of the organization, contributed to the direct performance of the organization [21] and developed better managerial capabilities or not [22] [32].

### B. Requirements Engineering Process

Requirements engineering process also known as requirements analysis; it is one of the main critical stages of any system development process as it is the first phase [26], in which the requirements is gathered analyzed, documented, along with defining the purpose and the scope of the software system are discovered and documented [23]; as well the stakeholders are identified with their demands and desires [24]. Its main objective is to provide a common understanding and vision of the system to be developed between technical and business stakeholders [25].

### C. Phases of Requirements Engineering Process

Requirements engineering is divided into two main parts that aren’t completely separated which are requirement development and requirements management [26]; it consists of Requirements Elicitation, Analysis, specification and validation shown in Fig. 1 [28], and Requirements management is focusing on the minimization of the disruptive impact on the project by accommodating the very real changes [27].

### D. Requirements Quality

As requirements started with the stakeholders’ intentions and it express the needs and expectations from different stakeholders; it is very difficult to measure or quantify how good the requirements are written. Several books describe how to write good requirements specification [29-30]. According to (ISO-IEEE) a good requirement should fulfill these quality criteria: “Correctness, unambiguousness, completeness, consistence, prioritization, verifiability, modifiability and traceability” [31] Then it was replaced by “ISO-IEEE 29148 which introduces feasibility, necessity, free of implementation, and singularity as new characteristics for requirements, while removing prioritization, correctness, and modifiability” [32].

ISO/IEEE 29148:2011 divided the quality characteristics into two categories: first, the characteristics that should be followed by an individual requirement, which are: atomic, complete, consistent, feasible, unambiguous, verifiable, traceable, necessary and should be free of implementation details [32].

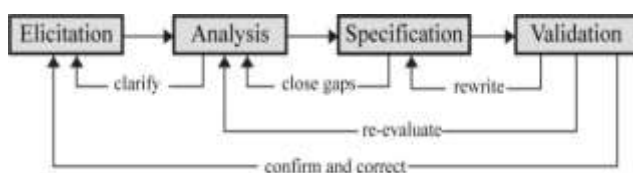


Fig. 1. Phases of Requirements Engineering.

### E. Requirements Volatility

Requirements volatility relates to changes, deletions, addition and adjustments to requirements during the life cycle of the system development. Handling requirements volatility leads to additional work in coding and design, and it has a huge effect on budget, schedule and performance of the project; which may intimidate the project success [33-35]. Although some of these changes are essential for the effective development of the project [36]. But managing requirements volatility is a challenging job [37]. This paper analyzes the volatility of requirements along with the development phases after the requirements phase freeze. According to Mohammad D., (2016) study; maintenance phase has the maximum number of changes requests while design phase has the least number of changes requests [38]. Although agile development is known for its flexibility with requirements volatility but it also has an impact on the continuous change of requirements [39].

### F. Requirements Volatility Factors

There are two main factors that lead to volatility of requirements; the first one is Business environment changes; the changes by external factors affecting the development of project as regulation and laws of government, sources of finances, organizational policies, technological factors, taxes laws and changes in management [40]. And the second is the Development environment changes; the changes during the development of the project, as number of requirements errors, evolving technological needs with the users, turnover of the team members and missing requirements. And there are several other factors as requirements uncertainty caused by lack of detailed information, changing user needs, communication issues and stakeholders dependencies [39]. Requirements Volatility Causes can be categorized to Human errors [41], Process Errors, and Documentation Errors.

### G. Requirements Volume

Measuring the size of software is different from software effort estimation. Software size estimate represents the set of deliverables that should be implemented, it also consider calibration size adjustment and the other software that will be integrated with the new system.

Software size can be driven from several techniques; the main two methods are: 1) expert judgment/analogy, in which an expert gives an estimate for the software based on his previous knowledge and experience [42-43]; 2) functional analysis which is based on requirements specifications and it represents the number of functionalities with several techniques IFPUG [44], NESMA [45] and COSMIC FFP [46].

#### 1) Research Hypotheses:

H1: Quality of requirements has a positive influence on the success of the whole information systems’ development.

H2: Requirements volume affects requirements specifications quality.

H3: There is a negative effect of requirements volatility on the quality of requirements.



### III. RESEARCH METHODS

The research design represents the framework that helps to deeply understand the research problem and plan the strategy, and its main purpose is to control and validate the study through examining the research questions [47]. While this study purposes to identify the effect of requirements quality and requirements volatility on the degree of projects' success, this study used a descriptive method to provide a comprehensive and in-depth interpretation of the characteristics of a good requirement and to identify the success measures of information systems' project within the Egyptian ICT sector. The researcher believes that studying the Egyptian ICT sector can provide new insights and results, given the continuous growth of Egyptian IT investment. This study followed a quantitative design through conducting questionnaire and distributing it in different system development companies around Cairo and Giza, to examine the relation between the research variables. The research population was determined to include all different roles that involved in the requirements of system development in Egyptian Companies. The measurement instruments were derived from earlier validated scales and were revised through piloting by academic experts. The first edition of the questionnaire was pre-tested by seven different experts in software development companies. Some modifications were applied on the pre-test findings in order to improve and clarify the questions.

#### A. Survey Design

The questionnaire inquired about respondents' background information and profile using six questions tackling "age, gender, occupation, years of experience and the number of requirements he worked on". The measures' of validity and reliability data are presented in the results section. Then it was divided into three thematic blocks: i) project success factors, ii) requirements quality criteria, ii) requirements volatility measurement, all related to the respondent's recent project. These variables were encompassed 33 statements under study of "5 points Likert scales whereby "1" represents "strongly disagree" and "5" represents "strongly agree"". The criteria of success of Information systems projects were defined and validated by Serrador & Turner, [2014] [48] and Fernández, D.M., Wagner, S., [2017] [49], based on 4 items Scope, Time, Cost, Quality, Stakeholders Satisfaction that were represented by 12 questions. The quality of requirements was measured with eight factors defined and validated by IEEE: Atomic, Complete, Consistent, Feasible, Unambiguous, Verifiable, Necessary, Implementation free and Traceable [31-32], the questions were validated by Abelein, U (2015) [50].

#### B. Data Collection

The questionnaire were distributed online through google forms, it targeted the IT practitioners from different organizations with different industries, who were, due to their job, able to give an expert judgment on how projects are seen and rated in their organization. This research followed convenience non-probability sampling in which the questionnaire was filled by 400 respondents from the software industry chosen randomly from different organization. The research population was determined to include all different

roles that involved in the requirements of system development. The sample was randomly selected to represents the population of SD companies in Egypt. The target sample was based on the role of the respondent, in which the respondent should be an employee from the IT sector, who was due to his/her job, able to have an expert judgment on their recent projects. It was tried to mainly recruit IT project managers. Requirements engineers, system analysts, developers, quality managers and executives from the IT sector for participation in the survey. In total, 400 usable returns were won. These were evaluated with Python and R.

### IV. RESULTS

#### A. Descriptive Statistics Results

The respondents' demographics are presented in Table I. As shown the respondents to this survey varied widely in relation to their primary position within the organization. However, middle-aged respondents who worked in senior positions on high-value projects were particularly strongly represented, so it was construed that the responses provided were based primarily on respondents' experience. Respondents' roles that stood out in the sample were: Scrum Master/Project Manager, with 110 of the respondents; Requirements Engineer with 100 respondents; Quality Manager/Test Manager and executive managers with 50 respondents each, followed by Programmers/Developers and Business/ system Analysts with 30 respondents each. Out of the total survey participants, most of the respondents have 2 to 5 years of working experience in system development industry (32.5%) followed 6 to 9 years (27.5%) and more than 10 years (25%), this reflects on a better understanding of project life cycle. Only a minority had less than 2 years and more (15%) working experience Information technology industry; which represents that the responses of the questionnaire is based on the responders' experiences, as years of experience play vital role in decision making.

#### B. Factorial Analysis

Factor analysis method were applied towards observed variables to find out subsets of variables, to describe the validity correlated variables and to test the proposed hypotheses based on validated and reliable item. The researcher generated the factor analysis using programming language R which provides a wide array of statistical techniques to capture the right model for your data in order to assign each item to its construct.

Table II presents the demonstration of the Factor Analysis results. As shown the sampling acceptability equals 0.718 that is well-accepted (as it exceeds 0.6). The sphericity test by Bartlett was found to be significant as shown as well (value is less than 0.05). A total of 3 factors measuring the effect of requirements quality on the managerial capabilities were extracted, this is similar as the number of proposed constructs.

#### C. Reliability Test

The reliability test was conducted to ensure that all the variables are internal. Cronbach's alpha was used to calculate the reliability since it is the common reliability measure; as well it is described as one of the most significant and widespread statistics in research involving Likert-type scales

and it is used for dichotomous and continuously scored variables [59]. It was generated using R studio to perform the reliability test and to assign each item on its construct, as there are no Python packages with the required functionality.

According to Bernard (2017), for all variables Cronbach's alpha value must be above 0.6. Preferably it should be greater than 0.7 [57]. Concluded from the reliability assessment shown in Table III that all the items in this study had good internal consistency and were highly reliable and this implied helping us to carry out more statistical analysis needed.

TABLE I. SAMPLE CHARACTERISTICS AND PROFILE

| Gender                       | Frequency | Percentage | Valid percentage |
|------------------------------|-----------|------------|------------------|
| Male                         | 270       | 67.5%      | 67.5%            |
| Female                       | 130       | 32.5%      | 32.5%            |
| Age Group                    | Frequency | Percentage | Valid percentage |
| 20-29                        | 180       | 45%        | 45%              |
| 30-39                        | 190       | 47.5%      | 47.5%            |
| 40-49                        | 20        | 5%         | 5%               |
| More than 50                 | 10        | 2.5%       | 2.5%             |
| Primary Role                 | Frequency | Percentage | Valid percentage |
| Scrum Master/Project Manager | 110       | 27.5%      | 27.5%            |
| Quality Manager/Test Manager | 50        | 12.5%      | 12.5%            |
| Software Tester              | 10        | 2.5%       | 2.5%             |
| Programmer/Developer         | 30        | 7.5%       | 7.5%             |
| Business/ system Analyst     | 30        | 7.5%       | 7.5%             |
| Requirements Engineer        | 100       | 25.0%      | 25.0%            |
| Requirements Process Owners  | 20        | 5.0%       | 5.0%             |
| Executive Manager            | 50        | 12.5%      | 12.5%            |
| Years of experience          | Frequency | Percentage | Valid percentage |
| Less than 2                  | 60        | 15%        | 15%              |
| 2-6 years                    | 130       | 32.5%      | 32.5%            |
| 6-10 years                   | 110       | 27.5%      | 27.5%            |
| More than 10 years           | 100       | 25%        | 25%              |

TABLE II. EFA RESULTS OF STUDY CONSTRUCTS

|  |                    |         |
|--|--------------------|---------|
| “KMO and Bartlett's Test Kaiser-Meyer-Olkin Measure of Sampling Adequacy.” | 0.718              |         |
| “Bartlett's Test of Sphericity”  | Approx. Chi-Square | 2489.22 |
|  | df                 | 66      |
|  | Sig.               | 0.000   |

TABLE III. “CRONBACH'S ALPHA TEST”

| Cronbach's Alpha | Cronbach's Alpha Based on Standardized Items | N of Items | Lower alpha | Upper alpha |
|------------------|--|------------|-------------|-------------|
| 0.801137         | 0.87   | 33         | 0.79        | 0.82        |

#### D. Hypothesis Testing

Using correlation analysis, to identify and examine the relation between different variables. Correlation coefficient results from this analysis that calculates the linear relation between two variables. Correlation coefficient calculation is used to figure out how deep a relationship is between data points.

Table IV summarizes the study findings of the correlations between the variables. The data in the table was generated by the researcher using Python. According to Table IV, there is a strong positive association between requirement quality and success of IS projects ( $r = 0.706356$ ,  $p < 0.01$ ). **Therefore, H1 is accepted.** This correlation indicates that the higher the quality of the requirements, the more successful the projects of the ISs.

Requirements Volume is weakly negative correlated with the success of ISS' projects. This relation can be Negligible ( $r = -0.140016$ ,  $p < 0.01$ ). This indicates there is no relation between requirements volume and the success of ISS' projects. **H2 is rejected.**

There is a moderate negative relationship between requirements volatility and success of ISS' projects ( $r = -0.373626$ ,  $p < 0.01$ ). This indicates that the higher the requirements' volatility, the lower the quality of requirements specifications. **H3 is accepted.**

The entire hypotheses with confidence level 99% which means even if the number of participants changes it won't affect the results.

#### E. Regression Analysis

Regression analysis is a statistical test in order to examine the data collected from questionnaire to define and quantify the impact of variables on each other. Regression analysis demonstrates the amount of change of the dependent variable (Success of IS projects) and the independent variables (Requirements Volatility and quality of requirements). Regression calculates the impact of variables by percentage.

Linear Regression was used that was introduced by Sir Francis Galton, to capture the effect of uncertainty in recognizing the relationship between two variables by Kumari K., et al., 2018 [58]; all the data below generated by the researcher using Python.

TABLE IV. THE CORRELATION BETWEEN PROJECT SUCCESS AND THE PREDICTOR VARIABLES (N=400)

|                     |                         | Success of ISS' Projects |
|---------------------|-------------------------|--------------------------|
| Pearson Correlation | Requirements Quality    | 0.706356                 |
| Sig. (2-tailed)     |                         | .000                     |
| Pearson Correlation | Requirements Volatility | -0.468581                |
| Sig. (2-tailed)     |                         | .000                     |
| Pearson Correlation | Requirements Volume     | -0.140016                |
| Sig. (2-tailed)     |                         | .000                     |

\*\*\*\*Confidence level 99%, significance level of P-value  $\leq 0.01$ , T-value  $\pm 2.58$ \*\*

Two different simple models were made with the data to know the relation between the Success of information systems' projects (dependent variable) and Quality of Requirements (independent variable) and Requirements Volatility shown in the supplementary material.

## V. CONCLUSION

The purpose of this research was to examine and analyze requirements specifications quality and the percentage of requirements volatility impact on the success of development of information systems' projects, using the respondents' empirical experiences. Briefly, this research leads to a deeper understanding of how organizations in real life evaluate the success of projects and how they measure the quality of requirements with regard to software development. These findings were based on the 400 respondents who answered the questionnaire, the respondents were chosen randomly from different Egyptian organizations within the software industry; the sample was based on the respondents' role; to benefit from their real life empirical experience. The findings obtained from this research indicate that the organizations with the finest requirement specifications quality are the organizations with higher success rates in software development, which congruent with the results in [51-54]. Moreover, Organizations with lower percentage of requirements volatility accomplished more software development success, yet it shouldn't be presumed that the use of high quality requirements specifications alone ensures the accomplishment of such success.

Furthermore, the size of requirements affects neither the requirements quality nor the success percentage of the SD project.

Finally, it is found that requirements volatility and have statistically moderate negative relationship; which congruent with the results in [16, 18-20, 55-56]. Based on the analysis results an initial model was developed by Python to predict the probability of project success, given the percentage of requirements quality and Requirements volatility from validation tools.

Future work is divided into two direction, first the extend of data collection to include different Geographic regions to gain further insights and provide a better understanding on the effect of poorly written requirements. Second, investigating how to assure the quality of requirements based on the 8 characteristics of ISO/IEEE that affects the success of information systems projects.

## REFERENCES

- [1] Shiffman, J. (2020). The Role of National Health Information Systems in the Response to COVID-19, from <https://coronavirus.jhu.edu/from-our-experts/the-role-of-national-health-information-systems-in-the-response-to-covid-19>.
- [2] Digital technologies critical in facing COVID-19 pandemic | UN DESA Department of Economic and Social Affairs. (2020, April). Retrieved June 01, 2020, from <https://www.un.org/development/desa/en/news/policy/digital-technologies-critical-in-facing-covid-19-pandemic.html>
- [3] Gartner, Inc., [Gartner Says Global IT Spending to Grow 3.7% in 2020], [October, 2019].
- [4] Standish Group (2014)–“The Chaos Report on Software Projects 2014, In: Project Smart. The Standish Group, USA. Available: < <http://www.projectsmart.co.uk/white-papers/chaos-report.pdf> >.
- [5] Egypt Ministry of Communication and Information Technology. (2019). ICT Sector Performance 2018-2019.
- [6] BS 6079-2:2000 (2015), BSI Standards Publication Project management – Part 2: Vocabulary, BS 6079-2:2000, March 2000, available at: <http://shop.bsigroup.com/en/ProductDetail/?pid=00000000030029005>.
- [7] PMI. A guide to the project management body of knowledge - PMBoK. 4th ed. PMI Standards Committee, Upper Daryby, 2008.
- [8] Edwita, A., Sensuse, D. I., & Noprisson, H. (2017, October). Critical success factors of information system development projects. In 2017 International Conference on Information Technology Systems and Innovation (ICITSI) (pp. 285-290). IEEE. doi:10.1109/icitsi.2017.8267958.
- [9] Goncalves, A., Oliveira, P. M., & Varajão, J. (2018, June). Success factors of information technology and information systems projects—A literature review. In 2018 13th Iberian Conference on Information Systems and Technologies (CISTI) (pp. 1-7). IEEE. doi:10.23919/cisti.2018.8398634.
- [10] Tsoy, M., & Staples, D. S. (2020, January). Exploring Critical Success Factors in Agile Analytics Projects. Proceedings of the 53rd Hawaii International Conference on System Sciences. doi:10.24251/hicss.2020.122.
- [11] Ramos, P. A., & Mota, C. M. (2016). Exploratory Study Regarding How Cultural Perspectives Can Influence the Perception of Project Success in Brazilian Companies. *Production*, 26(1), 105-114. doi:10.1590/0103-6513.173114.
- [12] Bodicha, H. H. (2015). How to Measure the Effect of Project Risk Management Process on the Success of Construction Projects: A Critical Literature Review. *The International Journal of Business & Management*, 3(12), 99-112.
- [13] Beleiu, I., Crisan, E., & Nistor, R. (2015). Main Factors Influencing Project Success. *Interdisciplinary Management Research*, 11, 59-72.
- [14] Bermejo, P. H. de S., Zambalde, A. L., Tonelli, A. O., Souza, S. A., Zuppo, L. A., & Rosa, P. L. (2014). Agile Principles and Achievement of Success in Software Development: A Quantitative Study in Brazilian Organizations. *Procedia Technology*, 16, 718–727. doi:10.1016/j.protcy.2014.10.021.
- [15] Serrador, P., & Turner, R. (2015). The Relationship between Project Success and Project Efficiency. *Project Management Journal*, 46(1), 30–39. doi:10.1002/pmj.21468.
- [16] Sanchez, O. P., Terlizzi, M. A., & de Moraes, H. R. de O. C. (2017). Cost and time project management success factors for information systems development projects. *International Journal of Project Management*, 35(8), 1608–1626. doi:10.1016/j.ijproman.2017.09.007.
- [17] De Wit, A. (1988). Measurement of project success. *International Journal of Project Management*, 6(3), 164–170. doi:10.1016/0263-7863(88)90043-9.
- [18] Agarwal, N., & Rathod, U. (2006). Defining “success” for software projects: An exploratory revelation. *International Journal of Project Management*, 24(4), 358–370. doi:10.1016/j.ijproman.2005.11.009.
- [19] Saleh, M., Baharom, F., Mohamed, S. F. P., & Ahmad, M. (2018). A Systematic Literature Review of Challenges and Critical Success Factors in Agile Requirement Engineering. 242-247. presented at Knowledge Management International Conference (KMICe): Access: <http://www.kmice.cms.net.my/ProcKMICe/KMICe2018/pdf/CR64.pdf>.
- [20] Misra, S. C., Kumar, V., & Kumar, U. (2009). Identifying some important success factors in adopting agile software development practices. *Journal of Systems and Software*, 82(11), 1869–1890. doi:10.1016/j.jss.2009.05.052.
- [21] Badewi, A. (2016). The impact of project management (PM) and benefits management (BM) practices on project success: Towards developing a project benefits governance framework. *International Journal of Project Management*, 34(4), 761–778. doi:10.1016/j.ijproman.2015.05.005.
- [22] Davis, K. (2016). A method to measure success dimensions relating to individual stakeholder groups. *International Journal of Project Management*, 34(3), 480–493. doi:10.1016/j.ijproman.2015.12.009.
- [23] Garcia Alcazar, E., & Monzon, A. (n.d.). A process framework for requirements analysis and specification. Proceedings Fourth

- International Conference on Requirements Engineering. ICRE 2000. (Cat. No.98TB100219). doi:10.1109/icre.2000.855582.
- [24] Pohl, K. (2010). Fundamentals of Requirements Validation. Requirements Engineering, 511–536. doi:10.1007/978-3-642-12578-2\_13.
- [25] Paetsch, F., Eberlein, A., & Maurer, F. (n.d.). Requirements engineering and agile software development. WET ICE 2003. Proceedings. Twelfth IEEE International Workshops on Enabling Technologies: Infrastructure for Collaborative Enterprises, 2003. doi:10.1109/enabl.2003.1231428.
- [26] Dick, J., Hull, E., & Jackson, K. (2017). Requirements Engineering in the Problem Domain. Requirements Engineering, 9-11. doi:10.1007/978-3-319-61073-3\_5.
- [27] Heikkila, V. T., Damian, D., Lassenius, C., & Paasivaara, M. (2015). A Mapping Study on Requirements Engineering in Agile Software Development. 2015 41st Euromicro Conference on Software Engineering and Advanced Applications. doi:10.1109/seaa.2015.70.
- [28] Wiegers, K. E., & Beatty, J. (2013). Software requirements. Redmond, WA: Microsoft Press.
- [29] Davis, A. (2013). Just enough requirements management: where software development meets marketing. Addison-Wesley.
- [30] Corbin, J. (1991). Exploring Requirements: Quality before Design. Donald C. Gause, Gerald M. Weinberg. The Library Quarterly, 61(2), 236–237. doi:10.1086/602347.
- [31] ISO/IEEE 830-1998, “Recommended Practice for Software Requirements Specifications”, IEEE. doi= 10.1109/IEEESTD.1998.88286, 1998.
- [32] ISO/IEC/IEEE International Standard - Systems and software engineering -- Life cycle processes -- Requirements engineering. (n.d.). doi:10.1109/ieeestd.2018.8559686.
- [33] Chari, K., & Agrawal, M. (2017). Impact of incorrect and new requirements on waterfall software project outcomes. Empirical Software Engineering, 23(1), 165–185. doi:10.1007/s10664-017-9506-4.
- [34] Dasanayake, S., Aaramaa, S., Markkula, J., & Oivo, M. (2019). Impact of requirements volatility on software architecture: How do software teams keep up with ever-changing requirements? Journal of Software: Evolution and Process, 31(6). doi:10.1002/smr.2160.
- [35] Alsalemi, A. M., & Yeoh, E.-T. (2017). A Systematic Literature Review of Requirements Volatility Prediction. 2017 International Conference on Current Trends in Computer, Electrical, Electronics and Communication (CTCEEC), 55–64. doi:10.1109/ctceec.2017.8455174.
- [36] AlSanad, A., & Chikh, A. (2015). The Impact of Software Requirement Change – A Review. Advances in Intelligent Systems and Computing, 803–812. doi:10.1007/978-3-319-16486-1\_80.
- [37] Akbar, M. A., Sang, J., Khan, A. A., & Hussain, S. (2019). Investigation of the requirements change management challenges in the domain of global software development. Journal of Software: Evolution and Process, 31(10). doi:10.1002/smr.2207.
- [38] D. Aljohani, M., & J. Qureshi, M. R. (2016). Management of Changes in Software Requirements during Development Phases. International Journal of Education and Management Engineering, 6(6), 12–26. doi:10.5815/ijeme.2016.06.02.
- [39] Aaramaa, S., Dasanayake, S., Oivo, M., Markkula, J., & Saukkonen, S. (2017). Requirements volatility in software architecture design: an exploratory case study. Proceedings of the 2017 International Conference on Software and System Process - ICSSP 2017, 40–49. doi:10.1145/3084100.3084105.
- [40] Khan, F., Benslimane, Y., & Yang, Z. (2019). Managing Information Systems Requirements Volatility in Development Projects: Mapping Research and Surveying Practices. 2019 IEEE International Conference on Industrial Engineering and Engineering Management (IEEM), 800–804. doi:10.1109/ieem44572.2019.897895.
- [41] Askarinejadmiri, Z., Azim, A., Zulzallil, H., & Tieng, K. (2017). Impact Propagation of Human Errors on Software Requirements Volatility. International Journal of Advanced Computer Science and Applications, 8(2), 227–237. doi:10.14569/ijacsa.2017.080230.
- [42] Keung, J. (2009). Software Development Cost Estimation Using Analogy: A Review. 2009 Australian Software Engineering Conference, 327–338. doi:10.1109/aswec.2009.32.
- [43] Jørgensen, M. (2004). A review of studies on expert estimation of software development effort. Journal of Systems and Software, 70(1-2), 37–60. doi:10.1016/s0164-1212(02)00156-5.
- [44] The Function Point Counting Practices Manual version 4.3.1, Published by the International Function Point Users Group (IFPUG), January 2010, www.ifpug.org.
- [45] Software engineering. NESMA functional size measurement method version 2.1. Definitions and counting guidelines for the application of function point analysis. (n.d.). doi:10.3403/03235541.
- [46] Software engineering. COSMIC-FFP. A functional size measurement method. (n.d.). doi:10.3403/02770143u .
- [47] Research Design and Sampling Techniques. (2016). Field Methods in Archaeology, 33–52. doi:10.4324/9781315428413-7.
- [48] Serrador, P., & Rodney Turner, J. (2014). The Relationship between Project Success and Project Efficiency. Procedia - Social and Behavioral Sciences, 119, 75–84. doi:10.1016/j.sbspro.2014.03.011.
- [49] Fernández, D.M., Wagner, S., Kalinowski, M. et al. Naming the pain in requirements engineering. Empir Software Eng 22, 2298–2338 (2017). https://doi.org/10.1007/s10664-016-9451-7.
- [50] Abelein, U., & Paech, B. (2013). Understanding the Influence of User Participation and Involvement on System Success – a Systematic Mapping Study. Empirical Software Engineering, 20(1), 28–81. doi:10.1007/s10664-013-9278-4.
- [51] Tamai, T., & Kamata, M. I. (2009). Impact of Requirements Quality on Project Success or Failure. Design Requirements Engineering: A Ten-Year Perspective, 258–275. doi:10.1007/978-3-540-92966-6\_15.
- [52] Hull, E., Jackson, K., & Dick, J. (2010). Requirement engineering. Springer Science & Business Media, London, 6-8. doi.org/10.1007/978-1-84996-405-0.
- [53] Kalinowski, M., Felderer, M., Conte, T., Spínola, R., Prikladnicki, R., Winkler, D., Wagner, S. (2015). Preventing Incomplete/Hidden Requirements: Reflections on Survey Data from Austria and Brazil. Software Quality. The Future of Systems- and Software Development, 63–78. doi:10.1007/978-3-319-27033-3\_5.
- [54] Hairul Nizam Md Nasir, M., & Sahibuddin, S. (2015). How the PMBOK Addresses Critical Success Factors for Software Projects: A Multi-round Delphi Study. Journal of Software, 10(11), 1283–1300. doi:10.17706/jsw.10.11.1283-1300.
- [55] Peña, M., & Valerdi, R. (2014). Characterizing the Impact of Requirements Volatility on Systems Engineering Effort. Systems Engineering, 18(1), 59–70. doi:10.1111/sys.21288.
- [56] Khan, A., Kumar, C., Shameem, M., & Chandra, B. (2019). Impact of Requirements Volatility and Flexible Management on GSD Project Success: A Study Based on the Dimensions of Requirements Volatility. International Journal of Agile Systems and Management, 12(4), 1. doi:10.1504/ijasm.2019.10018758.
- [57] Bernard, H. R. (2017). Research methods in anthropology: Qualitative and quantitative approaches. Rowman & Littlefield.
- [58] Kumari, K., & Yadav, S. (2018). Linear regression analysis study. Journal of the Practice of Cardiovascular Sciences, 4(1), 33. doi:10.4103/jpcs.jpccs\_8\_18.
- [59] Taber, K. S. (2017). The Use of Cronbach’s Alpha When Developing and Reporting Research Instruments in Science Education. Research in Science Education, 48(6), 1273–1296. doi:10.1007/s11165-016-9602-2

# Dissemination and Implementation of THK-ANEKA and SAW-Based Stake Model Evaluation Website

Dewa Gede Hendra Divayana<sup>1</sup>  
Department of IT Education  
Universitas Pendidikan Ganesha  
Singaraja, Indonesia

I Putu Wisna Ariawan<sup>2</sup>  
Department of Mathematics  
Education, Universitas Pendidikan  
Ganesha, Singaraja, Indonesia

Agus Adiarta<sup>3</sup>  
Department of Electrical Education  
Universitas Pendidikan Ganesha  
Singaraja, Indonesia

**Abstract**—The purpose of this study was to provide information about the dissemination and implementation of the THK-ANEKA and SAW-based Stake model evaluation website at Vocational Schools of IT in Bali. THK is an acronym for Tri Hita Karana. ANEKA is an acronym for Akuntabilitas, Nasionalisme, Etika publik, Komitmen mutu, dan Anti korupsi (in Indonesian) or Accountability, Nationalism, Public ethics, Quality commitment, and Anti-corruption (in English). SAW is an acronym for Simple Additive Weighting. This study used a development approach with the Borg and Gall model which consists of 10 development stages. Research in 2020 was focused on the dissemination and implementation stages. The research location was at several Vocational Schools of IT in Bali Province. The subjects involved in assessing website implementation were 110 respondents. The tool used to assess was a questionnaire. The analysis technique was carried out by interpreting the effectiveness level of dissemination and implementation. It was a reference to the eleven scale effectiveness standard. The research results showed that the dissemination and implementation of the THK-ANEKA and SAW-based Stake model evaluation website at Vocational Schools of IT in Bali had gone well. It was able to be seen from the documentary evidence of the dissemination activities implementation. The percentage results of the website implementing effectiveness were 88.973% and the simulation results of implementing the SAW method which was already accurate. It showed the evaluation aspects that support the realization of positive morals and students' learning quality.

**Keywords**—Evaluation website; stake model; THK; ANEKA; SAW

## I. INTRODUCTION

Evaluation activities are very important to do to determine the effectiveness of computer learning implementation at the Vocational Schools of IT. Several evaluation models that can be used to evaluate the computer learning implementation include: CIPP [1,2]; CSE-UCLA [3]; Formative-Summatif [4]; Discrepancy [5]; and Countenance [6]. However, not all of these models can produce accurate recommendations. The expected accurate recommendations are related to aspects that support positive moral improvement and the students' learning quality in the computer learning process. One effort that can be made to obtain these accurate recommendations is to present a web-based evaluation application. This web-based evaluation application can integrate the Stake evaluation model with the THK concept, the ANEKA concept, and the SAW method.

The Stake evaluation model [7-11] is one of the evaluation models used to provide recommendations based on a description and *judgment matrix*. The THK (*Tri Hita Karana*) concept is one of the Balinese local wisdom that teaches people to recognize the three causes of happiness. The three causes of happiness [12-14], included: *Parahyangan* (good relationship with God), *Pawongan* (good relationship with fellow human beings), and *Palemahan* (good relationship with nature and the environment). ANEKA is a concept that teaches internalizing the values of a positive attitude and self-quality that must be possessed by a civil servant in Indonesia. It is as a foundation for carrying out his/her professionalism as a good servant of the country. ANEKA consist of several components [15,16], included: accountability, nationalism, public ethics, quality commitment, and anti-corruption. SAW (Simple Additive Weighting) is one of the methods in the Multi Criteria Decision Making (MCDM) [17-20], which is how it works to determine the assessment score based on the multiplication results of each alternative with the decision-maker weight.

Aspects of the Stake model were used as the basic criteria for measurement in evaluating the computer learning process at Vocational Schools of IT. ANEKA components were internalized into the description matrix which contained in the Stake model. The aim was to ensure the positive attitude and students' learning quality in the computer learning process had been in accordance with the context, process, and impact variables in the description matrix. THK components were internalized into a judgment matrix in the Stake model with the aim of being used as a main basic in determining recommendations. The SAW method was used to determine the dominant aspects that need to be encouraged to realize students' learning quality and positive moral improvement.

The THK-ANEKA and SAW-based Stake model evaluation website can be said to run optimally if it has been disseminated and implemented. Therefore, it is necessary to conduct the dissemination and implementation of this website on a wider scale. Based on these, then the right question for this research was "What are the dissemination and implementation results of the THK-ANEKA and SAW-based Stake model evaluation website at Vocational Schools of IT (case study in Bali Province)?"

Several previous studies had provided a stimulus and effect for the realization of this research. It was like the research conducted in 2018 by Ihsan and Furnham [21], which

showed the existence of several technologies that can be used as a source for assessing personality. Some of the technologies referred to included: social media, wearable technology, mobile phone, gamification, video resume, and automated personality testing. The limitation of Ihsan and Furnham's research was that it only introduced some of the technologies used for personality assessment, but it had not yet explained in detail how the technology works. Besides, Ihsan and Furnham's research only focuses on personality assessments based only on the affective domain and it had not based on cognitive and psychomotor domains. Research was conducted in 2017 by Boitshwarelo, Reedy and Billany [22] demonstrated the use of online tests to measure 21<sup>st</sup> century learning outcomes. The limitation of Boitshwarelo, Reedy and Billany's research was that it had not been discussed in detail about measuring learning outcomes in the affective and psychomotor domains. Their research only focuses on the cognitive domain as measured by using an online test. Research was conducted in 2018 by Kyllonen and Kell [23] showed a test measuring cognitive ability and personality measurement. Measurement of cognitive abilities was measured using cognitive tests, such as multiple choice and essays. Personality measurement used attitude scale questionnaires. The limitation of Kyllonen and Kell's research was that it had not shown any measurement in the psychomotor domain. Research in 2015 by Mariš [24] showed that there were character measurements based on the individual character dimension scores. The limitation of Mariš's research was that it had not been shown the measurement of cognitive and psychomotor abilities in individuals. Research in 2018 by Elmahdi, Al-Hattami, and Fawzi [25] showed a formative assessment of the student learning process used Plickers technology. The limitation of the research of Elmahdi, Al-Hattami, and Fawzi was that it had not specifically shown any assessment in the affective and psychomotor domains, because they focus on cognitive assessments. Research in 2018 was conducted by Daniawan [26] showed the use of *the SAW* method in evaluating lecturer performance in teaching. The similarity between Daniawan's research and this research was that both of them apply the SAW method in making decisions. Daniawan's research limitation was that it had not shown specific criteria for measuring the cognitive domain. Daniawan only focused on showing ten criteria in the teaching process which focused more on the affective and psychomotor domains.

Based on the research question and previous research that had provided a stimulus, then the authors were interested in conducting more in-depth research. It was related to dissemination activities and the implementation of the THK-ANEKA and SAW-based Stake evaluation website at several Vocational Schools of IT in Bali Province.

## II. METHOD

This research was development research that had carried out from 2018 to 2020. The model used in this development research was Borg and Gall [27-29] which consists of 10 stages of development. Five stages which were carried out in 2018, included: 1) research and field data collection,

2) research planning, 3) design development, 4) preliminary field test, 5) the main product revision. Two stages which had carried out in 2019, included: 1) main field test and 2) operational product revision. Three stages which had carried out in 2020, included: 1) operational field testing, 2) final product revision, 3) dissemination and implementation of the final product.

Based on the research questions previously disclosed, so the discussion in this paper focused on the dissemination and implementation stages of the final product. There were 110 respondents involved in the dissemination and implementation stage of the THK-ANEKA and SAW-based Stake model evaluation website. The 110 respondents consist of 80 students and 30 teachers from Vocational Schools of IT in Bali Province.

The tool used to obtain quantitative data in dissemination and evaluation website implementation was the questionnaires. The research location was carried out in several Vocational Schools of IT in 6 regencies on Bali Province, included: Tabanan, Buleleng, Klungkung, Gianyar, Denpasar, and Badung. The analysis technique used in this research was descriptive quantitative by interpreting the results of the effectiveness level from dissemination and implementation. It was based on the effectiveness standard which refers to the eleven's scale. The formula used to determine the effectiveness level of dissemination and implementation can be seen in equation (1) [30,31], while the standard of effectiveness which refers to the eleven's scale [32] can be seen in Table I.

The effectiveness level of dissemination and implementation

$$= \frac{f}{N} * 100\% \quad (1)$$

Notes:

f = the acquisition value total.

N = the maximum value total.

TABLE I. ELEVEN-SCALE EFFECTIVENESS STANDARDS

| Percentage of Effectiveness | Category of Effectiveness |
|-----------------------------|---------------------------|
| 0-4                         | Poor                      |
| 5-14                        | Very Bad                  |
| 15-24                       | Bad                       |
| 25-34                       | Very Less                 |
| 35-44                       | Less                      |
| 45-54                       | Elementary                |
| 55-64                       | Enough                    |
| 65-74                       | Intermediate              |
| 75-84                       | Advanced                  |
| 85-94                       | Good                      |
| 95-100                      | Excellent                 |

### III. RESULTS AND DISCUSSION

Before showing the implementation results of the *THK-ANEKA* and *SAW*-based *Stake* model evaluation website, it was necessary to carry out dissemination activities to users. Dissemination activities were carried out by holding online workshops through zoom media and direct assistance to schools. The workshop and mentoring activities can be seen in Fig. 1. Details of the material provided in the dissemination activities can be seen in Table II.



Fig. 1. Dissemination Activities.

TABLE II. MATERIALS PROVIDED AT DISSEMINATION

| No | Materials  |
|----|--|
| 1  | Introduction to the purpose and benefits of the <i>THK-ANEKA</i> and <i>SAW</i> -based <i>Stake</i> model evaluation website |
| 2  | Procedures for managing the login form   |
| 3  | Procedures for managing the main menu form   |
| 4  | Procedures for managing the input indicator form   |
| 5  | Procedures for managing the weight input form  |
| 6  | Procedures for managing the antecedents form in the <i>description matrix</i>  |
| 7  | Procedures for managing the transactions form in the <i>description matrix</i>   |
| 8  | Procedures for managing the outcomes form in the <i>description matrix</i>   |
| 9  | Procedures for managing the <i>judgment</i> form matrix  |
| 10 | Procedures for managing the recommendation form  |
| 11 | Procedures for managing the decision form  |

The successful implementation of the *THK-ANEKA* and *SAW*-based *Stake* model evaluation website at several Vocational Schools of IT in Bali was able to be obtained from the assessment results of 110 respondents (30 teachers and 80 students). The assessment results of all respondents can be seen in Table III. The assessment activities documentation of evaluation website implementation can be seen in Fig. 2.

The successful implementation evidence of the evaluation website also was obtained from the results of *SAW* method calculation accuracy in addition to the assessment results from the 110 respondents. The *SAW* calculation process can be carried out if simulation data are provided (can be seen in Table IV) and the weight of decision-makers (can be seen in Table V).

TABLE III. RESPONDENTS ASSESSMENT RESULTS TO THE IMPLEMENTATION OF *THK-ANEKA* AND *SAW*-BASED *STAKE* MODEL EVALUATION WEBSITE

| No | Respondents | Items- |   |   |   |   |   |   |   |   |    |    |    |    |    |    |    |    |    |    | Σ | Percentage of Effectiveness (%) |        |
|----|-------------|--------|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|---|---------------------------------|--------|
|    |             | 1      | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 |   |                                 | 20     |
| 1  | Teacher-1   | 5      | 4 | 5 | 4 | 4 | 5 | 4 | 4 | 5 | 5  | 5  | 5  | 4  | 5  | 5  | 5  | 4  | 5  | 4  | 5 | 92                              | 92.000 |
| 2  | Teacher-2   | 5      | 4 | 5 | 4 | 5 | 4 | 4 | 4 | 4 | 5  | 4  | 5  | 4  | 5  | 4  | 5  | 4  | 5  | 4  | 5 | 89                              | 89.000 |
| 3  | Teacher-3   | 4      | 5 | 4 | 4 | 5 | 4 | 5 | 4 | 4 | 4  | 5  | 5  | 5  | 4  | 5  | 4  | 5  | 4  | 5  | 4 | 89                              | 89.000 |
| 4  | Teacher-4   | 5      | 4 | 5 | 4 | 4 | 5 | 5 | 4 | 5 | 4  | 4  | 5  | 5  | 4  | 5  | 4  | 4  | 4  | 5  | 4 | 89                              | 89.000 |
| 5  | Teacher-5   | 4      | 4 | 4 | 4 | 5 | 4 | 4 | 4 | 5 | 4  | 5  | 4  | 4  | 4  | 5  | 4  | 4  | 4  | 5  | 4 | 85                              | 85.000 |
| 6  | Teacher-6   | 4      | 4 | 5 | 4 | 4 | 5 | 5 | 5 | 4 | 4  | 4  | 5  | 5  | 5  | 4  | 4  | 5  | 5  | 4  | 4 | 89                              | 89.000 |
| 7  | Teacher-7   | 4      | 5 | 5 | 4 | 5 | 4 | 4 | 4 | 5 | 4  | 5  | 4  | 4  | 4  | 5  | 4  | 4  | 4  | 5  | 4 | 87                              | 87.000 |
| 8  | Teacher-8   | 5      | 4 | 4 | 4 | 5 | 5 | 5 | 4 | 4 | 4  | 5  | 5  | 5  | 4  | 4  | 4  | 4  | 4  | 4  | 4 | 87                              | 87.000 |
| 9  | Teacher-9   | 5      | 5 | 4 | 5 | 4 | 4 | 5 | 5 | 4 | 5  | 4  | 4  | 5  | 5  | 4  | 4  | 5  | 4  | 4  | 5 | 90                              | 90.000 |
| 10 | Teacher-10  | 4      | 4 | 5 | 4 | 5 | 5 | 4 | 4 | 4 | 5  | 4  | 5  | 4  | 4  | 5  | 5  | 4  | 5  | 5  | 4 | 89                              | 89.000 |
| 11 | Teacher-11  | 5      | 4 | 4 | 5 | 4 | 4 | 5 | 5 | 4 | 5  | 4  | 5  | 4  | 5  | 4  | 4  | 5  | 4  | 4  | 5 | 89                              | 89.000 |
| 12 | Teacher-12  | 4      | 5 | 5 | 4 | 5 | 5 | 4 | 4 | 4 | 4  | 5  | 4  | 4  | 4  | 5  | 5  | 4  | 5  | 5  | 4 | 89                              | 89.000 |
| 13 | Teacher-13  | 4      | 4 | 4 | 5 | 5 | 4 | 5 | 5 | 5 | 5  | 4  | 5  | 4  | 5  | 4  | 4  | 5  | 4  | 4  | 5 | 90                              | 90.000 |

| No | Respondents | Items- |   |   |   |   |   |   |   |   |    |    |    |    |    |    |    |    |    |    |    | Σ  | Percentage of Effectiveness (%) |
|----|-------------|--------|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|----|---------------------------------|
|    |             | 1      | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |    |                                 |
| 14 | Teacher-14  | 4      | 4 | 4 | 4 | 5 | 5 | 4 | 4 | 4 | 4  | 4  | 4  | 4  | 4  | 5  | 5  | 4  | 5  | 5  | 4  | 86 | 86.000                          |
| 15 | Teacher-15  | 4      | 4 | 5 | 5 | 5 | 4 | 5 | 5 | 5 | 5  | 4  | 4  | 4  | 5  | 4  | 4  | 5  | 4  | 4  | 5  | 90 | 90.000                          |
| 16 | Teacher-16  | 4      | 5 | 4 | 4 | 4 | 5 | 4 | 4 | 5 | 4  | 5  | 5  | 4  | 4  | 5  | 5  | 5  | 5  | 5  | 4  | 90 | 90.000                          |
| 17 | Teacher-17  | 4      | 4 | 5 | 5 | 5 | 4 | 5 | 5 | 5 | 5  | 4  | 4  | 4  | 5  | 4  | 4  | 4  | 4  | 5  | 5  | 90 | 90.000                          |
| 18 | Teacher-18  | 5      | 5 | 4 | 4 | 4 | 5 | 4 | 4 | 5 | 4  | 4  | 4  | 5  | 5  | 5  | 5  | 4  | 5  | 5  | 4  | 90 | 90.000                          |
| 19 | Teacher-19  | 4      | 4 | 5 | 5 | 5 | 5 | 5 | 5 | 4 | 5  | 5  | 5  | 4  | 4  | 4  | 5  | 5  | 4  | 4  | 4  | 91 | 91.000                          |
| 20 | Teacher-20  | 5      | 5 | 4 | 4 | 5 | 4 | 4 | 5 | 4 | 5  | 5  | 4  | 5  | 5  | 5  | 4  | 4  | 5  | 5  | 4  | 91 | 91.000                          |
| 21 | Teacher-21  | 5      | 5 | 5 | 5 | 5 | 4 | 5 | 4 | 5 | 4  | 4  | 5  | 4  | 4  | 4  | 4  | 5  | 4  | 4  | 4  | 89 | 89.000                          |
| 22 | Teacher-22  | 5      | 4 | 4 | 5 | 5 | 4 | 5 | 4 | 4 | 5  | 5  | 5  | 5  | 5  | 4  | 4  | 5  | 5  | 5  | 4  | 92 | 92.000                          |
| 23 | Teacher-23  | 5      | 4 | 5 | 4 | 5 | 5 | 4 | 4 | 5 | 4  | 4  | 4  | 4  | 5  | 4  | 5  | 4  | 4  | 5  | 5  | 89 | 89.000                          |
| 24 | Teacher-24  | 5      | 4 | 4 | 5 | 4 | 4 | 5 | 4 | 5 | 5  | 5  | 4  | 5  | 4  | 5  | 4  | 5  | 5  | 4  | 5  | 91 | 91.000                          |
| 25 | Teacher-25  | 4      | 4 | 5 | 4 | 5 | 4 | 4 | 4 | 4 | 4  | 5  | 5  | 4  | 4  | 5  | 4  | 4  | 4  | 4  | 5  | 86 | 86.000                          |
| 26 | Teacher-26  | 4      | 4 | 4 | 5 | 4 | 4 | 5 | 5 | 4 | 5  | 4  | 4  | 4  | 5  | 4  | 4  | 5  | 5  | 4  | 4  | 87 | 87.000                          |
| 27 | Teacher-27  | 4      | 4 | 4 | 4 | 4 | 5 | 4 | 4 | 4 | 5  | 4  | 4  | 5  | 4  | 5  | 4  | 4  | 5  | 4  | 4  | 85 | 85.000                          |
| 28 | Teacher-28  | 4      | 5 | 4 | 5 | 4 | 4 | 5 | 5 | 5 | 4  | 4  | 4  | 4  | 4  | 4  | 4  | 5  | 5  | 5  | 4  | 88 | 88.000                          |
| 29 | Teacher-29  | 5      | 4 | 4 | 4 | 5 | 5 | 4 | 4 | 4 | 5  | 4  | 5  | 5  | 4  | 5  | 4  | 5  | 4  | 4  | 5  | 89 | 89.000                          |
| 30 | Teacher-30  | 4      | 5 | 4 | 5 | 4 | 5 | 5 | 5 | 4 | 4  | 4  | 4  | 4  | 4  | 4  | 4  | 4  | 5  | 4  | 4  | 86 | 86.000                          |
| 31 | Student-1   | 5      | 4 | 4 | 4 | 4 | 4 | 4 | 5 | 4 | 5  | 4  | 5  | 5  | 5  | 4  | 4  | 5  | 5  | 5  | 4  | 89 | 89.000                          |
| 32 | Student-2   | 4      | 5 | 4 | 5 | 4 | 4 | 5 | 4 | 4 | 5  | 5  | 4  | 4  | 5  | 4  | 5  | 4  | 5  | 5  | 4  | 89 | 89.000                          |
| 33 | Student-3   | 5      | 4 | 5 | 4 | 4 | 4 | 5 | 4 | 4 | 4  | 5  | 5  | 5  | 4  | 4  | 4  | 5  | 4  | 4  | 4  | 87 | 87.000                          |
| 34 | Student-4   | 4      | 5 | 4 | 5 | 4 | 5 | 4 | 4 | 4 | 4  | 4  | 4  | 4  | 5  | 4  | 4  | 5  | 5  | 4  | 5  | 87 | 87.000                          |
| 35 | Student-5   | 4      | 4 | 5 | 4 | 5 | 4 | 5 | 4 | 5 | 4  | 4  | 5  | 5  | 4  | 5  | 4  | 4  | 4  | 4  | 4  | 87 | 87.000                          |
| 36 | Student-6   | 5      | 5 | 4 | 5 | 4 | 4 | 4 | 4 | 4 | 4  | 4  | 5  | 4  | 5  | 5  | 5  | 4  | 4  | 4  | 4  | 87 | 87.000                          |
| 37 | Student-7   | 4      | 4 | 5 | 4 | 4 | 5 | 4 | 4 | 4 | 5  | 4  | 5  | 4  | 4  | 5  | 5  | 4  | 4  | 5  | 4  | 87 | 87.000                          |
| 38 | Student-8   | 5      | 5 | 4 | 5 | 4 | 5 | 4 | 5 | 4 | 4  | 5  | 4  | 4  | 5  | 4  | 4  | 4  | 4  | 4  | 5  | 88 | 88.000                          |
| 39 | Student-9   | 4      | 4 | 5 | 4 | 5 | 4 | 4 | 4 | 5 | 5  | 4  | 5  | 4  | 4  | 5  | 4  | 5  | 4  | 4  | 5  | 88 | 88.000                          |
| 40 | Student-10  | 5      | 5 | 4 | 5 | 4 | 5 | 4 | 5 | 4 | 4  | 4  | 4  | 4  | 4  | 4  | 4  | 4  | 5  | 4  | 4  | 86 | 86.000                          |
| 41 | Student-11  | 4      | 4 | 5 | 4 | 4 | 4 | 4 | 4 | 5 | 5  | 5  | 4  | 4  | 5  | 4  | 4  | 5  | 5  | 5  | 4  | 88 | 88.000                          |
| 42 | Student-12  | 5      | 4 | 4 | 5 | 5 | 4 | 4 | 5 | 4 | 4  | 4  | 5  | 4  | 4  | 4  | 5  | 4  | 5  | 5  | 4  | 88 | 88.000                          |
| 43 | Student-13  | 4      | 4 | 5 | 4 | 4 | 5 | 4 | 4 | 5 | 4  | 4  | 5  | 5  | 4  | 4  | 4  | 5  | 4  | 4  | 4  | 86 | 86.000                          |
| 44 | Student-14  | 5      | 4 | 4 | 5 | 5 | 5 | 5 | 4 | 4 | 5  | 4  | 4  | 5  | 4  | 4  | 4  | 4  | 4  | 5  | 5  | 89 | 89.000                          |
| 45 | Student-15  | 4      | 4 | 5 | 4 | 4 | 5 | 5 | 4 | 5 | 5  | 5  | 4  | 4  | 4  | 4  | 5  | 4  | 5  | 4  | 4  | 88 | 88.000                          |
| 46 | Student-16  | 5      | 4 | 4 | 5 | 5 | 4 | 4 | 4 | 4 | 5  | 5  | 4  | 4  | 4  | 5  | 4  | 5  | 4  | 5  | 4  | 88 | 88.000                          |
| 47 | Student-17  | 4      | 5 | 5 | 4 | 5 | 4 | 4 | 4 | 5 | 4  | 4  | 4  | 5  | 4  | 4  | 5  | 4  | 5  | 4  | 5  | 88 | 88.000                          |
| 48 | Student-18  | 5      | 4 | 4 | 5 | 5 | 5 | 5 | 5 | 4 | 4  | 5  | 4  | 4  | 5  | 5  | 4  | 5  | 4  | 5  | 5  | 92 | 92.000                          |
| 49 | Student-19  | 4      | 4 | 5 | 4 | 4 | 4 | 5 | 4 | 4 | 4  | 4  | 5  | 5  | 5  | 4  | 5  | 4  | 4  | 5  | 5  | 88 | 88.000                          |
| 50 | Student-20  | 4      | 5 | 4 | 5 | 5 | 4 | 4 | 5 | 5 | 4  | 5  | 4  | 4  | 4  | 4  | 4  | 4  | 5  | 4  | 4  | 87 | 87.000                          |
| 51 | Student-21  | 5      | 4 | 5 | 5 | 4 | 5 | 5 | 4 | 4 | 5  | 4  | 5  | 4  | 4  | 4  | 4  | 5  | 5  | 5  | 4  | 90 | 90.000                          |
| 52 | Student-22  | 5      | 4 | 4 | 4 | 5 | 4 | 4 | 5 | 5 | 4  | 5  | 4  | 5  | 4  | 4  | 5  | 4  | 4  | 4  | 5  | 88 | 88.000                          |
| 53 | Student-23  | 4      | 5 | 5 | 5 | 4 | 5 | 5 | 4 | 4 | 5  | 4  | 5  | 5  | 5  | 5  | 4  | 5  | 4  | 5  | 5  | 93 | 93.000                          |
| 54 | Student-24  | 5      | 4 | 4 | 5 | 5 | 4 | 4 | 5 | 4 | 4  | 5  | 4  | 5  | 5  | 4  | 5  | 4  | 5  | 4  | 4  | 89 | 89.000                          |
| 55 | Student-25  | 4      | 5 | 4 | 4 | 5 | 5 | 5 | 5 | 4 | 4  | 5  | 5  | 4  | 4  | 5  | 4  | 5  | 5  | 5  | 4  | 91 | 91.000                          |



| No | Respondents | Items- |   |   |   |   |   |   |   |   |    |    |    |    |    |    |    |    |    |    |    | Σ  | Percentage of Effectiveness (%) |
|----|-------------|--------|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|----|---------------------------------|
|    |             | 1      | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |    |                                 |
| 56 | Student-26  | 5      | 4 | 4 | 5 | 4 | 4 | 5 | 4 | 4 | 4  | 4  | 5  | 5  | 4  | 4  | 4  | 4  | 4  | 4  | 5  | 86 | 86.000                          |
| 57 | Student-27  | 4      | 5 | 4 | 4 | 4 | 5 | 4 | 5 | 4 | 4  | 5  | 4  | 4  | 4  | 4  | 4  | 4  | 4  | 5  | 5  | 86 | 86.000                          |
| 58 | Student-28  | 5      | 4 | 5 | 5 | 4 | 4 | 4 | 4 | 5 | 5  | 4  | 5  | 4  | 4  | 5  | 4  | 4  | 5  | 4  | 4  | 88 | 88.000                          |
| 59 | Student-29  | 4      | 5 | 4 | 4 | 5 | 4 | 4 | 4 | 4 | 4  | 5  | 4  | 5  | 4  | 4  | 5  | 5  | 4  | 5  | 4  | 87 | 87.000                          |
| 60 | Student-30  | 5      | 4 | 5 | 5 | 5 | 4 | 5 | 4 | 5 | 5  | 4  | 5  | 5  | 5  | 4  | 4  | 5  | 5  | 4  | 5  | 93 | 93.000                          |
| 61 | Student-31  | 4      | 5 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 5  | 5  | 4  | 4  | 4  | 4  | 5  | 4  | 4  | 5  | 5  | 86 | 86.000                          |
| 62 | Student-32  | 5      | 4 | 5 | 5 | 4 | 4 | 5 | 4 | 5 | 4  | 4  | 4  | 4  | 4  | 5  | 4  | 4  | 4  | 5  | 5  | 88 | 88.000                          |
| 63 | Student-33  | 4      | 5 | 4 | 5 | 4 | 5 | 4 | 5 | 4 | 5  | 4  | 4  | 5  | 4  | 4  | 5  | 4  | 5  | 4  | 4  | 88 | 88.000                          |
| 64 | Student-34  | 5      | 5 | 5 | 4 | 5 | 4 | 5 | 4 | 5 | 4  | 5  | 5  | 4  | 5  | 5  | 4  | 5  | 4  | 5  | 4  | 92 | 92.000                          |
| 65 | Student-35  | 4      | 5 | 4 | 5 | 4 | 5 | 4 | 5 | 4 | 5  | 5  | 4  | 4  | 4  | 5  | 5  | 4  | 5  | 4  | 5  | 90 | 90.000                          |
| 66 | Student-36  | 5      | 5 | 5 | 5 | 5 | 4 | 5 | 5 | 5 | 4  | 5  | 4  | 4  | 5  | 4  | 4  | 5  | 4  | 5  | 5  | 93 | 93.000                          |
| 67 | Student-37  | 4      | 4 | 4 | 4 | 5 | 5 | 5 | 4 | 4 | 5  | 4  | 5  | 5  | 4  | 5  | 4  | 4  | 4  | 5  | 4  | 88 | 88.000                          |
| 68 | Student-38  | 5      | 4 | 5 | 4 | 4 | 5 | 5 | 4 | 5 | 4  | 5  | 5  | 4  | 5  | 4  | 5  | 4  | 5  | 5  | 5  | 92 | 92.000                          |
| 69 | Student-39  | 4      | 5 | 4 | 5 | 5 | 4 | 4 | 4 | 4 | 5  | 4  | 4  | 5  | 4  | 5  | 5  | 4  | 4  | 5  | 5  | 89 | 89.000                          |
| 70 | Student-40  | 5      | 4 | 5 | 4 | 4 | 5 | 4 | 4 | 5 | 5  | 4  | 5  | 4  | 4  | 4  | 5  | 4  | 5  | 4  | 4  | 88 | 88.000                          |
| 71 | Student-41  | 5      | 5 | 4 | 4 | 5 | 5 | 5 | 4 | 5 | 4  | 4  | 4  | 5  | 4  | 5  | 5  | 4  | 4  | 5  | 4  | 90 | 90.000                          |
| 72 | Student-42  | 4      | 4 | 4 | 5 | 4 | 5 | 5 | 4 | 4 | 4  | 5  | 4  | 4  | 5  | 4  | 4  | 4  | 5  | 5  | 5  | 88 | 88.000                          |
| 73 | Student-43  | 5      | 5 | 5 | 4 | 5 | 4 | 4 | 4 | 4 | 5  | 5  | 5  | 5  | 5  | 5  | 4  | 4  | 4  | 5  | 5  | 92 | 92.000                          |
| 74 | Student-44  | 4      | 4 | 4 | 5 | 4 | 4 | 4 | 4 | 4 | 4  | 5  | 5  | 4  | 5  | 5  | 4  | 5  | 5  | 4  | 4  | 87 | 87.000                          |
| 75 | Student-45  | 5      | 5 | 5 | 5 | 5 | 4 | 4 | 4 | 5 | 5  | 4  | 4  | 5  | 4  | 4  | 4  | 5  | 4  | 5  | 4  | 90 | 90.000                          |
| 76 | Student-46  | 4      | 4 | 4 | 5 | 5 | 4 | 5 | 4 | 4 | 5  | 5  | 4  | 4  | 4  | 4  | 4  | 5  | 4  | 4  | 5  | 87 | 87.000                          |
| 77 | Student-47  | 5      | 5 | 5 | 4 | 4 | 4 | 4 | 5 | 5 | 4  | 4  | 4  | 4  | 5  | 4  | 5  | 5  | 5  | 4  | 5  | 90 | 90.000                          |
| 78 | Student-48  | 4      | 4 | 4 | 5 | 4 | 5 | 4 | 4 | 5 | 4  | 5  | 4  | 4  | 4  | 5  | 4  | 5  | 5  | 4  | 4  | 87 | 87.000                          |
| 79 | Student-49  | 5      | 5 | 5 | 5 | 5 | 5 | 5 | 4 | 4 | 5  | 4  | 4  | 5  | 5  | 5  | 5  | 4  | 4  | 4  | 5  | 93 | 93.000                          |
| 80 | Student-50  | 4      | 4 | 4 | 4 | 4 | 5 | 5 | 4 | 5 | 5  | 5  | 4  | 5  | 5  | 5  | 5  | 5  | 5  | 5  | 5  | 93 | 93.000                          |
| 81 | Student-51  | 4      | 5 | 5 | 4 | 5 | 4 | 4 | 4 | 4 | 5  | 5  | 4  | 4  | 5  | 4  | 4  | 4  | 5  | 4  | 5  | 88 | 88.000                          |
| 82 | Student-52  | 4      | 4 | 4 | 5 | 4 | 4 | 4 | 4 | 5 | 4  | 4  | 4  | 5  | 5  | 5  | 4  | 5  | 4  | 4  | 5  | 87 | 87.000                          |
| 83 | Student-53  | 5      | 5 | 5 | 4 | 5 | 4 | 4 | 4 | 4 | 5  | 5  | 4  | 4  | 5  | 5  | 4  | 4  | 5  | 5  | 5  | 91 | 91.000                          |
| 84 | Student-54  | 4      | 4 | 4 | 5 | 5 | 4 | 5 | 4 | 5 | 4  | 4  | 4  | 5  | 4  | 4  | 4  | 4  | 4  | 4  | 5  | 86 | 86.000                          |
| 85 | Student-55  | 5      | 5 | 5 | 4 | 4 | 4 | 4 | 4 | 5 | 4  | 4  | 5  | 4  | 5  | 4  | 4  | 5  | 5  | 5  | 4  | 89 | 89.000                          |
| 86 | Student-56  | 4      | 4 | 4 | 5 | 5 | 5 | 4 | 4 | 5 | 5  | 5  | 4  | 4  | 5  | 5  | 5  | 5  | 5  | 5  | 4  | 92 | 92.000                          |
| 87 | Student-57  | 5      | 5 | 5 | 4 | 4 | 4 | 5 | 5 | 4 | 4  | 4  | 4  | 5  | 4  | 4  | 4  | 4  | 5  | 5  | 4  | 88 | 88.000                          |
| 88 | Student-58  | 4      | 5 | 4 | 5 | 5 | 5 | 4 | 4 | 5 | 5  | 5  | 4  | 4  | 5  | 5  | 5  | 5  | 4  | 4  | 4  | 91 | 91.000                          |
| 89 | Student-59  | 5      | 5 | 5 | 4 | 4 | 4 | 5 | 5 | 4 | 4  | 4  | 4  | 5  | 4  | 4  | 4  | 4  | 5  | 4  | 4  | 87 | 87.000                          |
| 90 | Student-60  | 4      | 5 | 4 | 4 | 5 | 5 | 4 | 4 | 5 | 4  | 4  | 4  | 4  | 5  | 4  | 4  | 4  | 5  | 5  | 5  | 88 | 88.000                          |
| 91 | Student-61  | 5      | 4 | 5 | 5 | 5 | 4 | 5 | 4 | 4 | 5  | 5  | 5  | 5  | 5  | 5  | 4  | 5  | 4  | 4  | 4  | 92 | 92.000                          |
| 92 | Student-62  | 5      | 5 | 4 | 5 | 4 | 4 | 4 | 5 | 5 | 4  | 4  | 4  | 4  | 5  | 5  | 4  | 4  | 5  | 5  | 5  | 90 | 90.000                          |
| 93 | Student-63  | 4      | 4 | 5 | 4 | 4 | 4 | 5 | 4 | 4 | 5  | 5  | 5  | 5  | 4  | 4  | 4  | 5  | 5  | 5  | 4  | 89 | 89.000                          |
| 94 | Student-64  | 5      | 5 | 4 | 4 | 4 | 5 | 4 | 4 | 5 | 4  | 4  | 4  | 4  | 5  | 4  | 4  | 4  | 5  | 5  | 4  | 87 | 87.000                          |
| 95 | Student-65  | 4      | 4 | 5 | 5 | 5 | 5 | 5 | 4 | 4 | 5  | 4  | 4  | 4  | 5  | 5  | 5  | 5  | 4  | 4  | 4  | 90 | 90.000                          |
| 96 | Student-66  | 5      | 5 | 4 | 4 | 4 | 5 | 5 | 4 | 5 | 5  | 4  | 5  | 5  | 4  | 4  | 4  | 4  | 5  | 4  | 4  | 89 | 89.000                          |
| 97 | Student-67  | 4      | 5 | 5 | 5 | 5 | 4 | 4 | 4 | 4 | 4  | 5  | 4  | 4  | 5  | 5  | 5  | 4  | 5  | 5  | 5  | 91 | 91.000                          |

| No             | Respondents | Items- |   |   |   |   |   |   |   |   |    |    |    |    |    |    |    |    |    |    |    | Σ  | Percentage of Effectiveness (%) |
|----------------|-------------|--------|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|----|---------------------------------|
|                |             | 1      | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |    |                                 |
| 98             | Student-68  | 4      | 5 | 4 | 4 | 5 | 5 | 4 | 5 | 5 | 5  | 4  | 5  | 5  | 4  | 4  | 4  | 5  | 4  | 4  | 4  | 89 | 89.000                          |
| 99             | Student-69  | 4      | 5 | 5 | 5 | 4 | 4 | 5 | 5 | 4 | 4  | 5  | 4  | 4  | 5  | 4  | 4  | 5  | 5  | 5  | 4  | 90 | 90.000                          |
| 100            | Student-70  | 5      | 4 | 4 | 4 | 5 | 5 | 4 | 5 | 5 | 5  | 4  | 5  | 5  | 4  | 5  | 5  | 4  | 5  | 4  | 5  | 92 | 92.000                          |
| 101            | Student-71  | 4      | 5 | 5 | 5 | 4 | 4 | 5 | 5 | 5 | 4  | 5  | 5  | 4  | 5  | 5  | 4  | 5  | 4  | 4  | 4  | 91 | 91.000                          |
| 102            | Student-72  | 5      | 4 | 4 | 4 | 5 | 5 | 4 | 5 | 5 | 5  | 5  | 5  | 4  | 5  | 4  | 4  | 4  | 5  | 5  | 5  | 92 | 92.000                          |
| 103            | Student-73  | 4      | 5 | 4 | 4 | 5 | 4 | 5 | 4 | 5 | 4  | 5  | 5  | 5  | 5  | 4  | 4  | 4  | 4  | 4  | 5  | 89 | 89.000                          |
| 104            | Student-74  | 5      | 4 | 4 | 4 | 5 | 5 | 5 | 4 | 5 | 5  | 4  | 4  | 4  | 4  | 5  | 4  | 5  | 5  | 5  | 4  | 90 | 90.000                          |
| 105            | Student-75  | 4      | 5 | 4 | 4 | 4 | 5 | 5 | 4 | 4 | 4  | 5  | 4  | 5  | 5  | 4  | 4  | 4  | 5  | 5  | 5  | 89 | 89.000                          |
| 106            | Student-76  | 5      | 4 | 4 | 5 | 5 | 4 | 4 | 4 | 5 | 5  | 4  | 5  | 4  | 4  | 5  | 4  | 4  | 5  | 4  | 5  | 89 | 89.000                          |
| 107            | Student-77  | 4      | 5 | 5 | 4 | 4 | 5 | 4 | 4 | 4 | 4  | 5  | 4  | 5  | 5  | 4  | 4  | 4  | 5  | 4  | 5  | 88 | 88.000                          |
| 108            | Student-78  | 5      | 4 | 4 | 4 | 5 | 5 | 5 | 4 | 5 | 5  | 4  | 5  | 5  | 4  | 5  | 5  | 5  | 5  | 4  | 4  | 92 | 92.000                          |
| 109            | Student-79  | 4      | 4 | 5 | 5 | 4 | 5 | 5 | 4 | 5 | 4  | 5  | 5  | 4  | 4  | 4  | 4  | 4  | 4  | 4  | 4  | 87 | 87.000                          |
| 110            | Student-80  | 5      | 5 | 4 | 4 | 5 | 4 | 4 | 4 | 5 | 5  | 5  | 4  | 4  | 5  | 5  | 5  | 4  | 5  | 5  | 4  | 91 | 91.000                          |
| <b>Average</b> |             |        |   |   |   |   |   |   |   |   |    |    |    |    |    |    |    |    |    |    |    |    | <b>88.973</b>                   |



Fig. 2. Assessment Activities to the Implementation of THK-ANEKA and SAW-based Stake Model Evaluation Website.

TABLE IV. SAW METHOD SIMULATION DATA

| No       | Aspects of <i>Tri Hita Karana</i>   | Components of ANEKA |             |               |                    |                 |
|----------|---|---------------------|-------------|---------------|--------------------|-----------------|
|          |   | Accountability      | Nationalism | Public Ethics | Quality Commitment | Anti-Corruption |
| <b>A</b> | <b>Parahyangan</b>  |                     |             |               |                    |                 |
| 1        | It is consistently carry out prayer activities before the lesson begins and after the end of the learning process | 4.386               | 4.455       | 4.500         | 4.432              | 4.386           |
| 2        | It is consistent respect for the way of prayer between students from different religions                          | 4.455               | 4.568       | 4.500         | 4.432              | 4.500           |
| <b>B</b> | <b>Pawongan</b>   |                     |             |               |                    |                 |
| 3        | It is maintain order in the learning process  | 4.364               | 4.386       | 4.386         | 4.341              | 4.477           |
| 4        | It is able to respect other people’s opinions   | 4.455               | 4.341       | 4.364         | 4.318              | 4.364           |
| 5        | It is able to work well together when completing group assignments  | 4.455               | 4.523       | 4.432         | 4.409              | 4.500           |
| 6        | It is always respect teachers and headmaster  | 4.568               | 4.364       | 4.455         | 4.341              | 4.477           |
| 7        | It is able to interact well and actively with all school members  | 4.477               | 4.386       | 4.477         | 4.455              | 4.500           |
| <b>C</b> | <b>Palemahan</b>  |                     |             |               |                    |                 |
| 8        | It is maintain the cleanliness of classrooms and the environment around the school consistently                   | 4.477               | 4.364       | 4.455         | 4.545              | 4.545           |
| 9        | It is able to maintain the cleanliness and facilities integrity to support the learning process                   | 4.545               | 4.455       | 4.409         | 4.409              | 4.500           |
| 10       | It is always obey the school rules  | 4.386               | 4.523       | 4.659         | 4.455              | 4.477           |

Based on the simulation data shown in Table IV and determining that all ANEKA components are included in the benefit attribute, the normalization calculation process can be carried out. The formula used for normalization calculations [17] refers to equation (2).

$$r_{ij} = \begin{cases} \frac{x_{ij}}{\text{Max}_i x_{ij}} & \text{if } j \text{ is benefit attribute} \\ \frac{\text{Min}_i x_{ij}}{x_{ij}} & \text{if } j \text{ is cost attribute} \end{cases} \quad (2)$$

Notes:

$r_{ij}$  = normalized performance rating score

$x_{ij}$  = attribute value of each criterion

Cost = if the lowest value is the best

Benefit = if the highest value is the best

Min  $x_{ij}$  = the lowest value of each criterion

Max  $x_{ij}$  = the highest value of each criterion

TABLE V. WEIGHTS FROM DECISION-MAKERS

| Components of ANEKA | Weights |
|---------------------|---------|
| Accountability      | 30%     |
| Nationalism         | 30%     |
| Public Ethics       | 30%     |
| Quality Commitment  | 30%     |
| Anti-Corruption     | 30%     |

The simulation calculation process can be explained as follows

$$\begin{aligned}
 r_{11} &= \frac{4.386}{\max\{4.386; 4.455; 4.364; 4.455; 4.455; 4.568; 4.477; 4.477; 4.545; 4.386\}} = \frac{4.386}{4.568} = 0.960 \\
 r_{21} &= \frac{4.455}{\max\{4.386; 4.455; 4.364; 4.455; 4.455; 4.568; 4.477; 4.477; 4.545; 4.386\}} = \frac{4.455}{4.568} = 0.975 \\
 r_{31} &= \frac{4.364}{\max\{4.386; 4.455; 4.364; 4.455; 4.455; 4.568; 4.477; 4.477; 4.545; 4.386\}} = \frac{4.364}{4.568} = 0.955 \\
 r_{41} &= \frac{4.455}{\max\{4.386; 4.455; 4.364; 4.455; 4.455; 4.568; 4.477; 4.477; 4.545; 4.386\}} = \frac{4.455}{4.568} = 0.975 \\
 r_{51} &= \frac{4.455}{\max\{4.386; 4.455; 4.364; 4.455; 4.455; 4.568; 4.477; 4.477; 4.545; 4.386\}} = \frac{4.455}{4.568} = 0.975 \\
 r_{61} &= \frac{4.568}{\max\{4.386; 4.455; 4.364; 4.455; 4.455; 4.568; 4.477; 4.477; 4.545; 4.386\}} = \frac{4.568}{4.568} = 1.000 \\
 r_{71} &= \frac{4.477}{\max\{4.386; 4.455; 4.364; 4.455; 4.455; 4.568; 4.477; 4.477; 4.545; 4.386\}} = \frac{4.477}{4.568} = 0.980
 \end{aligned}$$

|                |  |                       |           |
|----------------|--|-----------------------|-----------|
| $\Gamma_{81}$  | $\frac{4.477}{\max\{4.386; 4.455; 4.364; 4.455; 4.455; 4.568; 4.477; 4.477; 4.545; 4.386\}}$ | $\frac{4.477}{4.568}$ | $= 0.980$ |
| $\Gamma_{91}$  | $\frac{4.545}{\max\{4.386; 4.455; 4.364; 4.455; 4.455; 4.568; 4.477; 4.477; 4.545; 4.386\}}$ | $\frac{4.545}{4.568}$ | $= 0.995$ |
| $\Gamma_{101}$ | $\frac{4.386}{\max\{4.386; 4.455; 4.364; 4.455; 4.455; 4.568; 4.477; 4.477; 4.545; 4.386\}}$ | $\frac{4.386}{4.568}$ | $= 0.960$ |
| $\Gamma_{12}$  | $\frac{4.455}{\max\{4.455; 4.568; 4.386; 4.341; 4.523; 4.364; 4.386; 4.364; 4.455; 4.523\}}$ | $\frac{4.455}{4.568}$ | $= 0.975$ |
| $\Gamma_{22}$  | $\frac{4.568}{\max\{4.455; 4.568; 4.386; 4.341; 4.523; 4.364; 4.386; 4.364; 4.455; 4.523\}}$ | $\frac{4.568}{4.568}$ | $= 1.000$ |
| $\Gamma_{32}$  | $\frac{4.386}{\max\{4.455; 4.568; 4.386; 4.341; 4.523; 4.364; 4.386; 4.364; 4.455; 4.523\}}$ | $\frac{4.386}{4.568}$ | $= 0.960$ |
| $\Gamma_{42}$  | $\frac{4.341}{\max\{4.455; 4.568; 4.386; 4.341; 4.523; 4.364; 4.386; 4.364; 4.455; 4.523\}}$ | $\frac{4.341}{4.568}$ | $= 0.950$ |
| $\Gamma_{52}$  | $\frac{4.523}{\max\{4.455; 4.568; 4.386; 4.341; 4.523; 4.364; 4.386; 4.364; 4.455; 4.523\}}$ | $\frac{4.523}{4.568}$ | $= 0.990$ |
| $\Gamma_{62}$  | $\frac{4.364}{\max\{4.455; 4.568; 4.386; 4.341; 4.523; 4.364; 4.386; 4.364; 4.455; 4.523\}}$ | $\frac{4.364}{4.568}$ | $= 0.955$ |
| $\Gamma_{72}$  | $\frac{4.386}{\max\{4.455; 4.568; 4.386; 4.341; 4.523; 4.364; 4.386; 4.364; 4.455; 4.523\}}$ | $\frac{4.386}{4.568}$ | $= 0.960$ |
| $\Gamma_{82}$  | $\frac{4.364}{\max\{4.455; 4.568; 4.386; 4.341; 4.523; 4.364; 4.386; 4.364; 4.455; 4.523\}}$ | $\frac{4.364}{4.568}$ | $= 0.955$ |
| $\Gamma_{92}$  | $\frac{4.455}{\max\{4.455; 4.568; 4.386; 4.341; 4.523; 4.364; 4.386; 4.364; 4.455; 4.523\}}$ | $\frac{4.455}{4.568}$ | $= 0.975$ |
| $\Gamma_{102}$ | $\frac{4.523}{\max\{4.455; 4.568; 4.386; 4.341; 4.523; 4.364; 4.386; 4.364; 4.455; 4.523\}}$ | $\frac{4.523}{4.568}$ | $= 0.990$ |
| $\Gamma_{13}$  | $\frac{4.500}{\max\{4.500; 4.500; 4.386; 4.364; 4.432; 4.455; 4.477; 4.455; 4.409; 4.659\}}$ | $\frac{4.500}{4.659}$ | $= 0.966$ |
| $\Gamma_{23}$  | $\frac{4.500}{\max\{4.500; 4.500; 4.386; 4.364; 4.432; 4.455; 4.477; 4.455; 4.409; 4.659\}}$ | $\frac{4.500}{4.659}$ | $= 0.966$ |
| $\Gamma_{33}$  | $\frac{4.386}{\max\{4.500; 4.500; 4.386; 4.364; 4.432; 4.455; 4.477; 4.455; 4.409; 4.659\}}$ | $\frac{4.386}{4.659}$ | $= 0.941$ |
| $\Gamma_{43}$  | $\frac{4.364}{\max\{4.500; 4.500; 4.386; 4.364; 4.432; 4.455; 4.477; 4.455; 4.409; 4.659\}}$ | $\frac{4.364}{4.659}$ | $= 0.937$ |
| $\Gamma_{53}$  | $\frac{4.432}{\max\{4.500; 4.500; 4.386; 4.364; 4.432; 4.455; 4.477; 4.455; 4.409; 4.659\}}$ | $\frac{4.432}{4.659}$ | $= 0.951$ |
| $\Gamma_{63}$  | $\frac{4.455}{\max\{4.500; 4.500; 4.386; 4.364; 4.432; 4.455; 4.477; 4.455; 4.409; 4.659\}}$ | $\frac{4.455}{4.659}$ | $= 0.956$ |
| $\Gamma_{73}$  | $\frac{4.477}{\max\{4.500; 4.500; 4.386; 4.364; 4.432; 4.455; 4.477; 4.455; 4.409; 4.659\}}$ | $\frac{4.477}{4.659}$ | $= 0.961$ |
| $\Gamma_{83}$  | $\frac{4.455}{\max\{4.500; 4.500; 4.386; 4.364; 4.432; 4.455; 4.477; 4.455; 4.409; 4.659\}}$ | $\frac{4.455}{4.659}$ | $= 0.956$ |
| $\Gamma_{93}$  | $\frac{4.409}{\max\{4.500; 4.500; 4.386; 4.364; 4.432; 4.455; 4.477; 4.455; 4.409; 4.659\}}$ | $\frac{4.409}{4.659}$ | $= 0.946$ |
| $\Gamma_{103}$ | $\frac{4.659}{\max\{4.500; 4.500; 4.386; 4.364; 4.432; 4.455; 4.477; 4.455; 4.409; 4.659\}}$ | $\frac{4.659}{4.659}$ | $= 1.000$ |
| $\Gamma_{14}$  | $\frac{4.432}{\max\{4.432; 4.432; 4.341; 4.318; 4.409; 4.341; 4.455; 4.545; 4.409; 4.455\}}$ | $\frac{4.432}{4.545}$ | $= 0.975$ |
| $\Gamma_{24}$  | $\frac{4.432}{\max\{4.432; 4.432; 4.341; 4.318; 4.409; 4.341; 4.455; 4.545; 4.409; 4.455\}}$ | $\frac{4.432}{4.545}$ | $= 0.975$ |
| $\Gamma_{34}$  | $\frac{4.341}{\max\{4.432; 4.432; 4.341; 4.318; 4.409; 4.341; 4.455; 4.545; 4.409; 4.455\}}$ | $\frac{4.341}{4.545}$ | $= 0.955$ |
| $\Gamma_{44}$  | $\frac{4.318}{\max\{4.432; 4.432; 4.341; 4.318; 4.409; 4.341; 4.455; 4.545; 4.409; 4.455\}}$ | $\frac{4.318}{4.545}$ | $= 0.950$ |
| $\Gamma_{54}$  | $\frac{4.409}{\max\{4.432; 4.432; 4.341; 4.318; 4.409; 4.341; 4.455; 4.545; 4.409; 4.455\}}$ | $\frac{4.409}{4.545}$ | $= 0.970$ |
| $\Gamma_{64}$  | $\frac{4.341}{\max\{4.432; 4.432; 4.341; 4.318; 4.409; 4.341; 4.455; 4.545; 4.409; 4.455\}}$ | $\frac{4.341}{4.545}$ | $= 0.955$ |
| $\Gamma_{74}$  | $\frac{4.455}{\max\{4.432; 4.432; 4.341; 4.318; 4.409; 4.341; 4.455; 4.545; 4.409; 4.455\}}$ | $\frac{4.455}{4.545}$ | $= 0.980$ |
| $\Gamma_{84}$  | $\frac{4.545}{\max\{4.432; 4.432; 4.341; 4.318; 4.409; 4.341; 4.455; 4.545; 4.409; 4.455\}}$ | $\frac{4.545}{4.545}$ | $= 1.000$ |
| $\Gamma_{94}$  | $\frac{4.409}{\max\{4.432; 4.432; 4.341; 4.318; 4.409; 4.341; 4.455; 4.545; 4.409; 4.455\}}$ | $\frac{4.409}{4.545}$ | $= 0.970$ |

|           |   |  |   |                       |   |       |
|-----------|---|--|---|-----------------------|---|-------|
| $r_{104}$ | = | $\frac{4.455}{\max\{4.432; 4.432; 4.341; 4.318; 4.409; 4.341; 4.455; 4.545; 4.409; 4.455\}}$ | = | $\frac{4.455}{4.545}$ | = | 0.980 |
| $r_{15}$  | = | $\frac{4.386}{\max\{4.386; 4.500; 4.477; 4.364; 4.500; 4.477; 4.500; 4.545; 4.500; 4.477\}}$ | = | $\frac{4.386}{4.545}$ | = | 0.965 |
| $r_{25}$  | = | $\frac{4.500}{\max\{4.386; 4.500; 4.477; 4.364; 4.500; 4.477; 4.500; 4.545; 4.500; 4.477\}}$ | = | $\frac{4.500}{4.545}$ | = | 0.990 |
| $r_{35}$  | = | $\frac{4.477}{\max\{4.386; 4.500; 4.477; 4.364; 4.500; 4.477; 4.500; 4.545; 4.500; 4.477\}}$ | = | $\frac{4.477}{4.545}$ | = | 0.985 |
| $r_{45}$  | = | $\frac{4.364}{\max\{4.386; 4.500; 4.477; 4.364; 4.500; 4.477; 4.500; 4.545; 4.500; 4.477\}}$ | = | $\frac{4.364}{4.545}$ | = | 0.960 |
| $r_{55}$  | = | $\frac{4.500}{\max\{4.386; 4.500; 4.477; 4.364; 4.500; 4.477; 4.500; 4.545; 4.500; 4.477\}}$ | = | $\frac{4.500}{4.545}$ | = | 0.990 |
| $r_{65}$  | = | $\frac{4.477}{\max\{4.386; 4.500; 4.477; 4.364; 4.500; 4.477; 4.500; 4.545; 4.500; 4.477\}}$ | = | $\frac{4.477}{4.545}$ | = | 0.985 |
| $r_{75}$  | = | $\frac{4.500}{\max\{4.386; 4.500; 4.477; 4.364; 4.500; 4.477; 4.500; 4.545; 4.500; 4.477\}}$ | = | $\frac{4.500}{4.545}$ | = | 0.990 |
| $r_{85}$  | = | $\frac{4.545}{\max\{4.386; 4.500; 4.477; 4.364; 4.500; 4.477; 4.500; 4.545; 4.500; 4.477\}}$ | = | $\frac{4.545}{4.545}$ | = | 1.000 |
| $r_{95}$  | = | $\frac{4.500}{\max\{4.386; 4.500; 4.477; 4.364; 4.500; 4.477; 4.500; 4.545; 4.500; 4.477\}}$ | = | $\frac{4.500}{4.545}$ | = | 0.990 |
| $r_{105}$ | = | $\frac{4.477}{\max\{4.386; 4.500; 4.477; 4.364; 4.500; 4.477; 4.500; 4.545; 4.500; 4.477\}}$ | = | $\frac{4.477}{4.545}$ | = | 0.985 |

Based on the normalization results, then the conversion was carried out into matrix-R. The display of matrix-R can be seen in Fig. 3.

$$R = \begin{pmatrix} 0.960 & 0.975 & 0.966 & 0.975 & 0.965 \\ 0.975 & 1.000 & 0.966 & 0.975 & 0.990 \\ 0.955 & 0.960 & 0.941 & 0.955 & 0.985 \\ 0.975 & 0.950 & 0.937 & 0.950 & 0.960 \\ 0.975 & 0.990 & 0.951 & 0.970 & 0.990 \\ 1.000 & 0.955 & 0.956 & 0.955 & 0.985 \\ 0.980 & 0.960 & 0.961 & 0.980 & 0.990 \\ 0.980 & 0.955 & 0.956 & 1.000 & 1.000 \\ 0.995 & 0.975 & 0.946 & 0.970 & 0.990 \\ 0.960 & 0.990 & 1.000 & 0.980 & 0.985 \end{pmatrix}$$

Fig. 3. Matrix-R.

Based on the matrix-R and the weight from decision-makers shown in Table V, the ranking calculations can be performed. The formula used to calculate ranking [18] refers to equation (3).

$$V_i = \sum_{j=1}^n w_j r_{ij} \quad (3)$$

Notes:

- $V_i$  = rank for each alternative
- $w_j$  = weighted value of each criterion
- $r_{ij}$  = normalized performance rating score

The ranking calculating process can be explained as follows.

$$\begin{aligned} V_1 &= (0.30)(0.960) + (0.30)(0.975) + (0.30)(0.966) + \\ &\quad (0.30)(0.975) + (0.30)(0.965) = 1.4524 \\ V_2 &= (0.30)(0.975) + (0.30)(1.000) + (0.30)(0.966) + \\ &\quad (0.30)(0.975) + (0.30)(0.990) = 1.4719 \\ V_3 &= (0.30)(0.955) + (0.30)(0.960) + (0.30)(0.941) + \\ &\quad (0.30)(0.955) + (0.30)(0.985) = 1.4391 \\ V_4 &= (0.30)(0.975) + (0.30)(0.950) + (0.30)(0.937) + \\ &\quad (0.30)(0.950) + (0.30)(0.960) = 1.4317 \end{aligned}$$

$$\begin{aligned} V_5 &= (0.30)(0.975) + (0.30)(0.990) + (0.30)(0.951) + \\ &\quad (0.30)(0.970) + (0.30)(0.990) = 1.4631 \\ V_6 &= (0.30)(1.000) + (0.30)(0.955) + (0.30)(0.956) + \\ &\quad (0.30)(0.955) + (0.30)(0.985) = 1.4555 \\ V_7 &= (0.30)(0.980) + (0.30)(0.960) + (0.30)(0.961) + \\ &\quad (0.30)(0.980) + (0.30)(0.990) = 1.4614 \\ V_8 &= (0.30)(0.980) + (0.30)(0.955) + (0.30)(0.956) + \\ &\quad (0.30)(1.000) + (0.30)(1.000) = 1.4675 \\ V_9 &= (0.30)(0.995) + (0.30)(0.975) + (0.30)(0.946) + \\ &\quad (0.30)(0.970) + (0.30)(0.990) = 1.4630 \\ V_{10} &= (0.30)(0.960) + (0.30)(0.990) + (0.30)(1.000) + \\ &\quad (0.30)(0.980) + (0.30)(0.985) = 1.4747 \end{aligned}$$

Based on the ranking results, it can be determined the most dominant aspect recommendations in supporting the realization of positive moral improvement and student learning quality. The aspect referred to is C-10, namely the aspect of "it is always obey the school rules". This aspect was chosen because it had the highest compared to other aspects. The C-10 aspect is an aspect of the *Palemahan* component.

The dissemination activities that had been shown previously in Fig. 1 were carried out through two activities. The first activity was an online workshop on 11 materials related to the operation and management of *THK-ANEKA* and the *SAW-based Stake* model evaluation website. The second activity was assistance related to matters that were not clearly understood in the online workshop. It was discussed in-depth and directly through face to face at school.

Implementation of the *THK-ANEKA* and *SAW-based Stake* model evaluation website had been carried out well generally. The Evaluation website categorization had been classified as good and effective to determine appropriate and accurate recommendations. This recommendation was related to the

supporting aspects of increasing positive morale and student learning quality in computer learning at Vocational Schools of IT in Bali. It was reinforced from the effectiveness percentage results in the evaluation website implementing was 88.973%. When it is viewed from the effectiveness standard of the eleven's scale, it is classified in the good category because the percentage is in the range of 85% -94%.

The effectiveness percentage results were obtained from the respondent's assessment data on the website implementation by using a questionnaire containing 20 questions. Item-1 was about ease of website installation. Item-2 was about the website appearance. Item-3 was about the consistency of each layout form. Item-4 was about the suitability and accuracy of the login design. Item-5 was about the suitability and completeness of the features available on the main menu.

Item-6 was about the suitability and completeness from the features available on the input form of indicator and weight. Item-7 was about the suitability and completeness of the features available on the input form of evaluation aspect assessment data provided by the respondents. Item-8 was about the suitability and completeness of the features available on the evaluator data input form. Item-9 was about the suitability and completeness of the features available on the antecedent form located in the *description matrix*. Item-10 was about the suitability and completeness of the features available in the transaction form which was located in the *description matrix*. Item-11 was about the suitability and completeness of the features available in the form outcomes which were located in the *description matrix*.

Item-12 was about the suitability and completeness of the features available in the *judgment matrix* form had referred to the *Tri Hita Karana* and *ANEKA* aspects. Item-13 was about the suitability and completeness of the features available in the recommendation and decision form. Item-14 was about the suitability of evaluation aspects in the *accountability* section in the *description matrix* form. Item-15 was about the evaluation aspects suitability of the *nationalism* section in the *description matrix* form.

Item-16 was about the evaluation aspects suitability of the *public ethics* section in the *description matrix* form. Item-17 was about the evaluation aspects suitability of the *quality commitment* section in the *description matrix* form. Item-18 was about the evaluation aspects suitability of the *anti-corruption* section in the *description matrix* form. Item-19 was about features that make it easy to store data, edit, update, and delete. Item-20 was about the website accuracy in calculating the *SAW* method and showed the right recommendations.

This research had succeeded in being a solution to the limitations of Ihsan and Furnham's research [21]; Boitshwarelo, Reedy, and Billany's research [22]; Kyllonen and Kell's research [23]; Mariš's research [24]; and Elmahdi, Al-Hattami, and Fawzi's research [25]. The solution was the *Stake* model evaluation website implementation at Vocational Schools of IT in Bali. It was able to show an assessment of the affective domain through internalizing the *Tri Hita Karana* concept, cognitive and psychomotor assessments through internalizing the *ANEKA* concept. It was reinforced by the

research results of Divayana, Sudirtha, and Gading [33]. They showed that there was a *Countenance* evaluation model application design that was integrated with the *Tri Hita Karana* and *ANEKA* concept. It is used to measure the character aspects so the cognitive and psychomotor aspects of students in computer learning.

Another research result [34] that strengthens the position of this study is the research of Assielou *et al.* It showed that emotion (affective domain) can affect student performance (cognitive and psychomotor domains) in the learning process using *Intelligent Tutoring Systems*. The research conducted by Sökkhey and Okazaki [35] also strengthens the position of this study by showing the existence of a website-based decision support system. It was used to predict poor student performance in the learning process. The principle was the same with this research which also developed a website to evaluate student performance as a whole both from the moral side (affective domain) and from the learning quality side (cognitive and psychomotor aspects).

Although this research had succeeded in being a solution to the limitations found in the five previous studies, this research also has several limitations. The limitations of this research are: 1) The *THK-ANEKA* and *SAW* based-*Stake* model evaluation website has not been implemented at Vocational Schools of IT in all Indonesia regions; 2) This evaluation website has not been combined with robot technology so that the input activity indicators and evaluation weights are still done manually by evaluators or decision-makers.

#### IV. CONCLUSION

Generally, dissemination and implementation results of the *THK-ANEKA* and *SAW* based-*Stake* model evaluation website had been carried out well at Vocational Schools of IT in Bali Province. It was evident from the results of documentation in dissemination and implementation. The effectiveness percentage result of 88.973%, which is in the good category at the eleven's scale effectiveness standards indicated the success of evaluation website implementation. Likewise, the application simulation results of the *SAW* method in determining the dominant aspects of realizing positive moral improvement and student learning quality. Those had also proven the success of this evaluation website implementation. This research obstacle can be answered by doing the right work in the future. Some future work that can be done, included: 1) Dissemination and further implementation of evaluation website to several Vocational Schools of IT in western and eastern parts of Indonesia; 2) Development of evaluation website in the future is embedded in robotic technology so that the website will be more reliable in processing decision-making.

#### ACKNOWLEDGMENT

The authors express their sincere gratitude to the Directorate General of Research and Development, Ministry of Research and Technology of the Republic of Indonesia that had to provide the funding for this research. This research was able to be funded and completed on time based on the research contract No. 111/UN48.16/LT/2020.

REFERENCES

- [1] S.J. Hartati, N. Sayidah, and Muhajir, "The use of CIPP model for evaluation of computational algorithm learning program," IOP Conf. Series: Journal of Physics: Conf. Series, Vol. 1088, pp. 1–6, October 2018 [The 6th South East Asia Design Research International Conference (6th SEA-DRIC), Banda Aceh, Indonesia, p. 3, 2018].
- [2] D.T. Nkhosi, "The evaluation of a blended faculty development course using the CIPP framework," International Journal of Education and Development using Information and Communication Technology, Vol. 15, No. 1, pp. 245–254, 2019.
- [3] A. Hamid, T.M. Siregar, J. Purba, and B.A. Mukmin, "Evaluation of implementation of blended learning in Universitas Negeri Medan," Britain International of Linguistics, Arts and Education (BioLAE) Journal, Vol. 1, No. 2, pp. 224–231, 2019.
- [4] R. Donkin, and E. Askew, "An evaluation of formative (in-class) versus (e-learning) activities to benefit student learning outcomes in biomedical sciences," Journal of Biomedical Education, Vol. 2017, pp. 1–7, 2017.
- [5] D. Gunherani, W. Irawati, and A. Muhidin, "The evaluation of e-learning program at the University of Pamulang," Advances in Social Science, Education and Humanities Research, Vol. 335, pp. 710–716, 2019.
- [6] T.V. Thanabalan, S. Siraj, and N. Alias, "Evaluation of a digital story pedagogical module for the indigenous learners using the stake countenance model," Procedia-Social and Behavioral Sciences, Vol. 176, pp. 907–914, 2015.
- [7] R. Harjanti, Y. Supriyati, and W. Rahayu, "Evaluation of learning programs at elementary school level of 'Sekolah Alam Indonesia (SAL)': (evaluative research using countenance stake's model)," American Journal of Educational Research, Vol. 7, No. 2, pp. 125–132, 2019.
- [8] T. J. Gondikit, "The evaluation of Post PT3 program using stake's countenance model," Malaysian Journal of Social Sciences and Humanities, Vol. 3, No. 4, pp. 109–118, 2018.
- [9] I.P.M. Dewantara, "Stake evaluation model (countenance model) in learning process bahasa Indonesia at Ganesha university of educational," International Journal of Language and Literature, Vol. 1, No. 1, 19–29, 2017.
- [10] G. Fatima, M. Malik, A. Hussain Ch, and D.E. Nayab, "Antecedents of early childhood special education program: a stake's model perspective," Bulletin of Education and Research, Vol. 39, No. 1, pp. 275–290, 2017.
- [11] N. Komarasari, F. Dlis, and E. Utomo, "Implementation of the countenance stake model in evaluating the effectiveness of text-based Indonesian learning in junior high schools," East African Scholars Journal of Education, Humanities and Literature, Vol. 2, No. 2, pp. 52–55, 2019.
- [12] I.W. Sukarma, "Tri Hita Karana theoretical basic of moral Hindu," International Journal of Linguistics, Language and Culture (IJLLC), Vol. 2, No. 9, pp. 84–96, 2016.
- [13] T.G.R. Sukawati, "Establishing local wisdom values to develop sustainable competitiveness excellence," Journal of Management and Marketing Review, Vol. 2, No. 3, pp. 73–82, 2017.
- [14] I.G.A.A.O. Dewi, I.G.A.A.P. Dewi, K.T. Kustina, and G.D. Prena, "Culture of Tri Hita Karana on ease of use perception and use of accounting information system," International Journal of Social Sciences and Humanities, Vol. 2, No. 2, pp. 77–86, 2018.
- [15] M. Kamal, and J. Elim, "Implementation of project based learning model for anti corruption subject in fundamental training for BPKP's civil servant candidates of the millennials generation," Advances in Social Science, Education and Humanities Research, Vol. 262, pp. 114–122, 2018.
- [16] F.S. Hilyana, and M.M. Hakim, "Integrating character education on physics courses with schoology-based e-learning. Journal of Information Technology Education: Research, Vol. 17, pp. 577–593, 2018.
- [17] M. Muslihudin, Trisnawati, S. Mukodimah, W. Hashim, B. Ayshwarya, P.T. Nguyen, K. Shankar, S.K. Peteraitis, and A. Maselena, "Performance of SAW and WP method in determining the feasibility of motorcycle engineering workshop for competency test of vocational high school student," International Journal of Recent Technology and Engineering, Vol. 8, No. 2S2, 348–353, 2019.
- [18] N. Aminudin, M. Huda, A. Kilani, W. H.W. Embong, A.M. Mohamed, B. Basiron, S.S. Ihwani, S.S.M. Noor, K.A. Jasmi, J. Safar, N.L. Ivanova, A. Maselena, A. Triono, and Nungsiati, "Higher education selection using simple additive weighting," International Journal of Engineering & Technology, Vol. 7, No. 2.27, pp. 211–217, 2018.
- [19] T. Sagirani, M.G. Virawan, and V. Nurcahyawati, "Simple additive weighting method in the triage decision support system," International Journal of Scientific & Technology Research, Vol. 8, No. 12, pp. 3008–3012, 2019.
- [20] K.R. Zubaeti, A. Budiarto, and D. Maryono, "Simple additive weighting method in the development of a system assessing the feasibility of job training industry," Indonesian Journal of Informatics Education, Vol. 1, No. 2, pp. 17–28, 2017.
- [21] Z. Ihsan, and A. Furnham, "The new technologies in personality assessment: a review," Consulting Psychology Journal: Practice and Research, Vol. 70, No. 2, pp. 147–166, 2018.
- [22] B. Boitshwarelo, A.K. Reedy, and T. Billany, "Envisioning the use of online tests in assessing twenty-first century learning: a literature review," Research and Practice in Technology Enhanced Learning, Vol. 12, No. 16, pp. 1–16, 2017.
- [23] P.C. Kyllonen, and H. Kell, "Ability tests measure personality, personality tests measure ability: disentangling construct and method in evaluating the relationship between personality and ability," Journal of Intelligence, Vol. 6, No. 32, pp. 1–26, 2018.
- [24] L. Mariš, "The testing of the temperament and character inventory method in penitentiary environment," Transcom 2015, Žilina, Slovak Republic, pp. 1–5, June 2015.
- [25] I. Elmahdi, A.A. Hattami, and H. Fawzi, "Using technology for formative assessment to improve students' learning," TOJET: The Turkish Online Journal of Educational Technology, Vol. 17, No. 2, pp. 182–188, 2018.
- [26] B. Daniawan, "Evaluation of lecturer teaching performance using AHP and SAW methods," Bit-Tech, Vol. 1, No. 2, pp.30–39, 2018.
- [27] A. Said, and E. Syarif, "The development of online tutorial program design using problem-based learning in open distance learning system," Journal of Education and Practice, Vol. 7, No. 18, pp. 222–229, 2016.
- [28] S. T. Martaningsih, Soenarto, and E. Istiyono, "Evaluation Model of Career Counseling Program in Vocational High School", International Journal of Evaluation and Research in Education, Vol. 8, No. 2, pp. 318–329, 2019.
- [29] G. Setiadi, S. Joyoatmojo, Sajidan, and Soeharto, "The development of blended learning-based self-learning on classroom action research training material to improve teachers professionalism", International Journal of Education and Research, Vol. 4, No. 9, pp. 213–224, 2016.
- [30] Y. Maryansyah, "An analysis on readability of english reading texts for grade ix students at MTSN 2 Kota Bengkulu", Premise Journal, Vol. 5, No. 1, pp. 69–88, 2016.
- [31] F. Y. Ginting, "An analysis of students' ability in using punctuation marks in descriptive paragraph writing", Budapest International Research and Critics Institute-Journal, Vol. 1, No. 3, pp. 338–344, 2018.
- [32] A.A.G. Agung, I.G.P. Sudiarta, and D.G.H. Divayana, "The quality evaluation of school management model based on balinese local wisdom using weighted product calculation," Journal of Theoretical and Applied Information Technology, Vol. 96, No. 19, pp.6570-6579, 2018.
- [33] D.G.H. Divayana, I.G. Sudirtha, and I.K. Gading, "Application design of countenance evaluation based on Tri Hita Karana-Aneka for evaluating the students' computer capability and students' character," Cogent Psychology, Vol. 7, pp.1-18, 2020.
- [34] K. A. Assielou, C.T. Haba, B.T. Gooré, T.L. Kadjo, and K.D. Yao, "Emotional impact for predicting student performance in intelligent tutoring systems (ITS)," International Journal of Advanced Computer Science and Applications (IJACSA), Vol. 11, No. 7, pp. 219-225, 2020.
- [35] P. Sokkhey, and T. Okazaki, "Developing web-based support systems for predicting poor-performing students using educational data mining techniques," International Journal of Advanced Computer Science and Applications (IJACSA), Vol. 11, No. 7, pp. 23-32, 2020.

# Physically-Based Animation in Performing Accuracy Bouncing Simulation

Loh Ngiik Hoon<sup>1</sup>

Faculty of Applied and Creative Arts  
Universiti Malaysia Sarawak  
Kota Samarahan, Sarawak  
Malaysia

**Abstract**—This study investigates the use of physics formulas in achieving plausible bouncing simulation in animation. The need for physics animation was to produce visually believable animations adhering to the basic laws of physics. Based on the review, the creation of accurate timing simulation in bouncing dynamic was significantly difficult particularly in setting keyframes. It was comprehensible that setting the value of keyframes was unambiguous while specifying the timing for keyframes was a harder task and often time-consuming. The case study of bouncing balls' simulation was carried out in this research and the variables of mass, velocity, acceleration, force, and gravity are taking into consideration in the motion. However, the bouncing dynamic is a significant study in animation. It is often used and it shows many different aspects of animations, such as a falling object, walking, running, hopping, and juggling. Therefore, the physical framework was proposed in this study based on numerical simulations, as the real-time animation can be addressed for controlling the motion of bouncing dynamic object in animation. Animation based physics algorithm provided the animator the ability to control the realism of animation without setting the keyframe manually, to provide an extra layer of visually convincing simulation.

**Keywords**—Bouncing simulation; physics algorithm; physics animation; real time animation; animation

## I. INTRODUCTION

The concept of physically based in animation has been long established by Disney artists through “The Twelve Basic Principles of Animation”. The main purpose of the twelve basic principles was to produce more realistic animations adhering to the basic laws of physics [1]. In the context of 3D computer animation, realistic timing is extremely important to add a life-like quality to animate objects and give the animation some real-world authority. The proper timing is crucial to make the ideas readable. Consequently, animation artists carefully study the motion of the objects by adding quality and accuracy to generate realistic-looking animations. Thus, the concept of applying the laws of physics in animation has further gained importance to generate an accurate timing animation, and there is a need for consideration of physics motion in this field.

The concept of physics motion can be interpreted with the principle of bouncing ball simulation. The bouncing ball is the most basic and one of the most important animation exercises. Bouncing ball simulation causes the natural factors of motion, velocity, acceleration, mass, gravity, friction, elasticity, or

squash and stretch, and timing. Hence, this concept of a bouncing ball is often used in the animation as a reference because it integrates several fundamental concepts that animators apply to just about everything they animate. The author in [2] also explained, bouncing ball is a common model for numerous rhythmic tasks such as walking, running, hopping, and juggling, and it has been an extensive study which provides a theoretical basis for control of such rhythmic tasks in animation. Last but not least, realistic bouncing ball simulation showed the significance in the animation that a lot of physical measurements are required. Plus, timing an animation is often the most difficult part to set the spatial values of the keyframe in achieving realistic simulation. Most of the users are unable to imagine the timing and convey it using the provided interfaces. Therefore, the physics-based approach is a well-adapted concept to simulate believable animations. Based on the physical motion regarding numerical simulations framework, the animator is provided with the ability to control the realism motion of simulated object without setting the keyframe manually, by adding an extra layer of visually convincing animation. Hence, the different section of introduction, literature review, analysis physics motion of bouncing simulation, physics motion, real time dynamic bouncing ball, result and findings, implication, discussion, and conclusion are discussed in this paper.

## II. LITERATURE REVIEW

### A. Concept of Physically based Animation

Physically-based animation has emerged as a core area of computer graphics finding widespread application in the film and video game industries as well as in areas such as virtual surgery, virtual reality, and training simulations [3]. With the advance production technology, it allows designers to create animation by their own will with the greatest degree of freedom, but the products of the technology are not very good in terms of natural performance [4]. In order to fit into the current trend of fast increase of computer processing and user experience, physically-based animation is a well-adapted concept to simulate the realistic-looking animations with self-controllable performance. According to [5], the physics-based approach uses the law of physics to simulate motion and interaction with the environment. In his study, he pointed out a complete and effective system for animation should integrate key-framing and physics-based techniques. Key-framing allows objects or characters to perform unnatural tasks.



However, physics-based simulation models the object or character's interaction with the environment in a physical way and ensuring a realistic result. The author in [6] also described physical dynamics are based on two basic notions, which are material points and forces that induce movement. In the study, researchers explained clearly about these two notions that are included in their model and present a first simulation algorithm. The author in [7] also emphasized that physically-based models are well suited to simulate natural motion and flexible elastic objects. Moreover, [8] emphasized that physically-based animation tends to model only perfect worlds. Modeling a completely realistic physical world needs to talk in all external factors, including shifting winds, air humidity, different materials, and attraction forces from every object. All these factors will affect the moving objects in the real world. The study also showed that the existence of physics-based animation and alternative algorithms is a way to allow the animator to manipulate a rigid body during the simulation and have the computer to make the necessary adjustment to position and velocity. The results are presented in a good visual-oriented part and believability through several performances. From the findings, the researcher suggests that the animation developed with the reference to physics-based simulation to give better performance and even better visual results. The results also proved the inaccuracy will be minimal for a visual result, and thus the believability is retained by using the physically-based method.

### B. Bouncing Ball Animation

Based on [9], bouncing ball simulation is the standard animation test for all the beginners. It is one of the most important subjects to learn animation and is essential to know as an animator. From this simple simulation, principles of animation such as timing, squash and stretch, arcs, volume, and weight can be learned. The concept of the bouncing ball is shown in Fig. 1. The ball travels up and down through space while travels horizontally from left to right. With each bounce, the ball loses height because friction with the ground reduces the momentum of the initial force that set the ball in motion. If the ball is thrown upwards at an angle, it travels in a parabola. In addition, a ball bouncing on a hard surface that proceeds in a series of diminishing parabolas is caused by the energy lost on each bounce. Alternately, a cartoon character also bounces in much the same way as a ball, which is shown in Fig. 2.

According to [11], the bouncing ball simulation path has the similarity of the hop and jump of the character motion. Hence, this is to say that, the concept of bouncing ball simulation is critical to the animator as a reference for doing any kind of animation. To quote, [12] also stated that the movement is based on what happens in nature, however, simplified and exaggerated if necessary for dramatic effect. At the very least, the motion of an object should create based on the reality before it changed or modified to the motion we need. It means a realistic simulation of bouncing balls is essential as a guideline to the animators in animating their desired object motion.

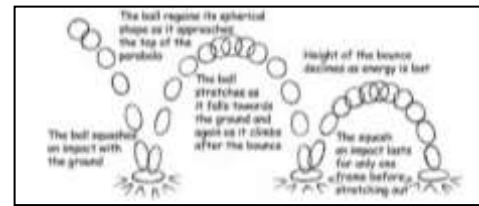


Fig. 1. The Concept of Bouncing Ball by [10].

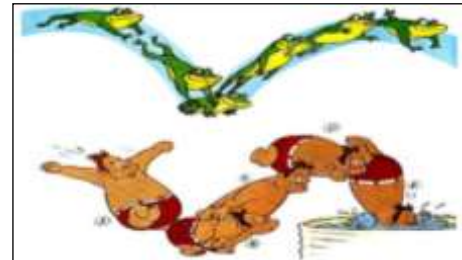


Fig. 2. Bouncing Ball and Character Motion Path by [11].

### III. ANALYSIS PHYSICS MOTION OF BOUNCING SIMULATION

Based on the review and information collected of physical measurement for bouncing balls simulation, it can be broken down into five distinct stages to analyze the details of physics calculation, to apply the accurate algorithm of bouncing simulation in animation.

In the first stage, when a ball is released from rest, it falls vertically downward under the influence of gravity ( $g$ ). The velocity ( $V$ ) points downward. The acceleration ( $a$ ) also points downward as shown in Fig. 3. When the object is in freefall, the magnitude of ( $a$ ) is equal to ( $g$ ), in the absence of air resistance. The acceleration due to gravity is  $g = 9.8 \text{ m/s}^2$  on earth.

In the second stage, the ball begins to make contact with the surface. The ball has slowed down. The velocity ( $V$ ) is still pointing downward. However, the ball has deformed sufficiently such that the acceleration ( $a$ ) is now pointing upward. This means that the ball has deformed enough until it is pushing against the surface with a force greater than its weight. As a result, the acceleration ( $a$ ) is pointing upward as shown in Fig. 4.

In the third stage, the ball has reached its maximum deformation when it barely touches the floor surface. As a result, the acceleration ( $a$ ) is still pointing upward, and the velocity ( $V$ ) is zero. The ball is preparing to re-bounce as shown in Fig. 5.

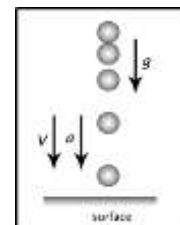


Fig. 3. The Ball in Vertically Downward Freefall.

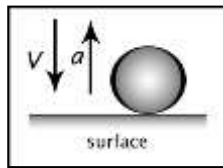


Fig. 4. The Ball Begins to Make Contact with the Surface.

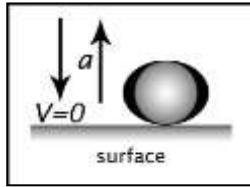


Fig. 5. The Ball Touches Floor Surface and Is Preparing to Re-Bounce.

In the fourth stage, the ball's velocity ( $V$ ) is increasing and pointing upward since the ball is now in the rebounding stage. As a result, the ball is less deformed than the previous stage, but is still deformed enough until it pushed against the surface with a force greater than its weight. This means that the acceleration ( $a$ ) is still pointing upward as shown in Fig. 6.

In the last stage, the ball has fully rebounded and has lifted off from the surface. The velocity ( $V$ ) is still pointing upward and the acceleration ( $a$ ) is pointing downward since the only force acting on the ball in this stage is gravity as shown in Fig. 7.

As shown in the above, a force would act between the ball and the floor. It changes the ball's velocity over some fairly small, but non-zero period. During this time, the ball would deform due to the force. The more rigid the ball material, the less the ball would deform, and the faster this collision would occur. In the limiting case, the ball is infinitely rigid, and cannot deform at all. To sum up, the collision is considered as occurring instantaneously in rigid body dynamics. It is the reason that the collision of the ball does not change the ball's shape in real motion. Thus, the deformation of the ball will be neglected in the factors of physics measurement.

From the above analysis, a free-fall motion under the influence of gravity can be described by the following basic motion equations:

$$\text{Force (F)} = \text{Mass (m)} \times \text{Acceleration (a)} \quad (1)$$

$$\text{Velocity (v)} = \text{Displacement (d)} / \text{Time (t)} \quad (2)$$

$$\text{Acceleration (a)} = \text{Velocity (v)} / \text{Time (t)} \quad (3)$$

The arrangement from the above basic equations can form four equations of linear motion which are important in solving the problems of linear motion,

$$V_1 = V_0 + gt \quad (4)$$

$$d = \frac{1}{2} (V_0 + V_1) t \quad (5)$$

$$d = V_0 t + \frac{1}{2} gt^2 \quad (6)$$

$$V_1^2 = V_0^2 + 2gd \quad (7)$$

, where  $g$ =the constant acceleration of gravity,  $V_0$ =initial velocity,  $V_1$ =final velocity,  $t$ =time and  $d$ =displacement.

From the above analysis, a realistic formula of bouncing ball's motion need to measure in terms of the force ( $F$ ), gravity ( $g$ ), distance=displacement ( $d$ ), time ( $t$ ), acceleration ( $a$ ), mass ( $m$ ), and speed=velocity ( $v$ ). All the physical quantity is dealing with the formula in physics law of motion. Hence, the basic physics' motion equation is very important to measure the movement of the bouncing ball.

Furthermore, physics analysis for the dynamics of the bouncing ball model was analyzed according to the graph shown in Fig. 8.

Fig. 8 shows the dynamics of squash ball bounces on the cement floors. The graph consists of the x-axis and y-axis which correspond to the variables of height,  $h$  (cm), and time,  $t$  (s) of the ball bounce. The graph represents that a ball dropped from an initial height,  $h_0$ . When it bounces, it loses energy so that the next bounce height,  $h_1$ , is smaller. It also happens in the bounce height,  $h_2$ . The model used here implied that the fraction of energy lost on each bounce. The bounce decreases by a constant amount until the ball comes to rest. Hence, the energy loss can be expressed in terms of the coefficient of restitution, COR, defined in the case of a rigid surface, is shown below:

$$\text{COR} = v_i / v_{f_0} = \sqrt{h_1 / h_0} \quad (8)$$

where,  $h_0$  is the initial height and  $h_1$  is the rebound height,  $v_{f_0}$  is the initial speed of the ball and  $v_i$  is the rebound speed.

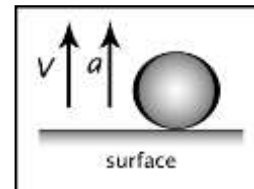


Fig. 6. The Ball is in the Rebounding Stage.

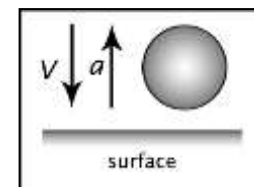


Fig. 7. The Ball Fully Rebounded and Lifted Off from the Surface.

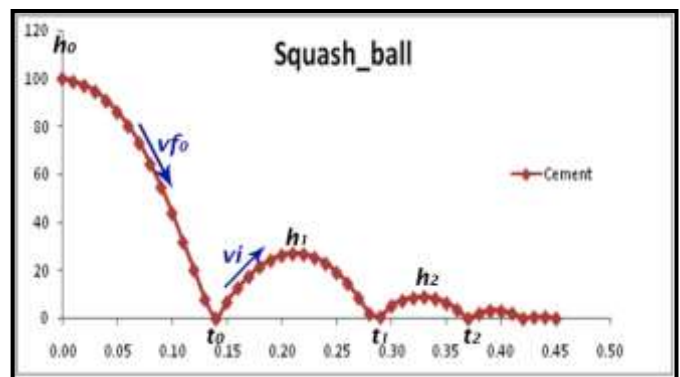


Fig. 8. The Dynamics of Bouncing Ball Model.

These parameters are important to measure the collision of the ball when it bounces and hits the floor surface. The formula of the coefficient of restitution is essential to measure the energy loss and how the ball bounces in the second height,  $h_1$ .

Besides, the motion in which a body is thrown or projected is called projectile motion. To calculate an accurate motion path of the bouncing ball simulation, projectile motion takes part in measuring the motion of the ball thrown. According to the projectile motion in Fig. 9, the physics analysis of a particle in free fall can be described by the equation from the projectile without air resistance,

$$V_i = -(COR)(V_{f_0}) \tag{9}$$

$$h_0 = V_i^2 / 2g \tag{10}$$

$$V_f = \sqrt{2gh_0} \tag{11}$$

$$t = \sqrt{2h_0/g} \tag{12}$$

$$T = 2(\sqrt{2h_0/g}) - 0.01 \tag{13}$$

$$h = h_0 - (0.5gt^2) \tag{14}$$

$$h_0 - (0.5 * g * (T - t)^2) \tag{15}$$

where,  $v_i$  = velocity reach top,  $V_{f_0}$ =initial velocity reach top,  $h_0$  =initial height,  $g$ =gravity,  $V_f$  =velocity reach top,  $t$ =time,  $T$ =total time,  $COR$ =coefficient of restitution.

The bouncing ball simulation addresses the needs of the projectile equation. Hence, the concept of the projectile equation is consistently as one of the parameters to formulate a realistic bouncing ball simulation.

In a nutshell, the summary of physics motion measurement for dynamic bouncing ball are presented in Table I.

From the above analysis, the basic physics motions, coefficient of restitution, and projectile equations are essential in measuring the dynamics of the bouncing ball's simulation. The algorithm carries on to the next phase in exploring the formula with the real-time bouncing simulation before applicable to the animation.

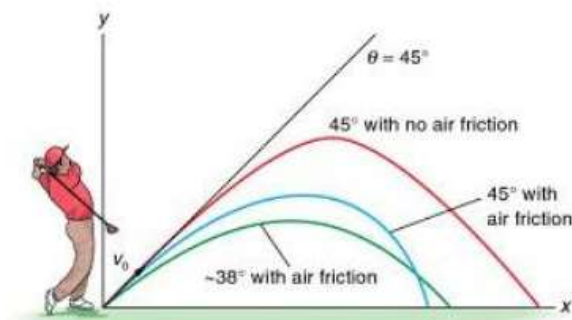


Fig. 9. The Effects of Air Resistance on Projectiles by [13].

TABLE I. SUMMARY OF PHYSICS MOTION MEASUREMENT FOR DYNAMIC BOUNCING BALL

| Dynamic of Bouncing Ball                 | Physics Motion                    | Formula Equations                                   |
|--|-----------------------------------|---|
| Energy Lost                              | Coefficient of Restitution        | $COR = \sqrt{h_1/h_0}$                              |
| Gravity                                  | Acceleration due Gravity on Earth | $g = 9.8 \text{ m/s}^2$                             |
| Velocity Reached Top                     | Projectile's Motion               | $V_i = -(COR)(V_{f_0})$                             |
| Initial height                           | Projectile's Motion               | $h_0 = V_i^2 / 2g$                                  |
| Velocity Reached Down                    | Projectile's Motion               | $V_f = \sqrt{2gh_0}$                                |
| Time of ball reached top at each bounce  | Projectile's Motion               | $t = \sqrt{2h_0/g}$                                 |
| Time of ball reached down at each bounce | Projectile's Motion               | $t = \sqrt{2h_0/g}$                                 |
| Total time reached of each half bounced  | Projectile's Motion               | $T = 2(\sqrt{2h_0/g}) - 0.01$                       |
| Height per Time                          | Projectile's Motion               | $h = h_0 - (0.5gt^2) / h_0 - (0.5 * g * (T - t)^2)$ |

#### IV. PHYSICS MOTION AND REAL TIME DYNAMIC BOUNCING BALL

In this phase, analysis and experiments are conducted on the physic formulas to define the algorithm to be considered and applied to simulate accurate bouncing animation. The summary of physics motion measurement for a dynamic bouncing ball in Table I was used as a guideline to obtain the bouncing ball's physic motion data. In the meantime, the physics factors of the ball's properties were considered in the parameters when doing the calculation. On the other hand, the real-time bouncing ball data will be gained based on the experiment, in which a squash ball will be used and dropped from the height of 1 meter to the hard surface floor, to determine the time and height of the ball will bounce. The experiments will be recorded via a video recorder, and a stopwatch will be used to time the number of seconds between the bounces of the ball with the height of the ball bounces. The height and times are recorded in the data chart in Microsoft Excel. The experiment will be repeated for a total of 10 trials for each ball to get accurate data, and the results would be shown in the graph. The progress measurement for getting real-time dynamic bouncing ball data is presented in Fig. 10.

As shown in Fig. 10, the videos recorded were rendered to frames by using Adobe Premium Pro. The frame rate of every video was rendered every 30 frames per second (30 fps). Every frame was calculated in the Autodesk AutoCAD to have an accurate measurement. All the data was then transferred to Microsoft Excel and the real motion of bouncing balls was analysed. For a real-time bouncing ball, a dynamic hysteresis curve was presented to show how energy is lost during and after the collision.

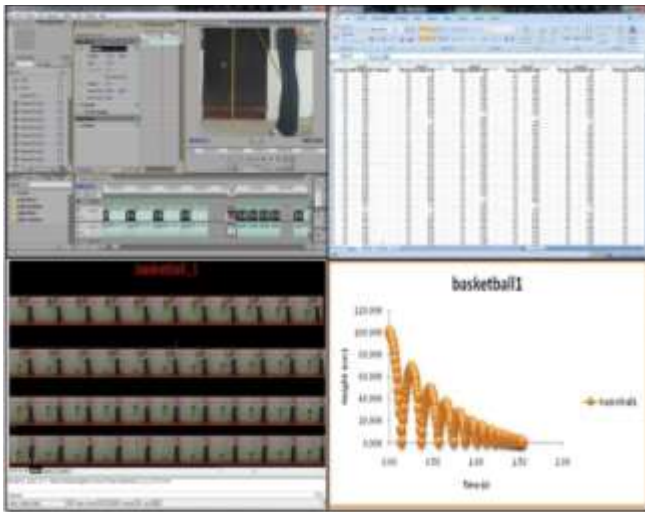


Fig. 10. Screen shot of Progress Measurement for Real Time Dynamic Bouncing Ball.

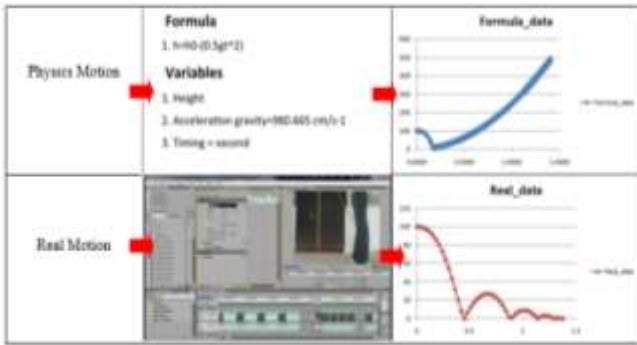


Fig. 11. Comparison between Both Physics and Real Time Data.

With the comparison between both physics and real-time data in Fig. 11, the accurate formulas and physics factors of the ball's properties will consider in the parameters and apply in animation.

First, the height of dynamic bouncing ball simulation is measured by using the formula of projectile motion as shown in Fig. 12. The formulas (14) are shown below,

$$\text{Height per Time, } h = h_0 - (0.5gt^2) \quad (16)$$

where,  $h$  = height,  $t$  = time and  $g$  = gravity.

The variables of the object's height, gravity acceleration, and time consider in this parameter. The result shows in the blue-colored graph in Fig. 12. The comparison of physics formula, the blue-colored graph has been done with real motion data, the red-colored graph in Fig. 12. The comparison shows that this physics parameter is not sufficient to simulate the realistic bouncing ball. Based on the graph, noticed that the height for the first ball bounce is correct but it is wrong starting from the second bounce. This parameter is just for measuring the correct height for the first bounce. Hence, the parameters of measuring the second bounce's height are needed.



Fig. 12. Measurement on the Height of Dynamic Bouncing Ball.

Next, two physics formulas of projectile motion were added to measure bounce's height of the bouncing ball. The formulas used are:

$$\text{Height, } h = h_0 - (0.5gt^2) \quad (17)$$

$$\text{Time reached of each half bounced, } T = 2(\sqrt{2h_0/g}) - 0.01 \quad (18)$$

$$\text{Height per Time, } h = h_0 - (0.5 * g * (T-t)^2) \quad (19)$$

where,  $h$  = height,  $t$  = time and  $g$  = gravity,  $h_0$  = initial height,  $T$  = total time.

The added formulas are shown in the red color in Fig. 13. The measurements in terms of the bounce's timing and height are separated. The timing is calculated when the highest point of the ball is achieved, in which the halfway between the two bounces. Thus, the accuracy of every ball bounce's height can be found. The comparison between physics motion in the blue-colored graph and real motion data in the red-colored graph is presented in Fig. 13. The comparison shows that this physics parameter is less successful to simulate the realistic bouncing ball. It does not consider the fraction of energy lost on each bounce. As a result, the height is the same in every ball bounce. Therefore, the coefficient bounce height is tested and added in the formula for the next experiment, which is shown in Fig. 14.

The formula of coefficient bounce height is added to measure the fraction of energy lost on each bounce. The formulas used are shown below:

$$\text{Height, } h = h_0 - (0.5gt^2) \quad (20)$$

$$\text{Time reached of each half bounced, } T = 2(\sqrt{2h_0/g}) - 0.01 \quad (21)$$

$$\text{Height per Time, } h = (h_0/h_1) - (0.5 * g * (T-t)^2) \quad (22)$$

where,  $h$  = height,  $t$  = time and  $g$  = gravity,  $h_0$  = initial height,  $T$  = total time.

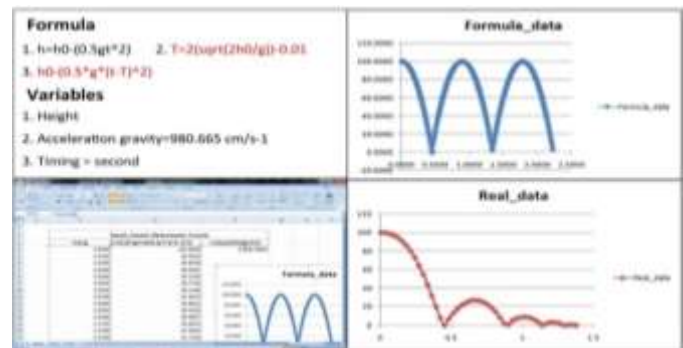


Fig. 13. Measurement on the Timing and Height of Dynamic Bouncing Ball.

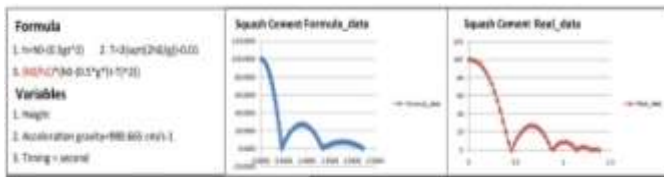


Fig. 14. Measurement on the Timing, Height and Coefficient Bounce Height of Dynamic Bouncing Ball.

The added formulas are shown in red color. The result is shown in the blue-coloured graph in Fig. 14. The comparison has been done between physics motion and real motion data in Fig. 14. From the comparison, the result shows that the graph of physics formula has a similar path compared to the real motion. But, it is still unsuccessful to measure the dynamics of the bouncing ball. It is because of the inaccuracy of the measurement on the bounce height in terms of the energy lost. Hence, the energy loss of dynamic bouncing ball is measured in the next experiment, by using the formula of the coefficient of restitution. The formulas are:

$$\text{Height, } h = h_0 - (0.5gt^2) \quad (23)$$

$$\text{Time reached of each half bounced, } T = 2(\sqrt{2h_0/g}) - 0.01 \quad (24)$$

$$\text{Height per Time, } h = h_0 - (0.5 * g * (T-t)^2) \quad (25)$$

$$\text{Velocity reached top, } V_i = - (COR) (Vf_0) \quad (26)$$

$$\text{Initial height, } h_0 = V_i^2 / 2g \quad (27)$$

$$\text{Velocity reached down, } Vf = \sqrt{2gh_0} \quad (28)$$

$$\text{Coefficient of Restitution, } COR = \sqrt{h_1/h_0} \quad (29)$$

where,  $h$  = height,  $t$  = time and  $g$  = gravity,  $h_0$  = initial height,  $h_1$  = second bounce height,  $T$  = total time,  $v_i$  = velocity reach top,  $Vf_0$  = initial velocity reach top,  $Vf$  = velocity reach top, and  $COR$  = coefficient of restitution.

Fig. 15 shows the measurement on the timing, height, and coefficient of restitution of dynamic bouncing ball. The formulas added are highlighted in red colour. This measurement has added the calculation on coefficient of restitution for energy loss. To calculate the coefficient of restitution for energy loss, the other formulas including velocity and initial height formula need to be added in to find the solution for it. Thus, these parameters have considered in the variables of the object's height, acceleration gravity, time and energy loss of bouncing ball. These elements which are important and essential in the real motion are highlighted in Fig. 15. The comparison between physics motion and real motion data indicated that these physics parameters can simulate the realistic bouncing ball. From the graph presented in Fig. 16, results proved that the formulas were validated to measure the bouncing ball's motion. Both motion paths achieved approximately the same result. However, there is a little difference between both bouncing ball's paths. It is because the real motion's data is gained via the average of 10 trial experiments for a type of bouncing ball. Next, the validated formulas from this experiment will be carried on to the next level, to transform them into algorithm and parameter in python script, in which applicable to animation software, Autodesk Maya to perform accurate bouncing simulation.

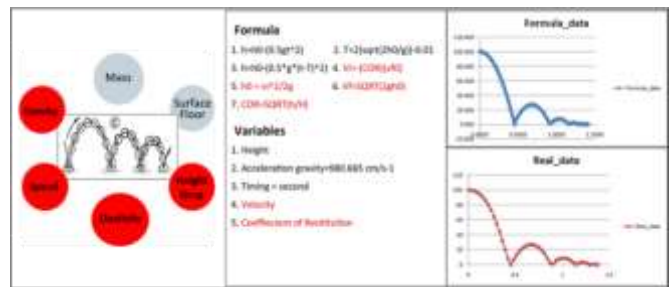


Fig. 15. Measurement on the Timing, Height and Coefficient of Restitution of Dynamic Bouncing Ball.

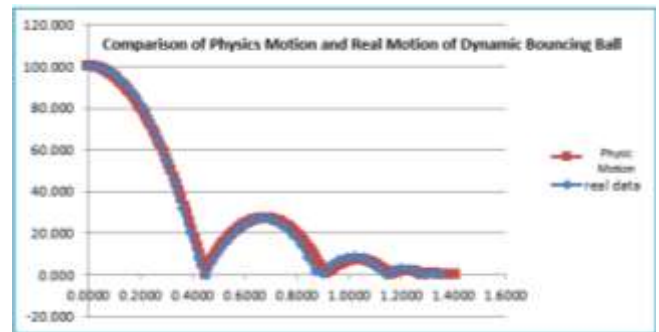


Fig. 16. Comparison of Physics Motion and Real Motion of Dynamic Bouncing Ball.

## V. RESULTS AND FINDINGS

Based on the research finding, the validated physics formulas can measure the accurate bouncing ball simulation. The summary of physics formulas and factors of ball's properties are considered as parameters will apply in animation are shown in Fig. 17.

According to the above analysis, the physics formulas are created regarding the physics factors of the object's properties. The factors of acceleration gravity, object's mass, floor surface, object's height, energy loss, speed, and time are included. All these factors play important roles in formulating a realistic dynamic of the bouncing ball. Apart from that, these parameters also involve the variables in terms of height, acceleration, timing, velocity, and coefficient of restitution for a dynamic bouncing ball. From research findings, it has been proven that the use of physics formulas would achieve realistic bouncing simulation. The physics formulas created from the above analysis are transformed into the Python scripting language as source code to drive the realistic motion of bouncing ball in Autodesk Maya, as shown in Fig. 18.

Furthermore, the comparisons are carried out between the real data via the previous experiments with the modelled bouncing ball simulation through the created system in Autodesk Maya. It is to validate the effectiveness of the applied physics formula to achieve accurate timing in bouncing ball animation. The comparisons are divided into two major elements include total time and the motion path for a bouncing ball's bounce. Every key-frame value will be calculated in the Autodesk Maya to gain accurate height and time of each modelled bouncing ball's motion. All the data are measured in Microsoft Excel and presented in the graphs. The comparison will be analysed in Table II.

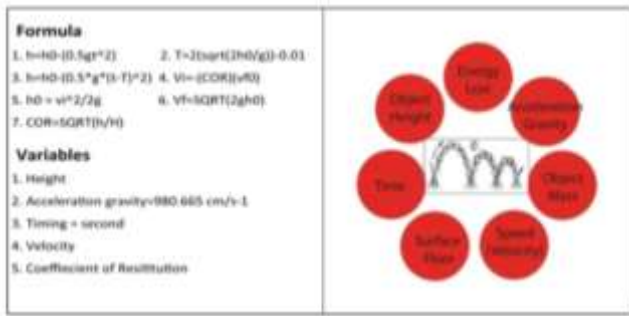


Fig. 17. Summary of Physics Formulas and Factors of Ball's Properties for Accurate Bouncing Ball Simulation.

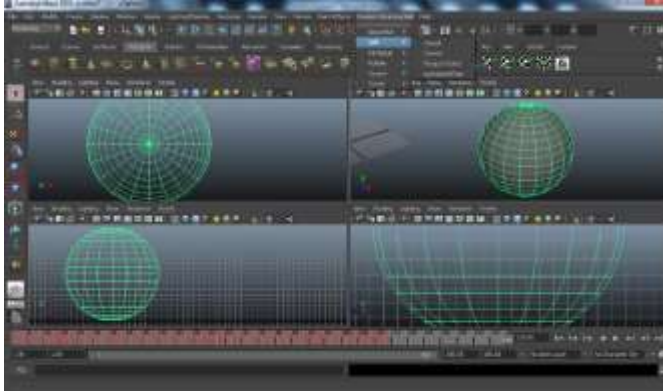


Fig. 18. Simulating Physics Based Bouncing Animation in Autodesk Maya.

TABLE II. COMPARISON BETWEEN MODELLED BOUNCING ANIMATION WITH REAL MOTION OF DYNAMIC BOUNCING BALL

| Ball   | Floor  | Bouncing Time (s)  |             | Comparison Graph |
|--------|--------|--------------------|-------------|------------------|
|        |        | Modelled Animation | Real Motion |                  |
| Squash | Cement | 1.10               | 1.11        |                  |

The graphs presented in Table II show the results of the comparison between the modeled and real simulation of bouncing balls. The graph consists of the x-axis and y-axis which correspond to the variables of height (cm) and time (s) of the ball bounce. The graph represents that a ball dropped from an initial height,  $h_0$ . When it bounces, it loses energy so that the next bounce height,  $h_1$ , is smaller. It also happens in the bounce height,  $h_2$ . The model used here implied that the fraction of energy lost on each bounce. The bounce decreases by a constant amount until the ball comes to rest. Based on the Table II, a squash ball was dropped from the height,  $h_0$  of 100cm, and reach the floor with the total time,  $t_0$  of approximately 0.45 second. When it bounces, it loses energy,

so the next bounce height,  $h_1$  is approximately 30cm, and reach the floor with total time,  $t_1$  of 0.85 second. It continue happens in the second bounce height,  $h_2$  with 10 cm and finally the ball comes to rest at the time of 1.11 second for real motion of bouncing ball. On the other hand, the simulation of bouncing ball animation comes to rest at the time of 1.10 second.

From Table II, the results prove that the physically-based animation system has been validated successfully to animate the bouncing ball's motion. The modeled bouncing balls' motion paths achieved approximately the same result as in real data. However, there are small differences between both the bouncing ball's motion path and the total time of bounce. The differences occur in the time of millisecond because the real motion's data is gained via an average of 10 trial experiments for a bouncing ball. However, the motion paths of modeled balls are still approximately matched in the range of real data.

As a result, it is proven that the physics-based formula can be applied to animation as algorithm animation. Besides, realistic simulation can be achieved through created physics-based animation. The result findings also showed that the use of physics formula would create realistic animation, and achieve both physics-based realism and user-specified expressive motion.

## VI. IMPLICATION AND DISCUSSION

This research studied physics algorithm formula and transformed into the Python scripting language as source code to drive accurate bouncing animation. The physically-based animation concept enables controllability on the dynamics bouncing ball in real-time without setting keyframes values in animation software. It offers controllability over a physical parameter using simple keyframe values, yet controllable animations. Besides, the application concept also allowed the animators to do the motion that would be very difficult to do by hand. Hence, it can help the animators to save more time and reduce workload while enabling faster output. The values in the keyframes tools can be changed and modified to achieve the user's desired motion. Animators are allowed to rearrange the object motion, guided by the pose specified at each keyframe. This function allows the user to further create a common model for numerous rhythmic tasks such as a falling object, walking, running, hopping, and juggling. This is because the bouncing ball simulation path has the similarity of the hop and jump of the character motion.

In addition, this research provided a physics object's factors to be considered in developing realistic bouncing ball animation, which is critical to the animator for doing any kind of animation. It is essential as a guideline to the animators in animating their desired object motion in real-time depending on their creativity. On the other hand, the algorithm can preserve the dynamics of motion and allowing animators to create motion libraries from a single input motion sequence. Once the original motion is fitted onto the model, the model can then be presented to the animator as a tool for generating a movement that meets the specifications of the given animation they are working on. The important factor of controlling real-time computer graphics is the combination of physics and

animation, which helps realistically imitate real world behavior and adding to the computer graphics' degree of realism.

However, the study of the formulas for bounce characteristics is limited for common types of hard surfaces, cement floor. The balls are bounced vertically without spin on a hard, horizontal surface. Based on the review, a strategic choice of the modeler can be studying first a simpler but related problem. In this case, this could be the free fall of an object without consideration of air resistance.

## VII. CONCLUSION

In conclusion, this study provided useful information on the development of physics-based animation. In computer animation, realistic timing is extremely important to add a life-like quality to animated objects, making their motions more interesting and able to convince the realism as well as convey the message to the audiences. However, timing of an animation is challenging to set the spatial values of the key-frame in achieving realistic simulation. It causes the burden of animation quality to rest entirely on the animator. The result findings of physics formulas in this paper provide animators a guideline that can be programmed into animation software as a library. Hence, animation such as objects falling, jumping, and bouncing animation can be produced easily without setting the keyframe value with real-time simulation.

Moreover, physically-based animation is an approach to make computer-generated virtual environments look realistic. It can be implemented in a production environment and interactive applications such as feature films, video games, or surgery simulation where the real-time simulation is required. The result of physics formulas in animation is useful to develop the games, especially the realism motion of algorithm is needed, such as shooter or bouncing game. The findings of the algorithm in this paper can be applied and transform into programming language to ensure the accurate movement of the object.

The development process of physically-based animation was explained in this paper which could be used as a reference in future research. Last but not least, future research can be expanded of the application of physics-based algorithm animation to a different type of ball's material and bounce surface. Besides, other objects or characters such as human or animal motion, human face expression, fluid, or fur object also can be applying with the physically-based concept. A physically-based animation approach helps the animators to

save more time and reduce the workload while enables faster outputs well as budget-saving. Moreover, the utilization of algorithm animation on computer animation can save the render space in computer's memories especially when rendering a thousand animal fur or human hair in animation. This approach should be generalized to deal with other objects, so that a wider variety of controllable physics-based realistic animation can be achieved in the future.

## ACKNOWLEDGMENT

This research was supported by the Research, Innovation & Enterprise Centre (RIEC), Universiti Malaysia Sarawak.

## REFERENCES

- [1] I. Gris, D. A. Rivera, & D. Novick, "Animation Guidelines for Believable Embodied Conversational Agent Gestures," In: R. Shumaker, S. Lackey, (eds) Virtual, Augmented and Mixed Reality. VAMR 2015.
- [2] A. K. Kamath, N. M. Singh, R. Pasumarthy, "Dynamics and Control of Bouncing Ball," International Conference on Intelligent Robotics and Manufacturing Automation, IRMA 2008, Vienna, Austria.
- [3] A. W. Bargteil, & T. Shinar, "An introduction to physics-based animation", SIGGRAPH '19: ACM SIGGRAPH 2019 Courses, Article No. 2, Pages 1-57, 2019.
- [4] Y. Yang, J. Yang, X. Zan, J. Huang, X. Zhang, "Research of Simulation in Character Animation Based on Physics Engine", International Journal of Digital Multimedia Broadcasting, vol. 2017, Article ID 4815932, 7 pages, 2017.
- [5] P. Faloutsos, "Physics-Based Animation and Control of Flexible Characters," Master's thesis, University of Toronto, 1995.
- [6] E. Promayon, P. Baconnier, & C. Puech, "Physically-Based Deformations Constrained in Displacements and Volume," In Eurographics '96, Blackwell Publishers, pp. C155-C162, 1996.
- [7] D. Terzopoulos, and K. Waters, "Techniques for Realistic Facial Modeling and Animation," Proc. Computer Animation '91, Geneva, Switzerland, Springer-Verlag, pp. 59-74, 1987.
- [8] S. Glimberg, & M. Engel, "Comparison of Ragdoll Method Physics-based Animation (Sid. 9 - 28)," Department of Computer Science, University of Copenhagen, 2007.
- [9] J. Lasseter, "Principles of Traditional Animation Applied to 3D Computer Animation," ACM Computer Graphics. Pixar, San Rafael, California, 1987.
- [10] C. Webster, "Animation: The Mechanics of Motion," London: Elsevier/Focal Press, 2005.
- [11] P. Blair, "Cartoon Animation," California: Walter T. Foster Publishing, 1994.
- [12] H. Whitaker, & J. Halas, "Timing for Animation," London: Elsevier/Focal Press, 2002.
- [13] D. L. James, & D. K. Pai, "ARTDEFO: Accurate Real Time Deformable Objects," University of British Columbia, 1999.

# Physiotherapy: Design and Implementation of a Wearable Sleeve using IMU Sensor and VR to Measure Elbow Range of Motion

Anzalna Narejo<sup>1</sup>, Attiya Baqai<sup>2</sup>, Neha Sikandar<sup>3</sup>, Absar Ali<sup>4</sup>, Sanam Narejo<sup>5</sup>

Department of Electronic Engineering, Mehran University of Engineering and Technology Jamshoro, Pakistan<sup>1, 2, 3, 4</sup>  
Department of Computer Systems Engineering, Mehran University of Engineering and Technology Jamshoro, Pakistan<sup>5</sup>

**Abstract**—Range of motion (RoM) is the measurement of angular movement of joints that defines the joints flexibility. It is crucial to measure RoM while performing musculo-skeletal diagnostics. The physiotherapy and the visits to hospitals can be very costly and demands a great deal of time; also most of the current digital instruments, used to measure RoM, are very expensive and hard to use. In this paper a digital wearable sleeve device is designed and tested which is cheap, time efficient and easy to use. The designed device is tested to be within 95 % agreement with Universal Goniometer (UG) when tested using Bland Altman Plots. Patients can take their measurements on their own and visualize results on their desktops or mobile phones. Patients also have graphical feedback, highlighting the extent of variation between their exercise performance and standard exercise. In addition to this; patients can also compare their current exercise from previous exercise using Kalmogorov-Simronov (K-S) test automatically. To make exercising more fun, we have developed 3D VR (Virtual Reality) gaming environment for elbow flexion, elbow supination and pronation and elbow extension exercises where patient can exercise in an interactive environment and visualize their progress side by side.

**Keywords**—Range of Motion (RoM); physiotherapy; Inertial Measurement Unit (IMU); Virtual Reality (VR)

## I. INTRODUCTION

Range of motion (RoM) is the measurement of angular movement of joints that defines the joints flexibility. These joint ranges of motion measurements are used to assess the patient's progress and to determine impairment ratings, when a patient is unable to return to his or her prior level of function. RoM measurement using manual methods is a time-consuming process [1]. Patients recovering from joint fractures and dislocations need to constantly visit physiotherapist and they need to maintain a regular exercise schedule. This can be very costly and demands a great deal of time also the current digital instruments, used to measure RoM, are very expensive and hard to use [2]. This paper discusses the design of a wearable digital device that measures RoM along with a user interface which helps patient exercise and monitors their progress in an easy and convenient manner without having to visit hospitals on regular basis. Patients can visualize their progress/results in graphical format. These graphs are created by statistically analyzing and comparing standard exercise and current exercise using Kalmogorov-Simronov (K-S) test. K-S test is a very efficient method to determine if two data sets are significantly different from each

other. Through a virtual reality application this system will provide patient an interactive environment to perform exercise; this will also isolate them from external disturbance. This is a complete system, through which patients will be able to perform physiotherapy at home without needing regular assistance from doctors. Compared to existing digital devices [3, 4] this system is cheaper and easy to use. The software applications instruct patients on how to perform exercises and display results on run time, making the system more user friendly than other ( purely hardware based) digital devices. This device will contribute to the health sector in a way that people would not need to ignore their physical health due to lack of time or money. It will also increase productivity of doctors and hospitals that will be able to give time to more patients as they would not need to perform complicated measurements manually.

The paper is organized as follows. Section 2 discusses the relevant work done in the field. Section 3 discusses the design of the digital sleeve, the experimental setup, and hardware and software implementation methodologies for the exercises of interest. Section 4 elaborates the results and discusses their analysis whereas Section 5 concludes the paper.

## II. RELATED WORKS

Visual estimation of RoM is the most preferred and commonly used method because of time constraints in occupational medical practice. This method can be used for all joints, and no additional equipment is needed but this method cannot be relied on to provide precise angular values and may give ambiguous results when assessor changes angle while visualizing RoM. Universal Goniometer (UG) is the most widely used device, to measure joints range of motion. It offers high accuracy, reliability in both inter-rater and intra-rater rehabilitation. It is, cheap, portable and noninvasive. But this instrument needs trained experts for evaluation. It is difficult to accurately position and requires clear visual estimation. In comparison to these manual methods digital photography provides slightly better accuracy, printable and savable records and ability to perform offsite measurement [1]. In this method accuracy is highly dependent on motion capture analysis and the whole camera setup is quite complex to arrange.

Laboratory based motion capture equipments are time consuming and costly thus they cannot be used for everyday



tasks. An easy to handle wearable device is a better alternative for day to day monitoring of joint RoM [5]. According to [6] most examiners adopt neutral-zero method while measuring RoM where the patient moves the distal segment away from a fixed starting position, around a certain axis of rotation.

Digital instruments provide a more objective and scientific assessment of patient's condition [7] but in order to be used in occupational practice these instruments must be easy to use and fast in application. A digital goniometer (Electro-goniometer) is an instrument which is similar to the electro potentiometer. It can also be used to measure joint range of motion, such that a change in joint position lead to the change in the resistance of the potentiometer and after some rectifications and calibrations, this resistance can be read as joint angle, the precision of electro goniometer is better than universal goniometer, but this precision depends upon the operator's ability to consistently place landmarks. This instrument is mostly used for clinical research. HALO [8] is another digital device which measures joint ROM. It is laser-guided digital goniometer. HALO uses lasers, magnetic system and accelerometer to guide alignment with anatomical landmarks. It has a digital display with memory feature. It is easy to handle but whenever, HALO is displaced off the horizontal plane the altered position of measuring system intermittently creates marked measurement errors, necessitating recurrent measurements.

Compared to these handheld instruments RoM measurement taken using Kinect are not tester dependent, thus Kinect provides better precision but, according to [9], Kinect cannot evaluate scapular motion and it cannot be used to measure neck and feet RoM.

The Inertial Measurement unit (IMU sensor) [10] contains an accelerometer, gyroscope, and magnetometer to calculate relative orientation in 3 dimensional space. IMU has already been used for RoM measurement of multiple joints in laboratory settings at the knee [11], cervical spine [12, 13] or shoulder. Although IMU provides high accuracy for continuous orientation estimation under ideal conditions, the presence of magnetometer may introduce orientation noise in sensor when placed in vicinity of ferromagnetic materials [14]. To eliminate this noise wooden chair, tables and couches were arranged [15]. These limitations are hard to avoid in occupational physiotherapy clinics however according to [16] considerable field accuracy has been achieved. The gyroscope component on the IMU tends to severely drift over time without the use of filtering or other navigation systems reference [17].

There is no rigid connection between IMU and human limb which may introduce soft tissue artifacts. In order to covert IMU orientation to anatomical angle the sensor requires anatomical calibration. To eliminate gyroscopic drift and to provide anatomical calibration a gradient descent algorithm can be used for IMU estimation [18]. The magnetometer and accelerometer data is used to estimate and compensate gyroscope error using quaternion representation of Euler's angles. This algorithm is designed for wearable inertial human motion tracking system in rehabilitation applications. The

algorithm has same level of accuracy as Kalman filter with static RMS error  $<0.8^\circ$  and dynamic RMS error  $<1.7^\circ$ .

In addition to these measurement techniques augmented reality is being frequently used in physiotherapy to provide patients an interactive virtual game-like environment in order to motivate them to exercise regularly. In [19] bio sensors along with Virtual Reality (VR) application and Kinect are used to provide Parkinson disease (PD) patient with an isolated rehabilitation environment to exercise while continuously monitoring their vital using biosensors.

This study discusses the design and implementation of a wearable digital device which is more time efficient than other manual devices and cheaper than other digital devices because it uses only one IMU sensor. This wearable device can be wirelessly connected to software applications which provide not only provide 3D visualization of arm movement but also provide analysis of patient performance and gives result on their improvement. Unlike other studies where VR application is used along with expensive hardware (like Kinect), this study uses Processing software to create a VR application which can be easily connected to the wearable device.

### III. MEASUREMENT OF ELBOW JOINT ROM USING WEARABLE DIGITAL DEVICE

#### A. Sleeve Design

In this study RoM measurements of elbow joint were taken using the digital wearable device. An IMU (Micro-processor Unit (MPU) 6050) sensor was sewn in a stretchable band as shown in Fig. 1. This band was worn over the distal segment (forearm in this case). The program was burned in node Microcontroller Unit (MCU) (ESP8266) module sewn in a band which was wrapped around the upper arm using Velcro strips, the subjects' angular motion was displayed on a small Organic Light-Emitting Diode (OLED) screen sewn adjacent to node MCU. As done in conventional examination, the subject started in a joint specific neutral starting-position and moved the adjacent, distal segment (forearm) to the end of range of motion, and the angle starting from stationary position to the fully extended distal segment was measured by both the UG and digital wearable device. Difference between the starting and current orientation of one IMU at the distal joint segment provided the angular measurement for IMU. The minimum and maximum angular values were noted and stored. The IMU was sewn in wearable band so examiner did not need to concern with manual handling of device.



Fig. 1. Wearable Sleeve.

## B. Study Design

In this study we have evaluated the validity, reliability and objectivity of this system by adopting the procedure in [20] but instead of adopting the visual estimation as comparison standard we have used UG for comparison. The validity is evaluated by comparing the measurement results to RoM measured using goniometer. The reliability is evaluated by analyzing repeatability under constant conditions. The objectivity is evaluated by analyzing the intra-rater agreements of measurements and examiner ratings.

Five subjects with age ranging from 23-80 years, weight 55-75 kg, and height 183-170cm volunteered to participate in the study. They gave their written consent to volunteer for this research study. Two of these volunteers were males and three were females. The subjects recruited were healthy without any known functional deficit. Functional deficit could occur due to a joint disease or recent joint injury. In order to ensure that subjects did not suffer from any musculo-skeletal complaints we checked patient's medical history of one month prior to examination. The joint examination using UG and the designed digital ROM measurement system was conducted by a physiotherapist. Examination rooms with a couch, digital wearable device and a laptop was provided to examiner.

## C. Examination Procedure

Fig. 2(a) shows how the stretchable band containing IMU was worn on forearm approximately one inch away from elbow joint for extension and flexion measurement. Fig. 2(b) shows position of band for supination and pronation measurement; the band was worn on forearm, one inch away from the wrist joint. Node MCU was connected to sensor. The examiner helped the patient in getting equipped. It took approximately one minute for subjects to get equipped and the removal took approximately thirty seconds. In order to avoid warming up or training effects each subject practiced the exercises three times. After warming up, each RoM measurement was repeated 5 times. The examiner measured RoM using UG while the IMU measured data simultaneously. All RoMs of elbow joint were examined actively and passively. For UG, measurement, the flexion and extension RoM were measured with shoulder in 90 degree forward flexion and forearm in maximum supination. The acromion and radial styloid process were landmarks for the goniometers' arms and the lateral epicondyle as the center of rotation. Supination and pronation were measured with a neutral position of the shoulder (0 degree shoulder abduction) and 90 degree of elbow flexion and a pencil placed over the distal palmar groove of the hand. The center of rotation for pronation and supination was over the head of the third metacarpal and the goniometers' arms were placed parallel to the humeral midline and parallel to the pencil. To achieve uniformity in participant's physical state all measurements of each participant were conducted in one day consecutively.

## D. Hardware

An IMU (MPU 6050) is used to measure angular data. The evaluation of sensor data was done using gradient descent algorithm which was implemented using arduino Integrated Development Environment (IDE). An ESP8266 module serves as controller and wifi module. The device was wirelessly

connected to processing three desktop applications, through which each subject's data for all exercises was automatically stored in Microsoft excel. The IMU was further calibrated by comparing anatomical angular readings against a geometric protector and recording change in IMU readings per 10 degrees. Then a simplified sensor change per degree algorithm was programmed in Node MCU and applied to sensor. IMU roll data was used for extension and flexion measurement while pitch data was used for supination and pronation assessment. A 0.96 inch programmable OLED display is used to display evaluated and calibrated RoM readings.

## E. Software

The arm movement can be visualized in Desktop, Mobile and VR application. The application is created using Processing 3 software. Fig. 3 shows few fundamental exercises for elbow physiotherapy. The software application designed provides game environments for these fundamental exercises of elbow extension, flexion, supination and pronation.

Fig. 4 shows 3D objects, being visualized in user interfaces which are created using Blender 3D. These 3D objects mimic user's arm movement in applications. Objects for forearm and upper arm are created separately because rotation is applied to the forearm object in Processing while the upper arm remains stationary.

Whenever user moves his arm, score is incremented; speed and time taken are also calculated. The arm under "instructional exercise" in Fig. 5 is used to instruct user on how to perform exercise. When user clicks the "Record data" button the exercise data is recorded in an excel file. If user wants to visualize recorded exercise, they can enter the name of exercise, click "load the excel file" and click "play recorded". User can switch between exercises by clicking "Exercise 1," "Exercise 2" or "Exercise 3".

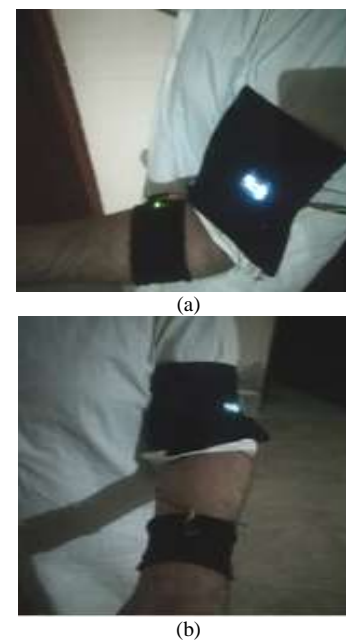


Fig. 2. Sensor Placement for Elbow (a) Extension and Flexion, (b) Supination and Pronation.

1) *Game 1*: Fig. 5 shows how Elbow flexion exercise is designed as a game 1, when the user moves the arm holding a dumbbell (or any weight in reality) from 130 degrees to 0 degrees the score is incremented. When the arm reaches 0 degrees a timer starts which counts to 30 seconds. When user completes this exercise five times the level is incremented.

2) *Game 2*: Elbow supination and pronation exercise is designed as game 2 as shown in Fig. 6, when user moves the arm holding drumstick (or any weight in reality) from 0 degrees to 90 or -90 degrees the score is incremented. When

the arm reaches 90 or -90 degrees the drumstick strikes the drum which changes color on every strike. When user completes this exercise five times the level is incremented.

3) *Game 3*: Fig. 7 shows how elbow extension exercise is designed as game 3, when user moves the arm holding box from 0 degrees to 140 degrees, while collecting the coins in box, the score is incremented. When user completes this exercise five times the level is incremented.

All these games can be played in VR mode in a separate VR application as shown in Fig. 8.

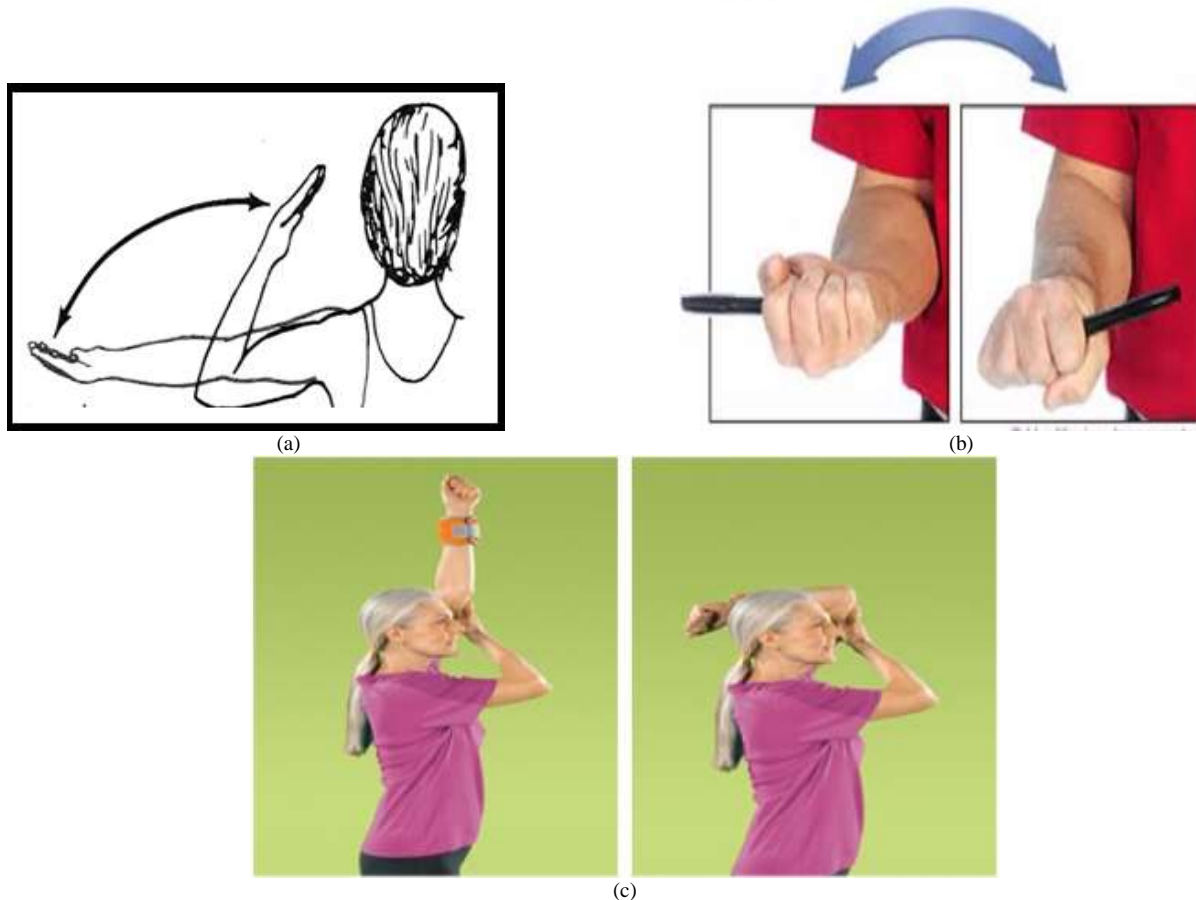


Fig. 3. Fundamental Exercise for (a) Elbow Flexion, (b) Elbow Supination and Pronation, (c) Elbow Extension.

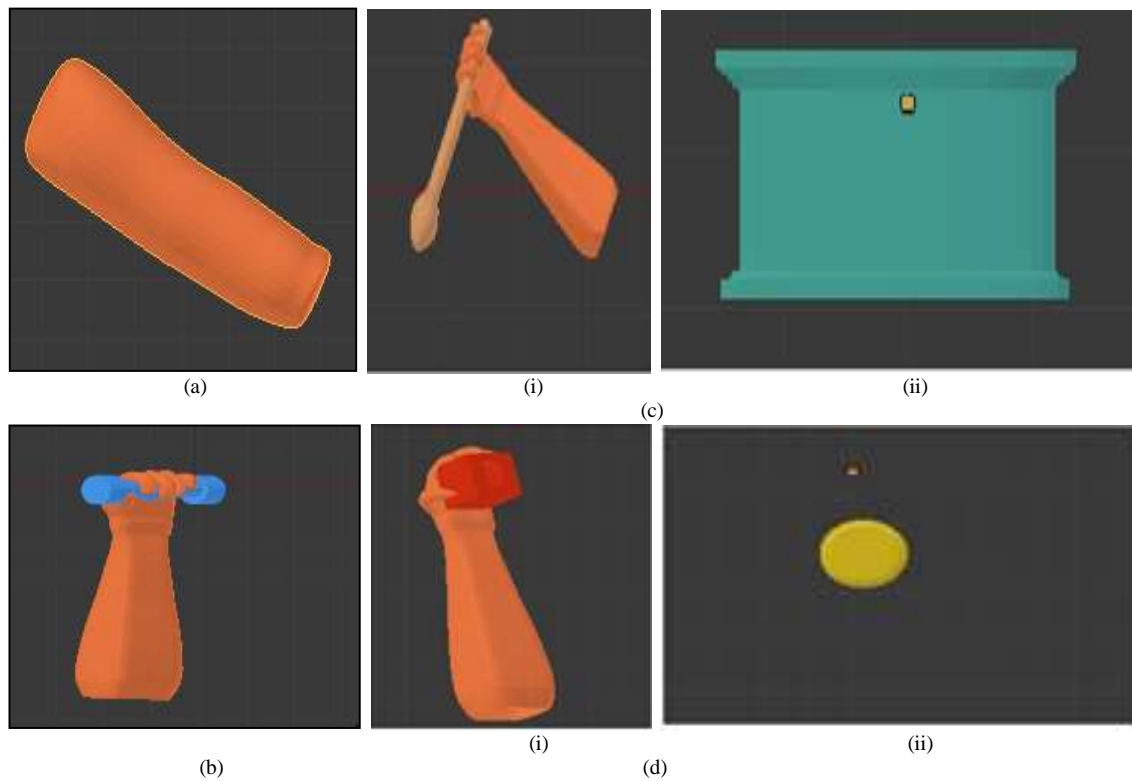


Fig. 4. 3D Objects for (a) Upper Arm, (b) Exercise 1 (Forearm Holding Dumbbell), (c) Exercise 2 ((i) Forearm Holding Drumstick (ii) Drum), (d) Exercise 3 ((i)Forearm Holding Box (ii) Coin).

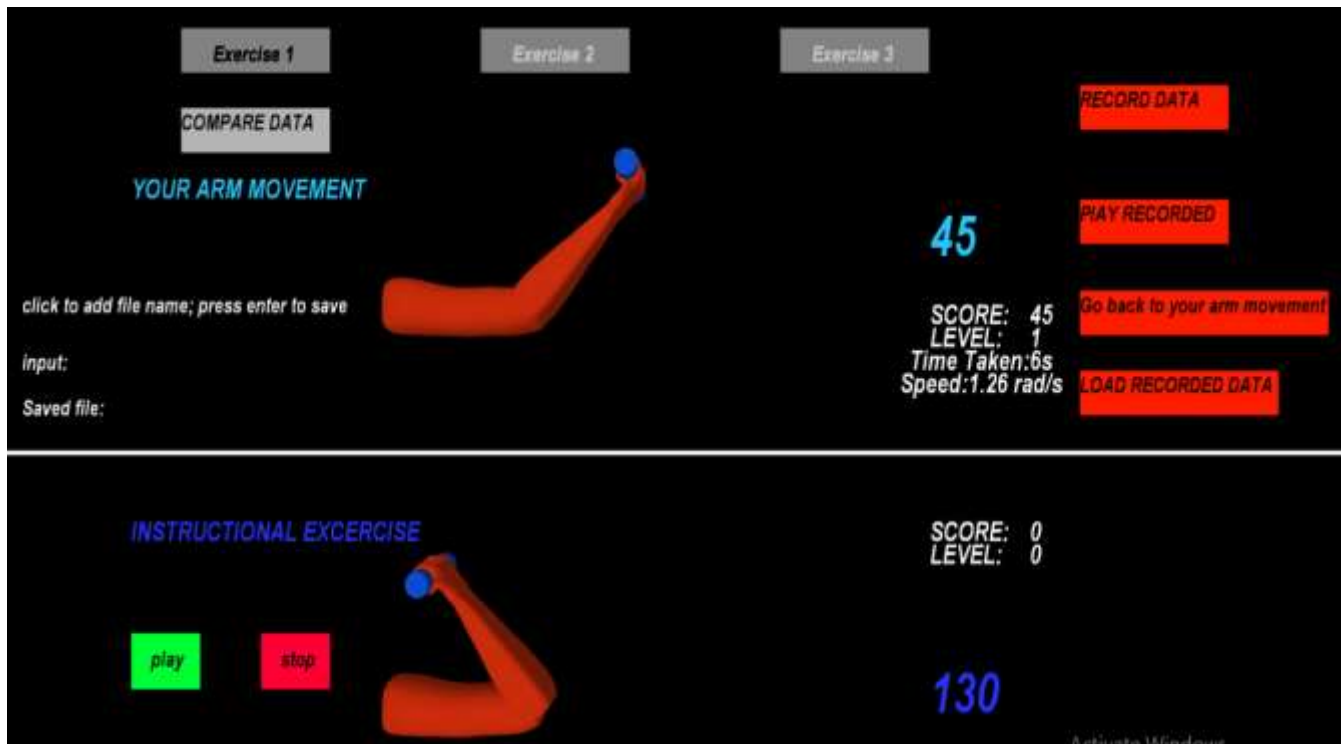


Fig. 5. Game 1: Elbow Flexion.

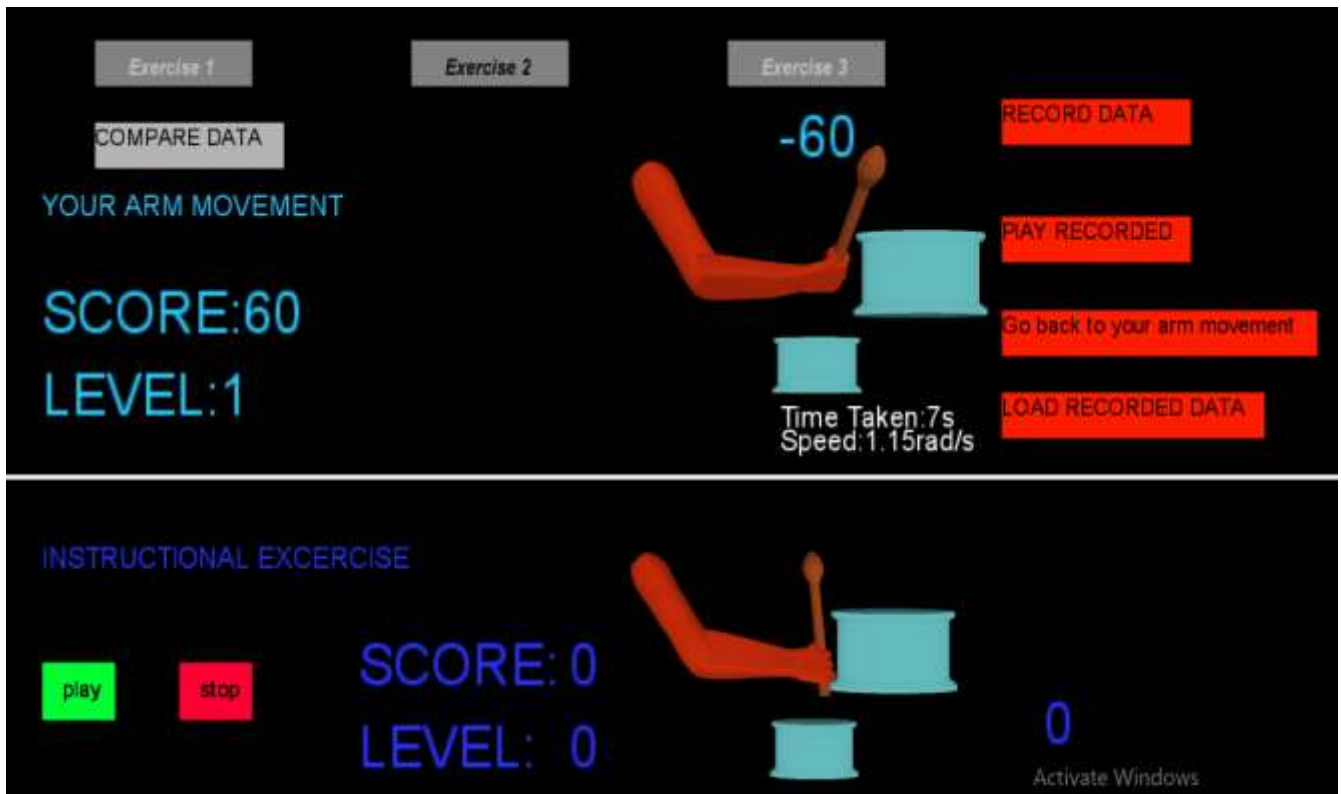


Fig. 6. Game 2: Elbow Supination and Pronation.

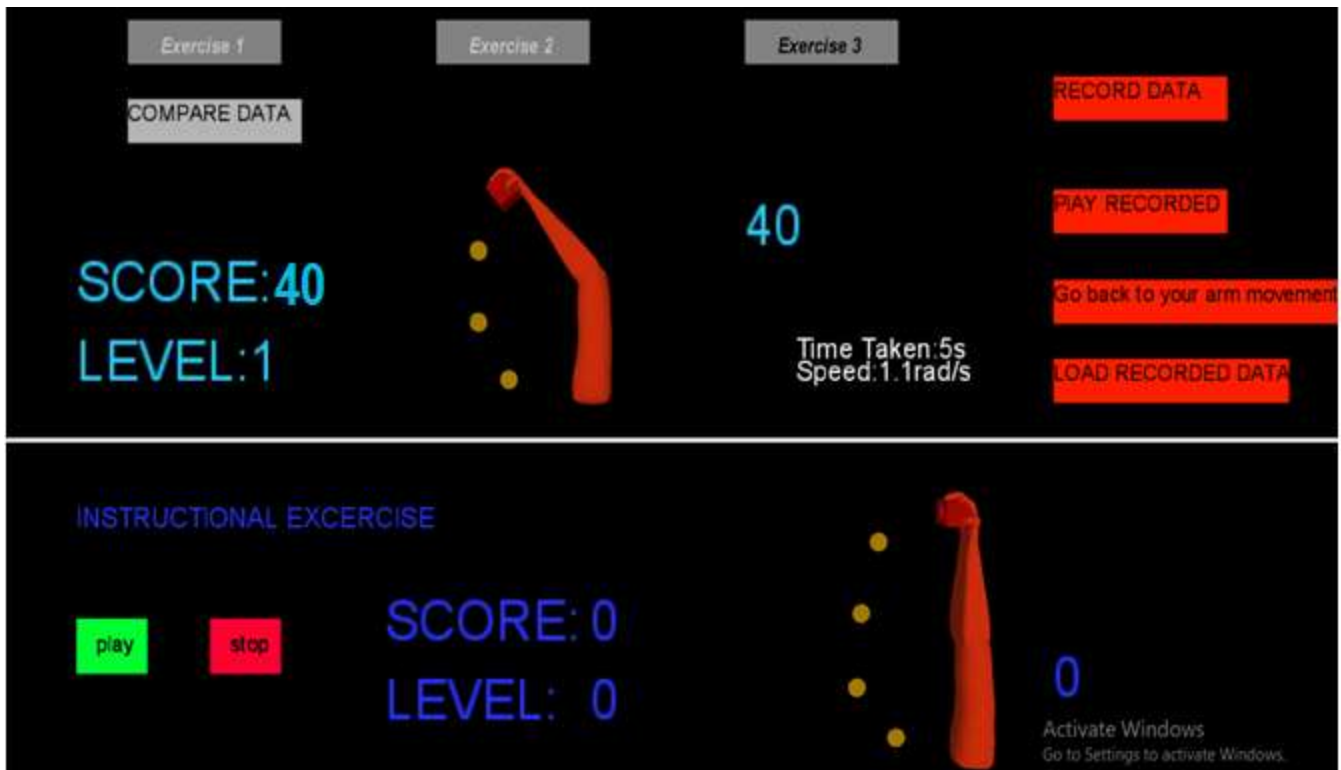


Fig. 7. Game 3: Elbow Extension.



Fig. 8. VR Application.

#### IV. IMPLEMENTATION RESULTS

##### A. Data Analysis

The IMU readings were compared with UG measurements. The data was stored in Microsoft excel. The validity of elbow RoM measurement using device was evaluated by comparing the UG and device measurement with standard elbow RoM values of healthy adult. For intra-rater reliability evaluation, under uniform condition, the mean standard deviation between five repetitions was considered.

##### B. Results

1) *Validity of measurement:* Table I shows the mean of elbow RoM values of all subjects as measured by the device and UG. Elbow extension, flexion and supination RoM values measured by UG and device lie in the range of expected RoM values of healthy adults but the measured RoM values of elbow pronation (66 degrees to 76 degrees) were below the expectations (90 degrees), the goniometer measurements were more close to the expected mark.

2) *Intra-rater repeatability:* The mean Standard Deviation (SD) between five repetitive measurements is given in Table II. The mean SD in device measurements for active RoM ranges from 1.35 degrees to 6.3 degrees and 1.1 degrees to 6 degrees for passive RoM measurement. This is approximately equal to the mean SD range of UG measurement 0.7 degrees to 6 degrees (for active RoM) and 1.1 degrees to 6.4 degrees (for passive RoM). But considering the whole range of measurements the device has higher SD as compared to UG.

3) *Agreement between UG and wearable device:* The agreement between validity and repeatability of both tools can be further be analysed using Bland Altman plot. The Bland Altman plot is used to assess agreement between two sets of measurements. Difference between measurements is plotted on y-axis while average of measurements is plotted on x-axis. Fig. 9 and Fig. 10 show Bland Altman plot for agreement between UG and the wearable device. The central Bias line indicate the average difference between measurement and if

the differences lie between Upper Limit of Agreement (LOA) and Lower Limit of Agreement then both measurement tools agree with each other, which means they can be used interchangeably with 95% probability of providing similar readings. The deviation between these from central bias indicates presence of systematic and random errors. These errors do not occur due to the measurement tools but they occur due to limitations of measurement procedure and human error. Any deviation outside the Upper and Lower LOA indicate lack of agreement between devices (less than 95% probability of similar readings) for that particular measurement.

4) *Testing on patient:* The purpose of this system is to make physiotherapy easy, so people who do not have any technical knowledge of physiotherapy can perform physiotherapeutic exercises in their homes while keeping a check on their performance. To help a common man figure out how well they performed, the current performance of patient is statistically compared against standard exercise data set (stored in application) using K-S test. Graph of both the standard and current data are plotted using cumulative frequency formula which relies on how many times a particular reading appears in data set. The frequency of appearance of these readings is added and plotted on y-axis while ROM readings are plotted on x-axis as shown in Fig. 11. The green curve represents standard performance while the blue curve represents patient's performance User just needs to click the "Compare Data" button (shown in Fig. 5) to make these comparisons. Whenever the button is clicked the software application provides a summary of patient's performance in comparison to their previous performance as shown in Fig. 12 and 13.

##### C. Discussion

The purpose of this study was to discuss the design and implementation of a digital wearable sleeve which can be used as a physiotherapeutic aid for doctors as well as patients who want to perform physiotherapeutic exercises from home. To

judge the accuracy and precision of device measurement, these measurements are compared with measurements taken using Universal Goniometer (UG). The accuracy and repeatability of measurement tool were found to be acceptable.

1) *Validity*: The accuracy of device is judged by comparing the mean of all device measurements for each elbow joint RoM with mean of all UG's measurements and with RoM values for healthy adults. The device gave valid joint angle measurement for all cases except elbow pronation where the device gives error of approximately +/-10 degrees because, when supination and pronation RoM is measured, the stretchable band containing IMU is worn close to wrist. Some of the subjects had thin arms which may have caused displacement of sensor during readings thus introducing error in measurements.

2) *Inter-rater repeatability*: The precision of device measurements for repetitive readings is judged through performing immediate repetitive measurements of each patient. It is necessary to maintain uniformity while taking consecutive measurements so the measurements do not deviate due to change in measurement procedure and deviations due to device limitations can be properly analysed. We tried to maintain uniformity by starting arm movement from horizontal position (neutral zero position) and keeping the rate of arm movement approximately constant for all readings. The SD between repetitive measurements of each patient is calculated. The extent of precision between measurements is given by mean SD for all readings. The less the deviation the better the repeatability of device. The device shows good repeatability for all measurement. The SD deviation for supination and pronation is higher than average but compared to UG's SD for those readings, it is acceptable.

3) *Agreement between UG and wearable device*: According to the Bland Altman plots shown in Fig. 9 and Fig. 10 the wearable device and UG can be used interchangeably and 95% of times they will give similar reading. Except for one measurement that lies below the lower LOA. This reading indicates that if device is used in place of UG for pronation measurement then the probability that it will give similar readings is less than 95%.

4) *Testing on patient*: We tested this device on a 27 year old female patient who suffered from elbow hemarthrosis. Initially she couldn't move her arm at all (which happens in hemarthrosis) so RoM for flexion and extension is 0 degree. After two weeks of injury, a physiotherapist prescribed patient to perform elbow flexion and extension exercises. According to the patient the device motivated her immensely to perform exercises regularly due to its game-like features and the graphical comparisons helped her keep a check on her performance without needing to contact a physiotherapist on regular basis. The patient was able flex their elbow upto 120 degrees within 2 weeks of exercising.

TABLE I. MEAN OF ROM READINGS OF SUBJECTS AND STANDARD DEVIATION (SD) BETWEEN MEANS

| ROM                             | Mean (SD) Device | Mean (SD) UG |
|---------------------------------|------------------|--------------|
| <b>Active ROM(in degrees)</b>   |                  |              |
| Extension L                     | 2(1)             | 0.5 (0.53)   |
| Extension R                     | 0.9(0.6)         | 0.45 (4.5)   |
| Flexion L                       | 135(2.5)         | 139.5 (3.8)  |
| Flexion R                       | 139.5(1.6)       | 138.7(4.5)   |
| Supination L                    | 95.3 (6.7)       | 97.4(2.3)    |
| Supination R                    | 99.75 (5)        | 98.6(7.5)    |
| Pronation L                     | 66.2(5.8)        | 77.2 (5.8)   |
| Pronation R                     | 66.85(9)         | 79.6 (3.5)   |
| <b>Passive ROM (in degrees)</b> |                  |              |
| Extension L                     | 4.8(1)           | 2.35(4.5)    |
| Extension R                     | 3.3(1)           | 2.95(1.65)   |
| Flexion L                       | 145.2(0.5)       | 143.95 (4.5) |
| Flexion R                       | 145.5 (1.5)      | 142.4(3.5)   |
| Supination L                    | 111.6 (8)        | 110.4(7)     |
| Supination R                    | 113.5(6)         | 110.8(9.14)  |
| Pronation L                     | 78(10.7)         | 89.6(12.7)   |
| Pronation R                     | 76.6(9.9)        | 90.4(9.3)    |

\*All examinations are done on elbow joint (L: left R: right)

TABLE II. MEAN OF STANDARD DEVIATION WITHIN FIVE IMMEDIATE EXAMINATION REPETITIONS OF ELBOW JOINT ROM EXAMINATION

| RoM                | Mean SD Device | Mean SD UG |
|--------------------|----------------|------------|
| <b>Active RoM</b>  |                |            |
| Extension L        | 1.35           | 0.7        |
| Extension R        | 1.16           | 1.09       |
| Flexion L          | 2.8            | 2.72       |
| Flexion R          | 2.7            | 2.56       |
| Supination L       | 4.3            | 2.5        |
| Supination R       | 5.7            | 6          |
| Pronation L        | 6.3            | 5          |
| Pronation R        | 3.51           | 3          |
| <b>Passive RoM</b> |                |            |
| Extension L        | 1.1            | 1.16       |
| Extension R        | 1.63           | 1.3        |
| Flexion L          | 2.13           | 2.16       |
| Flexion R          | 4.3            | 2.92       |
| Supination L       | 5.6            | 5          |
| Supination R       | 2.4            | 5.5        |
| Pronation L        | 6              | 5          |
| Pronation R        | 5.3            | 6.4        |

\*All examinations are done on elbow joint (L: left R: right).

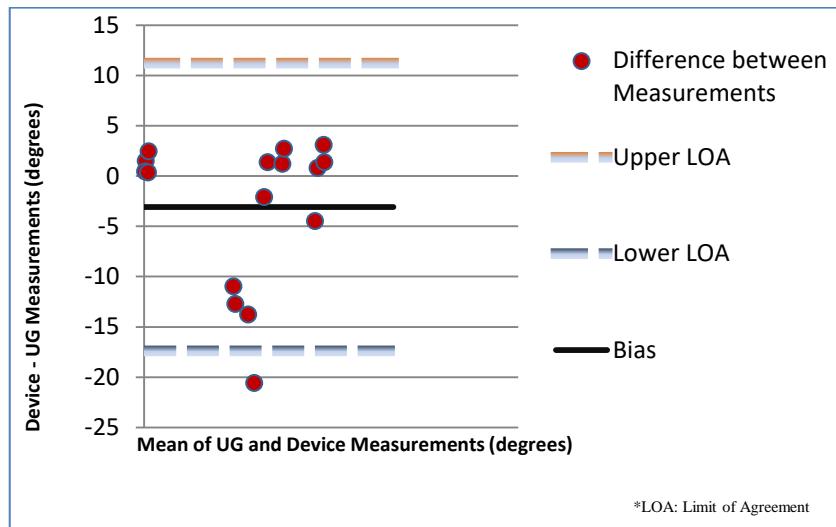


Fig. 9. Bland Altman Plot between difference of UG and Device Measurement and mean of Measurements.

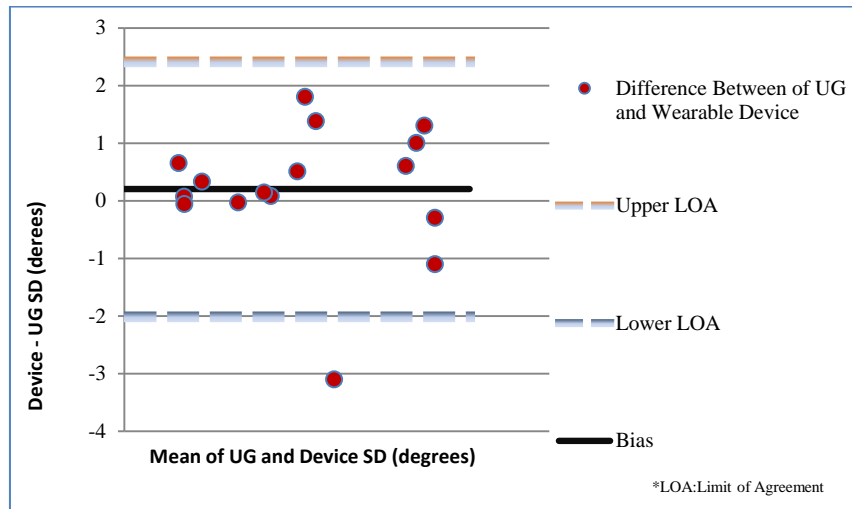


Fig. 10. Bland Altman Plot between difference of UG and Device SD and mean of SD.

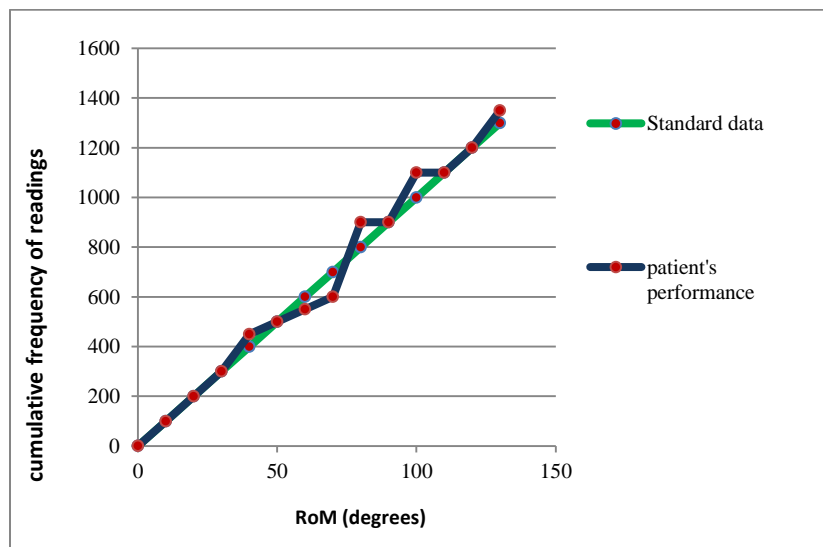


Fig. 11. Graphical Results of Elbow Flexion of a Healthy Subject.





Fig. 12. Comparison Results of Recovering Patient.

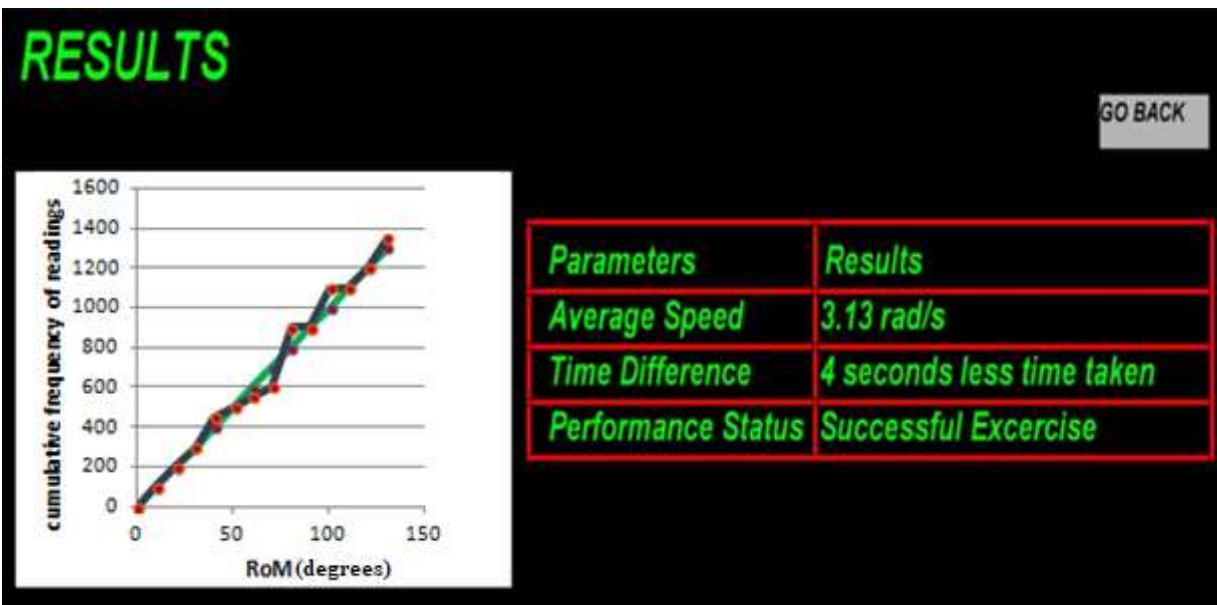


Fig. 13. Comparison Results of Recovered Patient.

## V. CONCLUSION

This wearable sleeve is a cheap, easy to use device which can be used in both clinical RoM measurements and for performing regular physiotherapeutic exercise. The device along with software application helps patient to maintain a regular exercise schedule. Compared to UG the accuracy and precision of device are good enough for performing clinical measurements. In future, the device could be further modified to cover wrist, hip, knee, shoulder, neck, spinal and feet joint. In addition to RoM measurement the features of device can be further extended to include muscle contraction measurements and muscle testing under load. By incorporating EMG (Electromyography) sensors in device, it can be used for detecting muscle pain. The software applications can be

further extended to cover muscular pain reduction exercises in addition to joint movement exercises.

There are few limitations in this study. First, only the elbow joint RoM was considered so no definite conclusion can be drawn about device performance for other joints. A second limitation is that the agreement between active and passive measurements relies highly on examiner so it is hard to draw a direct comparison, between active and passive measurements, which depends on device and not the examiner. Finally, the placement of IMU was also an issue. IMU was sewn in a stretchable band and depending on varying thickness of each patient's arm, it was hard to maintain a uniform position of sensor placement for all patients.

#### ACKNOWLEDGMENTS

The authors are thankful to the Department of Electronics Engineering, Mehran University of Engineering and Technology, Jamshoro, Sindh for facilitating this research. This research is funded by National Grassroots ICT Research Initiative (NGIRI) by Ignite Ministry of Science & Technology Pakistan and Dow University of Health & Sciences (DUHS) Karachi- in collaboration with Distinguished Innovations, Collaboration and Entrepreneurship (DICE-Health) Foundation USA.

#### REFERENCES

- [1] R. R. Russo et al., "Is digital photography an accurate and precise method for measuring range of motion of the shoulder and elbow?," *J. Orthop. Sci.*, vol. 23, no. 2, pp. 310–315, 2018.
- [2] J. Modest et al., "Self-measured wrist range of motion by wrist-injured and wrist-healthy study participants using a built-in iPhone feature as compared with a universal goniometer," *J. Hand Ther.*, pp. 1–7, 2018.
- [3] B. Huang et al., "Wearable stretch sensors for motion measurement of the wrist joint based on dielectric elastomers," *Sensors (Switzerland)*, vol. 17, no. 12, 2017.
- [4] B. Qi and S. Banerjee, "GonioSense: a wearable-based range of motion sensing and measurement system for body joints," *Proc. 22nd Annu. Int. Conf. Mob. Comput. Netw. - MobiCom '16*, pp. 441–442, 2016.
- [5] D. Saucier et al., "Closing the wearable gap—Part II: Sensor orientation and placement for foot and ankle joint kinematic measurements," *Sensors (Switzerland)*, vol. 19, no. 16, 2019.
- [6] Berryman Reese N, Bandy WD, Yates C. Joint range of motion and muscle length testing. St. louis, Missouri: Elsevier Health Sciences; 2009. <https://evolve.elsevier.com/cs/product/9781416058847>.
- [7] Line Blixt, Kari Nyheim Solbrække & Wenche Schrøder Bjorbækmo (2020) Becoming data. Patient perspectives on using an eTool in physiotherapy sessions, *Physiotherapy Theory and Practice*, DOI: 10.1080/09593985.2020.1790071.
- [8] S. Correll, J. Field, H. Hutchinson, G. Mickevicius, and A. Fitzsimmons, "Reliability and Validity of the Halo Digital Goniometer for Shoulder Range of Motion in Healthy Subjects," vol. 13, no. 4, pp. 707–714, 2018.
- [9] S. H. Lee et al., "Measurement of shoulder range of motion in patients with adhesive capsulitis using a Kinect," *PLoS One*, vol. 10, no. 6, 2015.
- [10] T.-L. Yoon, "Validity and Reliability of an Inertial Measurement Unit-Based 3D Angular Measurement of Shoulder Joint Motion," *J. Korean Phys. Ther.*, vol. 29, no. 3, pp. 145–151, 2017.
- [11] Favre Ja, Aissaoui Rab, Jolles BMc, de Guise JAB, Aminian Ka. Functional calibration procedure for 3d knee joint angle description using inertial sensors. *J Biomech.* 2008;42(14):2330–335.
- [12] Theobald PSa, Jones MDa, Williams JMb. Do inertial sensors represent a viable method to reliably measure cervical spine range of motion? *Man Ther.* 2012;17(1):92–6.
- [13] Jordan K, Dziedzic K, Jones PW, Ong BN, Dawes PT. The reliability of the three-dimensional fastrak measurement system in measuring cervical spine and shoulder range of motion in healthy subjects. *Rheumatology (Oxford)*. 2000;39(4):382–8.
- [14] Bachmann ER, Yun X, Peterson CW. An investigation of the effects of magnetic variations on inertial/magnetic orientation sensors. In: *IEEE International Conference on Robotics and Automation*. New Orleans, LA, USA: IEEE; 2004. p. 1115–1122. doi:10.1109/ROBOT.2004.1307974.
- [15] J. Ku, Y. J. Kang, J. Ku, and Y. J. Kang, "Novel Virtual Reality Application in Field of Neurorehabilitation Novel Virtual Reality Application in Field of Neurorehabilitation," vol. 11, no. 1, 2018.
- [16] Faber GS, Chang CC, Rizun P, Dennerlein JT. A novel method for assessing the 3-d orientation accuracy of inertial/magnetic sensors. *J Biomech.* 2013;46:2745–751.
- [17] F. Javier and S. Ortega, "Microelectromechanical Systems Inertial Measurement Unit As An Attitude And Heading Reference System," 2017.
- [18] S. O. H. Madgwick, A. J. L. Harrison and R. Vaidyanathan, "Estimation of IMU and MARG orientation using a gradient descent algorithm," 2011 IEEE International Conference on Rehabilitation Robotics, Zurich, 2011, pp. 1-7.
- [19] A. Baqai, K. Memon, and A. Rafique, "Interactive Physiotherapy: An Application Based on Virtual Reality and Bio-feedback," *Wirel. Pers. Commun.*, 2018.
- [20] C. Schiefer, T. Kraus, R. P. Ellegast, and E. Ochsmann, "A technical support tool for joint range of motion determination in functional diagnostics - An inter-rater study," *J. Occup. Med. Toxicol.*, vol. 10, no. 1, 2015.

# Outlier Detection using Nonparametric Depth-Based Techniques in Hydrology

Insia Hussain<sup>1</sup>

College of Computer Science and Information System  
Institute of Business Management, Karachi, Pakistan

**Abstract**—Several issues arise when extending the methods of outlier detection from a single dimension to a higher dimension. These issues include limited methods for visualization, marginal methods inadequacy, lacking a natural order and limitation in parametric modeling. The intension to overcome and address such limitations the nonparametric outlier identifier, based on depth functions, is introduced. These identifiers comprise of four threshold type outlyingness functions for outlier detection that are Mahalanobis distance, Tukey depth, spatial Mahalanobis depth, and projection depth. The object of the present research is the application of the proposed nonparametric technique in hydrology. The study is intended to be executed in two different frameworks that are multivariate hydrological data analysis and functional hydrological data analysis. The event of a flood is graphically represented by hydrograph whose components are used for computing flood characteristics that are peak( $p$ ) and volume( $v$ ). These characteristics are frequently employed for the various types of analysis in the multivariate study. Whereas, hydrograph is exhaustively employed in the analysis of functional data so that all the important information regarding flood event are not missed while analysis. The proposed technique in a multivariate framework is applied to the bivariate flood characteristics ( $p, v$ ) while in functional framework proposed approach is applied to the initial two scores of principal components denoted as ( $z_1, z_2$ ), since initial two principal components capture major variation of data employed for analysis.

**Keywords**—*Outlyingness functions; nonparametric techniques; flood characteristics; principal component scores; multivariate analysis; functional analysis*

## I. INTRODUCTION

The “outlier” observations in any data set is crucial to be detected and identified for nonparametric or parametric inferences. “Outliers” are the observations that are inconsistent or far from the majority of data points or within the chunk of data points with unusual behaviour. The presence of unusual observations in the data set acts as an outlier that can impact adversely the outcomes of estimation, inference, and testing procedures. Therefore, outliers are required to be identified and treated so that inferences are not violated due to unusual observations [1,2].

Outliers identified marginally suffer inadequacy of checking, in each coordinate, an outlier can find to be nonoutlying. Approaches that are algorithmic and take into account underlying geometry are required. A suitable function of outlyingness may be formulated with a threshold specified. A suitable choice can be Mahalanobis distance which is a

highly tractable function of outlyingness but constrained for having elliptical contours of symmetric outlyingness, even though whether the model under consideration is symmetric elliptically.

The author in [3] introduced a nonparametric technique which is based on functions of depth and orders the multidimensional data in center-outward. Higher depth represents higher centrality whereas lower depth greater outlyingness. One can associate with any depth function an equivalent function of outlyingness. For a suitable selection of depth function, actual geometrical structure and data shape are formed by equal outlyingness contours. In general, four different affine invariant functions of outlyingness were derived which are based on Mahalanobis distance outlyingness (MO), projection depth outlyingness (PO), halfspace or Tukey depth outlyingness (TO), and Spatial Mahalanobis outlyingness (SO). Related to these outlyingness functions the corresponding points are “outliers” having values of outlyingness exceed the constrained threshold of a particular function.

The nonparametric approaches introduced by [3] have been practiced by [4] and [5] in hydrology while [4] executed multivariate hydrological data analysis using two frequently employed flood characteristics; peak( $p$ ) & volume( $v$ ), for the identification of unusual observations i.e. outliers.

The author in [5] came up with groundbreaking research and extended the work of [4] by conducting functional hydrological data analysis. The nonparametric outlier identification technique was practiced in hydrology by [5] in such a way that the initial two scores of principal components were employed for the detection of outliers in a functional context. In multivariate analysis, employed flood characteristics are dependent and mutually correlated whereas scores of principal components employed in functional analysis are uncorrelated.

The execution of research in the functional framework follows the claim made by [5] that the characteristic of flood use in conducting the multivariate hydrological study are computed by subjective approach and do not encounter the complete series of employed data set, therefore, inferences of multivariate study suffer lack of authenticity. Hence it is crucial to conduct research in a functional framework so that authentic estimation regarding the associated risk of flood is obtained by incorporating complete phenomena produced through employed data series.

The objective carried by present research is the implementation of nonparametric techniques based on depth functions in both the context of a study that is a multivariate and functional framework using hydrological data of Kotri Barrage on Indus River in Pakistan.

## II. LITERATURE REVIEW

The methods going to be presented are based mainly on the statistical notion of depth functions. These functions provide convenient ranking tools for ordering data variables. Depth functions were initiatively practiced in hydrology by [6]. Several techniques of univariate analysis were extended to execute multivariate analysis developed through analogy. The variables that are dependent mutually affect the performance badly when analysing data component-wise, whereas moment-based techniques required the moment's existence.

Review in detail regarding techniques use for conducting classical multivariate analysis, it is referred to follow [7,8]. Techniques that are developed on the basis of depth, avoid the earlier drawbacks science depth functions are ordered using multivariate inward and outward ranking [9]. Indeed, techniques based on depth aren't component-wise, also, they are affine invariant and moment-free. Numerous techniques of outlier detection are enabled by ranking based on depth. The number of depth function formulas have been derived for executing the multivariate study. Depth region location inference considered by [3] is evaluated on sample space. Description of connection and general treatment related to multivariate quantile and centre ranked functions can be studied through [10,11]. For other inferential applications of depth see [12,13]. Numerous studies conducted in hydrology using various nonparametric approaches. The functions based on depth have been recently employed for the detection of outliers by [14,15]. According to [16], nonparametric models are suitable for capturing subtle aspects related to the frequency estimation of a flood. Flood inundation and flood damage were analysed using hydrologically distributed models through nonparametric techniques [17]. Similar other studies recently conducted in hydrology for outlier detection and risk estimation using nonparametric approaches are [18,19]. Characteristics of drought evaluation were assessed in a multivariate context implementing a nonparametric approach by [20-22]. Further research of [23] discussed data cleaning of water consumption and estimation of uncertainty regarding hydrologic modeling. Depth notion in regression was practiced and the performance of runoff model was evaluated, see work of [24-26]. Author in [27] used parametric and nonparametric multivariate approaches for designing rainfall framework whereas [28] applied rank-based nonparametric techniques to study trends of rainfall.

Multidimensional data is reduced by of analysis of functional principal component (AFPC) techniques to attain an easy approach for analyzing hydrological data. Notable work includes profile classification of streamflow, minimum indicators selection and functional data analysis application on streamflow are the studies executed on the basis of AFPC. Simulation of drought interval and drought changes were analysed by [29,30]. [31-33] studied rainfall variability modeling, pattern identification, and outlier detection. Other

relevant studies include work of [34-38], are also preferred for acquiring information about the useful application of AFPC in hydrology.

This paper is organized in such a way that the discussion regarding proposed methodologies is presented in Section 3. Section 4 provide description related to hydrological data employed for executing present research. Section 5 provides an application of the discussed methodology on employed hydrological data and obtained results are provided in Section 6 whereas Section 7 contain the conclusion drawn from the research.

## III. METHODOLOGY

This section contains methods for computing bivariate series of flood characteristic  $(p, v)$  and also bivariate series of principal component scores  $(z_1, z_2)$ . Both the computed series  $(p, v)$  and  $(z_1, z_2)$  are required for obtaining outliers in multivariate and functional context, respectively, using proposed threshold type nonparametric techniques which will also be discussed later in this section.

### A. Flood Characteristics

The flood peak  $(p)$  and volume  $(v)$  are the fundamental and most studied flood characteristics [39-41] and their computation based on the work of [41].

The bivariate series  $(p, v)$  are generated through hydrograph components using following formulas.

The flow peak series  $p_j$  is calculated as.

$$p_j = y_{hj}(t_k) \quad (1)$$

where  $y_{hj}(t_k)$  is the highest recorded observation of flow on a  $k$ th day in a  $j$ th year.

The flow volume series  $v_j$  is calculated as.

$$v_j = \sum_{t=SD_j}^{ED_j} y_j(t_k) - \frac{1}{2} (y_{ij}(t_k) + y_{fj}(t_k)) \quad (2)$$

where  $y_j(t_k)$  are the recorded observations of flow on a  $k$ th day in a  $j$ th year,  $y_{sj}(t_k)$  and  $y_{ej}(t_k)$  are the recorded observation of flow on starting  $(SD_j)$  and ending day  $(ED_j)$  respectively, in the  $k$ th year of flood time span.

### B. Analysis of Functional Principal Component

Analysis of principal component (APC) practices in a multivariate study for reducing the dimensionality through the computation of new variables which are the linear combination for original values so that the maximum of data variation could be captured. After the conversion of data as functions, analysis of functional principal component (AFPC) permits us to compute new functions so that special kind of variation for curve data could be revealed [5]. The AFPC method maximizes sample variance scores as orthonormal constraints. It divides the functional centred observations in orthogonal basis form and defined as follows.

Let functional observations be  $y_j(t), j = 1, \dots, n$  obtained after smoothing the discrete observations  $(y_j(t_1), \dots, y_j(t_T)), j = 1, \dots, n$ . By definition, the curve of mean is a same variation for most of the curves which can be

fixed by centering. Let  $(y_j^*(t) = y_j(t) - \bar{y}(t))_{j=1,\dots,n}$  be functional centered observations where  $\bar{y}(t)$  represents the function of mean for  $(y_1(t), \dots, y_n(t))$ . Now AFPC is applied to  $(y_j^*(t))_{j=1,\dots,n}$  for creating a set of small functions, known as harmonics which reveals the type of variation important for analysis. The first principal component  $(y_j^*(t))_{j=1,\dots,n}$  denoted as  $w_1(t)$  be a function so that variance regarding corresponding scores  $z_{j,1}$  of real value is as follows.

$$z_{j,1} = \int_C w_1(s) y_j^*(s) ds, j = 1, \dots, n \quad (3)$$

is maximized under  $\int_C w_1(s)^2 ds = 1$  constraint. The next  $w_l(t)$ ; a principal component computed by maximization of variance related to corresponding scores  $z_{j,l}$ :

$$z_{j,l} = \int_C w_l(s) y_j^*(s) ds, j = 1, \dots, n \quad (4)$$

under  $\int_C w_l(s) w_k(s) ds = 0, l \geq 2, l \neq k$  constraints.

### C. Detection of Outliers

The approaches for detection of outliers employed by [4] in the multivariate context was adapted by [5] in functional context; applying functions of outlyingness on the scores of initial two principal components. The purpose of this adaption is to create a comparison between multivariate and functional results.

Functions of outlyingness in a multivariate context were described and employed for detecting outliers. These functions have values ranging [0,1] interval. The outlyingness of a particular point is measured related to the whole sample. A value of outlyingness close to 1 shows high outlyingness, and a value close to 0 shows centrality. An observation is determined to be an outlier by defining a threshold i.e. the outlyingness value corresponds to an outlier must exceed their respective threshold values. Reference [3] introduced outlyingness functions which are based on the functions of depth, are going to be presented in the following section.

1) *Outlyingness functions:* A depth function is transformed to depth outlyingness for a F given distribution and  $x \in R^d$ . Reference [3] studied as follows.

a) *Half space*

$$O_{HO}(x, F) = 1 - 2HO(x, F) \quad (5)$$

b) *Mahalanobis*

$$O_{MO}(x, F) = d^2_{A(F)}(x, \mu(F)) / [1 + d^2_{A(F)}(x, \mu(F))] \quad (6)$$

c) *Projection*

$$O_{PO}(x, F) = PO(x, F) / [1 + PO(x, F)] \quad (7)$$

where  $HO(\cdot, F)$ ,  $d^2_{A(F)}(\cdot, \mu(F))$  and  $PO(\cdot, F)$  are given by [4], a location measure is  $\mu(F)$  and  $A(F)$  is non-singular measure of scale matrix.

Spatial

$$O_{SO}(x, F) = \|E(\text{Sign}(x - X))\| \quad (8)$$

d) *Spatial Mahalanobis*

$$O_{Ms}(x, F) = \left\| E[\text{Sign}(C^{-\frac{1}{2}}(x - X))] \right\| \quad (9)$$

where the Euclidean norm is  $\|\cdot\|$ ,  $F$ -distribution is  $X$  and the sign multidimensional function is  $\text{Sign}(\cdot)$  given by  $\text{Sign}(x) = x/\|x\|$  if  $x \neq 0$  and  $\text{Sign}(0) = 0$  also,  $C$  is any positive definite affine invariant  $d \times d$  symmetric matrix.

2) *Threshold:* An essential step in the detection of an outlier is the appropriate selection of the threshold. It relates to true positive and false positive rates.  $\alpha_n$  denoted for a false positive arbitrary rate which is defined as the proportion of misidentified nonoutliers as outliers. This constant relates closely to the  $\varepsilon_n$  true positive rate by which the theoretical proportion for real outliers are represented (also known as contaminants). Ideally,  $\alpha_n$  suppose to be smaller than  $\varepsilon_n$ . Reference [3] fixed the false outliers' ratio  $\delta = \alpha_n/\varepsilon_n$  and also used another coefficient  $\beta = \varepsilon_n\sqrt{n}$ , in order to define a threshold for the values of outlyingness as  $(1 - \alpha_n)$  quantile.

$$\rho_n = F_{O(x,F)}^{-1}(1 - \alpha_n) = F_{O(x,F)}^{-1}(1 - \delta\beta_n/\sqrt{n}) \quad (10)$$

where false positive rate  $\alpha_n$  is represented as  $\alpha_n = \delta\beta_n/\sqrt{n}$  and true positive rate  $\varepsilon_n$  represented as  $\varepsilon_n = n\varepsilon_n/n$ ; a number of true outliers are  $n\varepsilon_n$  and a number of observations are  $n$ , in such a way that  $\alpha_n < \varepsilon_n$ . For further calculations and applications, readers are referred to follow [4].

## IV. DATA DESCRIPTION

The major source of hydrological data is daily streamflow. The daily flow data series of the Kotri barrage are available from Sindh Irrigation department, Sindh Secretariat, Karachi, Pakistan.

A daily flow observations ( $m^3s^{-1}$ ) of Kotri barrage which is located between Jamshoro and Hyderabad in Sindh province on the Indus River, Pakistan. It has a discharge capacity of 875,000 cusecs (i.e. approximately  $24800 m^3s^{-1}$ ). Fig. 1 indicates the geographical location of the Kotri Barrage.

Some studies contain data of complete year while some consider section of a year having high flow observations. Hydrological data observations of the present study contain a duration of 6 months (i.e.  $T = 183$  days) per year spanning 1977 to 2017 (i.e.  $n=41$  years) since high flow period is observed during the months April to September, in Pakistan.

The series of observations are  $Y_j = (y_j(t_1), \dots, y_j(t_T))$ ,  $j = 1, \dots, n$ ,  $k = 1, \dots, T$ , where  $n=41$  years,  $T = 183$  days and  $y_j(t_k)$  is the recorded flow observation on  $t_k$  day in the  $j$ th year. Before any computation is performed the streamflow observations which are recorded on measurement scale in cusec (a volume flow rate) are required to be converted into cubic meter per second ( $m^3s^{-1}$ ).

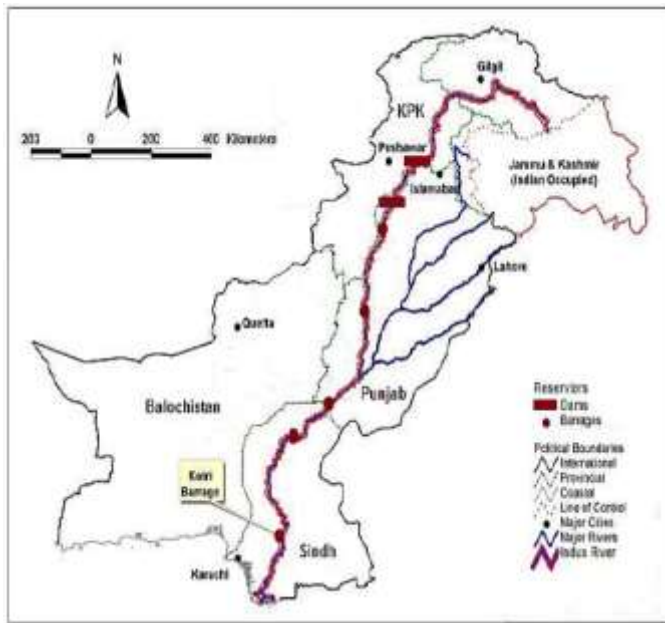


Fig. 1. Geographical Location of Kotri Barrage.

V. APPLICATION

The two most studied and examined characteristics of the flood that is peak (p) and volume (v) are focused here. The series of bivariate (p,v) are computed by using (1) and (2) and results are displayed in Table I.

According to [4], an approach developed by [3] are based on the function of depth outlyingness and the threshold corresponded. The four functions of depth outlyingness are evaluated for the (p,v) series of bivariate observation i.e., Mahalanobis (MO), Projection (PO), Spatial (SO) and Tukey (TO). The values of depth outlyingness correspond to each (p,v) observation for years 1977-2017 are reported in the last four columns of Table I. The thresholds correspond to each outlyingness functions are computed by selecting 15% false outlier ratio and the number of true outliers as 5, this selection is similar to the choices made by [4] in such a way that the outlyingness value corresponds to an outlier must exceed their respective threshold values.

Hence, 98% quantile is a corresponding threshold for the values of outlyingness. The computed values of the threshold for MO, PO, SO & TO are 0.9412, 0.9040, 0.9719, and 0.9444, respectively. The values of threshold approximately remain constant if the number of true outliers is considered greater than 5 with changed false outlier ratio i.e. 5%, 10% and 20%. The detected outliers correspond to MO, PO, SO & TO with respect to their respective threshold values are graphically displayed by Fig. 2.

Reference [5] employed the procedure for detecting outliers which are based on the function of depth outlyingness and the threshold corresponded. As discussed earlier and also practiced in preceding section, four functions of depth outlyingness are evaluated for the series of the bivariate score (z<sub>1</sub>, z<sub>2</sub>) i.e., Mahalanobis (MO), Projection (PO), Spatial (SO) and Tukey (TO).

TABLE I. MULTIVARIATE RESULTS FOR FLOOD PEAK AND VOLUME

| Year | Peak  | Volume | MO     | PO     | SO     | TO     |
|------|-------|--------|--------|--------|--------|--------|
| 1977 | 7490  | 248765 | 0.0979 | 0.5424 | 0.4134 | 0.4634 |
| 1978 | 15747 | 249063 | 0.8782 | 0.8631 | 0.4183 | 0.9512 |
| 1979 | 7342  | 305373 | 0.4843 | 0.7099 | 0.6352 | 0.7561 |
| 1980 | 5776  | 170479 | 0.0852 | 0.2978 | 0.0255 | 0.2195 |
| 1981 | 7149  | 246426 | 0.1473 | 0.5673 | 0.3586 | 0.5610 |
| 1982 | 5560  | 129340 | 0.1783 | 0.4059 | 0.2671 | 0.3171 |
| 1983 | 9367  | 260061 | 0.1161 | 0.5753 | 0.4844 | 0.5610 |
| 1984 | 7913  | 290839 | 0.2922 | 0.6491 | 0.5849 | 0.7073 |
| 1985 | 3662  | 126804 | 0.3419 | 0.5121 | 0.3348 | 0.5610 |
| 1986 | 10160 | 185277 | 0.6149 | 0.7608 | 0.1526 | 0.9024 |
| 1987 | 2771  | 128432 | 0.4982 | 0.6217 | 0.2893 | 0.9024 |
| 1988 | 14527 | 467773 | 0.6348 | 0.7848 | 0.7808 | 0.8049 |
| 1989 | 6276  | 112997 | 0.3567 | 0.6141 | 0.3900 | 0.6585 |
| 1990 | 6355  | 243994 | 0.3066 | 0.6250 | 0.3110 | 0.6585 |
| 1991 | 5309  | 276870 | 0.6496 | 0.7430 | 0.5363 | 0.9512 |
| 1992 | 15241 | 618581 | 0.8350 | 0.8484 | 0.8783 | 0.9024 |
| 1993 | 9617  | 217016 | 0.3713 | 0.6765 | 0.1981 | 0.7073 |
| 1994 | 19109 | 921882 | 0.9482 | 0.9043 | 0.9756 | 0.9512 |
| 1995 | 17998 | 483519 | 0.7882 | 0.8274 | 0.8288 | 0.8537 |
| 1996 | 8520  | 417460 | 0.7610 | 0.8073 | 0.7321 | 0.9024 |
| 1997 | 6898  | 145428 | 0.2501 | 0.5854 | 0.1765 | 0.4634 |
| 1998 | 6263  | 181396 | 0.0444 | 0.2874 | 0.1065 | 0.2195 |
| 1999 | 4133  | 59546  | 0.4171 | 0.5835 | 0.5856 | 0.8049 |
| 2000 | 1372  | 27595  | 0.5406 | 0.6543 | 0.8807 | 0.9512 |
| 2001 | 1969  | 39701  | 0.4927 | 0.6301 | 0.6815 | 0.8537 |
| 2002 | 2581  | 32254  | 0.4782 | 0.6272 | 0.7895 | 0.8537 |
| 2003 | 4171  | 146269 | 0.3006 | 0.4783 | 0.1627 | 0.5122 |
| 2004 | 898   | 30884  | 0.5784 | 0.6626 | 0.8236 | 0.9512 |
| 2005 | 6800  | 236405 | 0.1577 | 0.5614 | 0.2491 | 0.5122 |
| 2006 | 7922  | 154970 | 0.3857 | 0.6698 | 0.0733 | 0.7073 |
| 2007 | 3323  | 147582 | 0.4653 | 0.5966 | 0.1364 | 0.8049 |
| 2008 | 2882  | 87966  | 0.4016 | 0.5513 | 0.4880 | 0.6098 |
| 2009 | 2111  | 36592  | 0.4870 | 0.6291 | 0.7316 | 0.8049 |
| 2010 | 28244 | 694249 | 0.9404 | 0.9044 | 0.9267 | 0.9512 |
| 2011 | 4459  | 45005  | 0.5054 | 0.6305 | 0.6391 | 0.9512 |
| 2012 | 2115  | 22688  | 0.5078 | 0.6420 | 0.9725 | 0.9512 |
| 2013 | 8475  | 174738 | 0.3751 | 0.6731 | 0.0634 | 0.6585 |
| 2014 | 3005  | 24519  | 0.5024 | 0.6425 | 0.9248 | 0.9512 |
| 2015 | 14155 | 325111 | 0.6981 | 0.7957 | 0.6819 | 0.8537 |
| 2016 | 3257  | 86015  | 0.3600 | 0.5389 | 0.5355 | 0.5610 |
| 2017 | 5730  | 97637  | 0.3715 | 0.5963 | 0.4407 | 0.7561 |

Legend



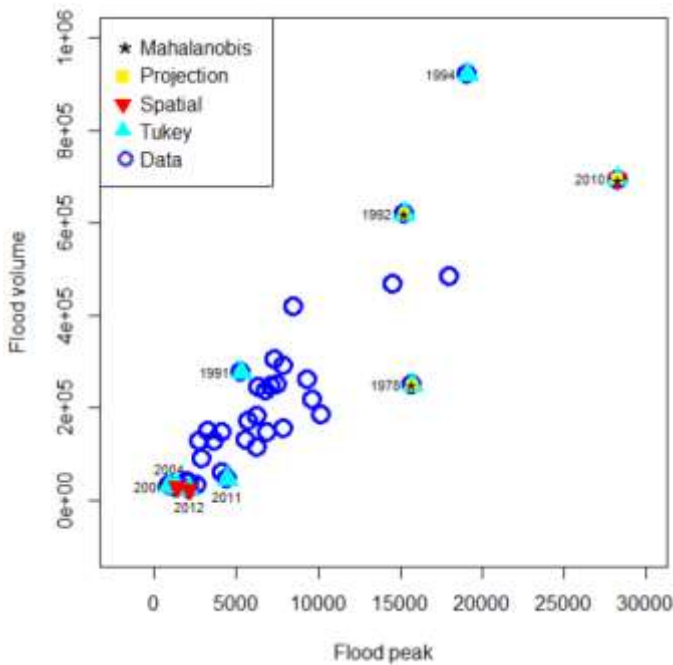


Fig. 2. Detected Outliers using Flood Peak and Volume.

The thresholds correspond to each outlyingness functions are computed by selecting 15% false outlier ratio and the number of true outliers as 5, this selection is similar to the choices made by [4] in such a way that the outlyingness value corresponds to an outlier must exceed their respective threshold values. Hence, 98% quantile is a corresponding threshold for the values of outlyingness. The computed values of the threshold for MO, PO, SO & TO are 0.9106, 0.8905, 0.9264, and 0.9444, respectively. The values of threshold approximately remain constant if the number of true outliers is considered greater than 5 with changed false outlier ratio i.e. 5%, 10% and 20%. The computed outlyingness values of MO, PO, SO & TO for years 1977-2017 are tabulated in Table II whereas Fig. 3 displays the detected outliers correspond to MO, PO, SO & TO with respect to their respective threshold values.

VI. RESULTS

A. Multivariate Result

The year 1994 contain outlyingness values greater than their respective threshold values by MO, PO & SO functions. Several years including years 1978, 1994, 2010 and 2012 are detected by TO function as outliers. The year 2010 is detected by MO and PO, and year 2012 is detected by SO functions as the closest value of outlyingness with respect to their threshold values. In addition, the year 1978 corresponds to the third highest MO and PO values whereas the year 2010 correspond the third highest SO value compare to their respective threshold values. Hence, it can objectively be inferred from Table I that the years 1994 and 2010 are identified as outliers by all the four functions of outlyingness. Whereas, the year 1978 is detected by the three and the year 2012 is detected by the two functions of outlyingness. For illustrative purpose a scatter plot constructed between bivariate (p,v) series (i.e. flood peak and flood volume) is

displayed through Fig. 2 so that the above interpretation can explicitly comprehensible. The years 1978, 1990, 1994, 2000, 2004, 2010, 2011, 2012 and 2014 computed as outliers by the outlyingness functions, among them the years 1978 and 1992 are present outside compare to the rest of the years whereas the years 1994 and 2010 are appear as outliers.

TABLE II. FUNCTIONAL RESULTS FOR PRINCIPAL COMPONENT (z<sub>1</sub>, z<sub>2</sub>)

| Year | z <sub>1</sub> | z <sub>2</sub> | MO     | PO     | SO     | TO     |
|------|----------------|----------------|--------|--------|--------|--------|
| 1977 | -2.29          | -2.339         | 0.1671 | 0.5215 | 0.3226 | 0.5122 |
| 1978 | 13.09          | -10.776        | 0.8342 | 0.8565 | 0.8512 | 0.9024 |
| 1979 | 6.218          | 7.085          | 0.6323 | 0.8256 | 0.6295 | 0.7561 |
| 1980 | -3.056         | 0.173          | 0.1077 | 0.2189 | 0.0381 | 0.1707 |
| 1981 | 9.388          | 8.516          | 0.7435 | 0.8548 | 0.7702 | 0.9512 |
| 1982 | -5.564         | 3.525          | 0.4119 | 0.5763 | 0.4923 | 0.7073 |
| 1983 | 7.676          | -2.302         | 0.4697 | 0.7249 | 0.6046 | 0.6585 |
| 1984 | -3.352         | -4.623         | 0.3994 | 0.6705 | 0.5732 | 0.8537 |
| 1985 | -9.07          | -1.06          | 0.5201 | 0.5805 | 0.7348 | 0.9512 |
| 1986 | -3.226         | -2.818         | 0.2465 | 0.5537 | 0.3994 | 0.6585 |
| 1987 | 1.832          | 5.992          | 0.4786 | 0.7791 | 0.5310 | 0.6585 |
| 1988 | 7.617          | -3.672         | 0.5177 | 0.7498 | 0.6238 | 0.7073 |
| 1989 | -5.09          | 0.182          | 0.2501 | 0.3459 | 0.2150 | 0.3171 |
| 1990 | 6.42           | -1.683         | 0.3743 | 0.6943 | 0.5049 | 0.5610 |
| 1991 | 20.652         | 9.365          | 0.8839 | 0.8928 | 0.9015 | 0.9512 |
| 1992 | 28.743         | 3.228          | 0.9157 | 0.8888 | 0.9422 | 0.9512 |
| 1993 | 7.174          | 7.721          | 0.6788 | 0.8378 | 0.6965 | 0.8049 |
| 1994 | 7.345          | -21.331        | 0.9217 | 0.8938 | 0.9077 | 0.9512 |
| 1995 | 9.233          | -6.894         | 0.6926 | 0.8068 | 0.7436 | 0.8049 |
| 1996 | 7.365          | -4.316         | 0.5350 | 0.7577 | 0.6280 | 0.7561 |
| 1997 | -4.721         | -0.359         | 0.2244 | 0.2994 | 0.2017 | 0.3659 |
| 1998 | 9.949          | 8.385          | 0.7490 | 0.8559 | 0.7998 | 0.9024 |
| 1999 | -6.452         | 1.793          | 0.3800 | 0.4861 | 0.4139 | 0.5610 |
| 2000 | -10.234        | 2.782          | 0.6053 | 0.6467 | 0.8525 | 0.9512 |
| 2001 | -5.921         | 6.952          | 0.6195 | 0.7290 | 0.7311 | 0.9512 |
| 2002 | -9.377         | 1.845          | 0.5479 | 0.5970 | 0.7410 | 0.8537 |
| 2003 | -3.783         | -0.194         | 0.1559 | 0.2193 | 0.0807 | 0.2195 |
| 2004 | -10.197        | 2.608          | 0.6002 | 0.6415 | 0.8266 | 0.9512 |
| 2005 | 0.865          | 2.131          | 0.1073 | 0.6408 | 0.3050 | 0.4634 |
| 2006 | -6.466         | -2.704         | 0.4169 | 0.6282 | 0.6128 | 0.8537 |
| 2007 | 0.884          | 6.854          | 0.5359 | 0.7897 | 0.5904 | 0.9024 |
| 2008 | -8.884         | 0.856          | 0.5077 | 0.5714 | 0.6634 | 0.7561 |
| 2009 | -8.824         | 1.963          | 0.5224 | 0.5829 | 0.6608 | 0.8049 |
| 2010 | -0.407         | -19.296        | 0.9007 | 0.8898 | 0.8743 | 0.9512 |
| 2011 | -6.425         | -1.158         | 0.3601 | 0.4990 | 0.4612 | 0.6098 |
| 2012 | -9.277         | -0.121         | 0.5251 | 0.5829 | 0.7346 | 0.9024 |
| 2013 | -6.19          | -2.373         | 0.3862 | 0.5996 | 0.5283 | 0.7561 |
| 2014 | -7.242         | 1.574          | 0.4233 | 0.5073 | 0.4995 | 0.6098 |
| 2015 | 1.404          | -1.802         | 0.0946 | 0.5918 | 0.3483 | 0.4634 |
| 2016 | -3.317         | 5.357          | 0.4567 | 0.7064 | 0.5341 | 0.9024 |
| 2017 | -6.487         | 0.934          | 0.3596 | 0.4502 | 0.3736 | 0.5610 |

Legend

■ Highest  
■ 2<sup>nd</sup> Highest  
■ 3<sup>rd</sup> Highest

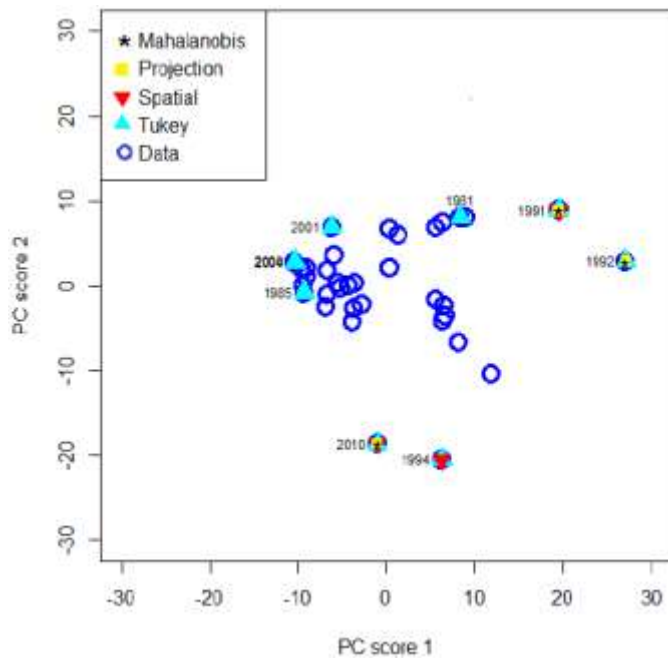


Fig. 3. Detected Outliers using Principal Component Scores.

### B. Functional Result

It is observed that the year 1994 contain outlyingness values greater than their respective threshold values by MO and PO functions whereas outlyingness value of the year 1992 is greater than the threshold value by SO function. Several years including 1991, 1992, 1994 and 2010 are detected by TO function as outliers. The year 1991 is detected by the PO, the year 1992 is detected by MO and the year 1994 is detected by the SO functions as a second highest outlyingness values compare to their respective threshold values. In addition, the year 2010 corresponds to the third highest MO and PO values, whereas the year 1991 corresponds to the third highest SO outlyingness value according to their respective threshold values.

Hence, it can distinctly be inferred from the values of Table II, the year 1994 is detected by all the four outlyingness functions as an outlier. Whereas the years 1991, 1992 and 2010 are identified as outliers by the three outlyingness functions. Above interpretation can better be comprehended by the scatter plot constructed between scores of initial two principal components (i.e. PC score 1 & score 2) and represented by Fig. 3 which reveals that the years 1981, 1985, 1991, 1992, 1994, 2000, 2001, 2004 and 2010 computed as outliers by the outlyingness functions, among them the years 1991 and 1992 are present outside compare to the rest of the years whereas the years 1994 and 2010 are appear as outliers.

The functional results are almost consistent with the results of the multivariate framework such that the years 1992, 1994 and 2010 have been detected as the most unusual flows in both the multivariate and functional context.

## VII. CONCLUSION

The nonparametric techniques based on depth function for outlier identifiers have been practiced in two different

frameworks of study that are multivariate hydrological data analysis and functional hydrological data analysis. The identification of outlier is essential for the appropriate selection of suitable hydrologic models so that risk associated with flood events can be authentically estimated. The methods employed in the present research are multivariate methods that are superior to previously practiced classical methods that were moment-based, follow normality assumption and component-wise techniques. The implemented techniques are based on depth function notion, free of moment, do not require normality assumption, and also affine invariant.

The proposed approaches have been implemented in two different frameworks of analysis. The intention of executing this study is to gauge the performance of proposed methodologies in both multivariate and functional context. The two most widely practice flood characteristics in hydrological analysis, peak ( $p$ ) & volume ( $v$ ) have been included to execute study in multivariate hydrological data analysis. Besides this, two initial scores of principal components ( $z_1, z_2$ ) used as a series of bivariate variables for executing functional hydrological data analysis since initial two principal components have a capability to capture major variation of data employed for analysis.

The outliers of both the framework are almost consistent but the results of functional analysis can be considered more reliable since it is based on complete information of flood hydrograph whereas flood characteristics ( $p, v$ ) are not able to generate hydrograph even though more than two characteristics of flood are included in study. Nevertheless, the multivariate results cannot be ignored and must be employed in a parallel complement to functional results so that dynamics of a hydrological event can be analysed to attain comprehensive information related to causes of flood.

## REFERENCES

- [1] V. Barnett, and T. Lewis, Outliers in Statistical Data, 3<sup>rd</sup> ed., John Wiley, Chichester, U.K, 1998.
- [2] V. Barnett, Environmental Statistics: Methods and Applications, John Wiley, Chichester, U.K, 2004.
- [3] X. Dang, and R. Serfling, "Nonparametric depth-based multivariate outlier identifiers, and masking robustness properties," Journal of Statistical Planning and Inference, vol. 140, no. 1, pp. 198–213, 2010. doi: 10.1016/j.jspi.2009.07.004.
- [4] F. Chebana and T. B. Ouarda, "Depth-based multivariate descriptive statistics with hydrological applications," Journal of Geophysical Research, vol. 116, D10120, 2011b. doi:10.1029/2010JD015338.
- [5] F. Chebana, S. Dabo-Niang and T. B. Ouarda, "Exploratory functional flood frequency analysis and outlier detection," Water Resources Research, vol. 48, no. 4, W04514, 2012. doi:10.1029/2011WR011040.
- [6] F. Chebana, and T. B. Ouarda, "Depth and homogeneity in regional flood frequency analysis," Water Resources Research, vol. 44, W11422, 2008. doi:10.1029/2007WR006771.
- [7] T. W. Anderson, An Introduction to Multivariate Statistical Analysis, 2nd ed., John Wiley, Chichester, U. K, 1984.
- [8] M. J. Schervish, "A review of multivariate analysis," Statistical Science, vol. 2, no. 4, pp. 413–417, 1987. doi:10.1214/ss/1177013111.
- [9] Y. Zuo, and R. Serfling, "General notions of statistical depth function," Annals of Statistics, vol. 28, no. 2, pp. 461–482, 2000b. doi:10.1214/aos/1016218226.
- [10] R.Y. Liu, J. M. Parelius, and K. Singh, "Multivariate analysis by data depth: Descriptive statistics, graphics and inference," Annals of Statistics, vol. 27, no. 3, pp. 783–858, 1999.



- [11] J. Zhang, "Some extensions of Tukey's depth function," *Journal of Multivariate Analysis* vol. 82, no. 1, pp. 134–165, 2002.
- [12] Y. Zuo, and R. Serfling, "On the performance of some robust nonparametric location measures relative to a general notion of multivariate symmetry," *Journal of Statistical Planning and Inference*, vol. 84, no. 1–2, pp. 55–79, 2000a. doi:10.1016/S03783758(99)00142-1.
- [13] C. H. Müller, "Depth estimators and tests based on the likelihood principle with application to regression," *Journal of Multivariate Analysis*, vol. 95, no. 1, pp. 153–181, 2005.
- [14] I. Hussain, and M. Uddin, "Functional and multivariate hydrological data visualization and outlier detection of Sukkur Barrage," *International Journal of Computer Applications*, vol. 178, no. 28, pp. 20-29, 2019. doi:10.5120/ijca2019919097.
- [15] I. Hussain, "Outlier detection using graphical and nongraphical functional methods in hydrology," *International Journal of Advanced Computer Science and Applications*, vol. 10, no. 12, pp. 438-445, 2019. doi: 10.14569/IJACSA.2019.0101259.
- [16] G. A. Griffiths, S. K. Singh, and A. I. McKerchar, "Flood frequency estimation in New Zealand using a region of influence approach and statistical depth functions," *Journal of Hydrology*, vol. 589, pp. 125-187, 2020. doi 10.1016/j.jhydrol.2020.125187.
- [17] M. Karamouz, F. ASCE, F. Ahmadvand, and Z. Zahmatkesh, "Distributed hydrologic modelling of coastal flood inundation and damage: Nonstationary approach," *Journal of Irrigation and Drainage Engineering*, vol. 143, no. 8, 2017. doi: 10.1061/(ASCE)IR.1943-4774.0001173.
- [18] W. Fan, L. Heng, D. Chao & D. Lieyun, "Knowledge representation using non-parametric Bayesian networks for tunneling risk analysis," *Reliability Engineering and System Safety*, Elsevier, vol. 191(C), 2019. doi: 10.1016/j.res.2019.106529.
- [19] L. Millán-Roures, I. Epifanio, and V. Martínez, "Detection of Anomalies in Water Networks by Functional Data Analysis," *Mathematical Problems in Engineering*, 2018. doi: org/10.1155/2018/5129735.
- [20] Y. Zhang, S. Huang, Q. Huang, G. Leng, H. Wang, and L. Wang, "Assessment of drought evolution characteristics based on a nonparametric and trivariate integrated drought index," *Journal of Hydrology*, vol. 579, 2019. doi 10.1016/j.jhydrol.2019.124230.
- [21] K. T. Peterson, V. Sagan, and J. J. Sloan, "Deep learning-based water quality estimation and anomaly detection using Landsat-8/Sentinel-2 virtual constellation and cloud computing," *GIScience & Remote Sensing*, vol. 57, no. 4, pp. 510-525, 2020. doi: 10.1080/15481603.2020.1738061.
- [22] J. Rhee, K. Park, S. Lee, S. Jang, and S. Yoon, "Detecting hydrological droughts in ungauged areas from remotely sensed hydro-meteorological variables using rule-based models," *Natural Hazards*, vol. 57, 2020. doi:10.1007/s11069-020-04114-5.
- [23] R. Padulano, G. D. Giudice, "A nonparametric framework for water consumption data cleansing: an application to a smart water network in Naples (Italy)," *Journal of Hydroinformatics*, vol. 22, no. 4, pp. 666–680, 2020. doi: org/10.2166/hydro.2020.133.
- [24] S. Samadi, D. L. Tufford, and G. J. Carbone, "Estimating hydrologic model uncertainty in the presence of complex residual error structures," *Stochastic Environmental Research and Risk Assessment*, vol. 32, pp. 1259–1281, 2018. doi: org/10.1007/s00477-017-1489-6.
- [25] Y. Zuo, "On general notions of depth for regression," arXiv e-prints, 2018.
- [26] S. Pool, M. Vis, and J. Seibert, "Evaluating model performance: towards a non-parametric variant of the Kling-Gupta efficiency," *Hydrological Sciences Journal*, vol. 63, no. 13-14, pp. 1941-1953, 2018. doi: 10.1080/02626667.2018.1552002.
- [27] M. A. Sherly, S. Karmakar, T. Chan, and C. Rau, "Design Rainfall Framework Using Multivariate Parametric-Nonparametric Approach," *Journal of Hydrologic Engineering*, vol. 21, no. 1, 2016. doi: org/10.1061/(ASCE)HE.1943-5584.0001256.
- [28] W.W.U.I. Wickramaarachchi, T.U.S. Peiris, and S. Samita, "Rainfall Trends in the North-Western and Eastern Coastal Lines of Sri Lanka Using Non – Parametric Analysis," *Tropical Agricultural Research*, vol. 31, no. 2, pp. 41-54, 2020. doi: 10.4038/tar.v31i2.8366.
- [29] U. Beyaztas, and Z. M. Yaseen, "Drought interval simulation using functional data analysis" *Journal of Hydrology*, vol. 579, 2019. doi: 10.1016/j.jhydrol.2019.124141.
- [30] J. Xia, P. Yang, C. Zhan and Y. Qiao, "Analysis of changes in drought and terrestrial water storage in the Tarim River Basin based on principal component analysis," *Hydrology Research*, vol. 50 no. 2, pp. 761–777, 2019. doi: 10.2166/nh.2019.033.
- [31] M. A. Hael, "Modeling of rainfall variability using functional principal component method: a case study of Taiz region, Yemen," *Modeling Earth Systems and Environment*, vol. 2, no. 7, 2020, doi: 10.1007/s40808-020-00876-w.
- [32] M. A. Hael, Y. Yongsheng, and B. I. Saleh, "Visualization of rainfall data using functional data analysis", *SN Applied Sciences*, vol. 2, no. 2, 2020. doi:10.1007/s42452-020-2238-x.
- [33] S.M. Shaharudin, N. Ahmad, N.H. Zainuddin, and N.S. Mohamed, "Identification of rainfall patterns on hydrological simulation using robust principal component analysis", *Indonesian Journal of Electrical Engineering and Computer Science*, vol. 11, no. 3, pp.1162-1167, 2018. doi: 10.11591/ijeecs.v11.i3.pp1162-1167.
- [34] J. Suhaila, Application of Functional Data Analysis in Streamflow Hydrograph. In: Kor LK., Ahmad AR., Idrus Z., Mansor K. (eds) *Proceedings of the Third International Conference on Computing, Mathematics and Statistics (iCMS2017)*. (2019) Springer, Singapore.
- [35] M. A. T. M. T. Rahman, S. Hoque, and A. H. M. Saadat, "Selection of minimum indicators of hydrologic alteration of the Gorai river, Bangladesh using principal component analysis," *Sustainable Water Resources Management*, 2017. doi: org/10.1007/s40899-017-0079-6.
- [36] S. Xiao, Z. Lu, and L. Xu, "Multivariate sensitivity analysis based on the direction of eigen space through principal component analysis," *Reliability Engineering & System Safety*, vol. 165, 2017. doi: 10.1016/j.res.2017.03.011.
- [37] S. K. Sharma, S. Tignath, S. Gajbhiye , and R. Patil, "Application of principal component analysis in grouping geomorphic parameters of Uttela watershed for hydrological modeling," *International Journal of Remote Sensing & Geoscience*, vol. 2, no. 6, 2013.
- [38] C. Gyamfi, J. M. Ndambuki, and R. W. Salim, "Simulation of sediment yield in a semi-arid River Basin under changing land use: an integrated approach of hydrologic modelling and principal component analysis," *Sustainability*, vol. 8, no. 11, 2016. doi: 10.3390/su8111133.
- [39] S. Yue, T. B. Ouarda, B. Bobée, P. Legendre, and P. Bruneau, "The Gumbel mixed model for flood frequency analysis," *Journal of Hydrology*, vol. 228, pp.88-100, 1999. doi:10.1016/S0022-1694(99)00168-7.
- [40] J. T. Shiau, "Return period of bivariate distributed extreme hydrological events," *Stochastic Environmental Research and Risk Assessment*, vol. 17, pp.42–57, 2003. doi:10.1007/s00477-003-0125-9.
- [41] S. Naz, M. J. Iqbal, S. M. Akhter, and I. Hussain, "The Gumbel mixed model for food frequency analysis of Tarbela," *The Nucleus*, 53(3), pp. 171-179, 2016.

# A Hybrid Approach to Enhance Scalability, Reliability and Computational Speed in LoRa Networks

S. Raja Gopal<sup>1</sup>

Research Scholar, Department of ECE  
Koneru Lakshmaiah Education Foundation

V S V Prabhakar<sup>2</sup>

Professor, Department of ECE  
Koneru Lakshmaiah Education Foundation

**Abstract**—The spreading out of Internet of Things (IoT) facilitates with new wireless communication Technologies. To have reliable communication for long duration, low power wide area can be aimed at sensor nodes. These days, LoRa has become a recognizable preference for IoT based solutions. In this paper LoRa network performance is analyzed. A hybrid technique is proposed to overcome the scalability, reliability and computational speed issues in LoRa network. In the proposed hybrid technique, a lightweight scheduling technique is used to address scalability and reliability issues and then pruning algorithm is also incorporated to enhance the computational speed of LoRa network for IoT applications. Further, LoRa network for IoT applications is analyzed using LoRaWAN NS-3 module and the parameters packet error ratio (PER), network throughput and fairness for improved reliability and scalability are illustrated. Simulation results are obtained for multiple gateways scenarios. The analysis presents that the LoRa network has addressed scalability and reliability issues using lightweight scheduling technique and computational speed is also enhanced using pruning algorithm. Therefore, the hybrid technique illustrated for LoRa network for IoT applications is in good agreement.

**Keywords**—Hybrid technique; Internet of Things (IoT); lightweight scheduling; LoRa; pruning algorithm; scalability; reliability

## I. INTRODUCTION

To accomplish enhanced horizontal integration among IoT service, data aggregation, smart management and protocol adaptation services are very much essential. An overview of IoT (Internet of Things) with importance on application issues, protocols and enabling technologies had been provided clearly in [1] by the author Ala Al-Fuqaha. The public availability of LoRa technology specifications is making it more momentum to balance current IoT standards as an enabler of smart city applications. LoRa supports a communications range of 5-15 km and thousands of devices can be connected to internet through a single LoRa gateway. It is a PHY (physical) layer technology built around CSS (Chirp Spread Spectrum) modulation. The major challenges faced by LoRa are: scalability, interference from other co-located networks and restrictions on duty cycle. Many researchers proposed various algorithms and designs to overcome these challenges and presented the performance evaluation of

different parameters by simulating in NS-3 (Network Simulator).

## II. RELATED WORK

Authors in [2] have implemented a LoRa deployment test using LoRa network and presented LoRa coverage analysis. In [3] authors have presented the evaluation of LoRa for wireless sensor networks. A comprehensive analysis of LoRa modulation including the frame format, receiver sensitivity, data rate, spreading factor for IoT applications is described in [4]. A relative revision of LPWAN technologies for deployment of large scale IoT and known and unknown facts of LoRa are presented in [5-6]. In [7], a new private LoRa network is designed and implemented for IoT applications. Authors also addressed the hardware design and implementation of Lora Gateway and also the performance evaluation of LoRa networks in typical environments. The authors Brecht Reynders et al. [8] have designed a new MAC protocol RS-LoRa and presented a two-step lightweight scheduling algorithm to improve scalability and reliability of LoRaWAN networks. In this protocol, the entire bandwidth is splitted into a synchronous downlink channel and many asynchronous uplink/downlink channels as shown in Fig. 1.

LoRa gateways coordinated the lightweight scheduling in two steps. In the first step, the gateway schedules nodes dynamically by specifying the allowed SF (spreading factors) and transmission powers on each channel in a coarse-grained manner. In the second step, according to coarse-grained scheduling information sent by gateway, a node decides its own spreading factor transmission power, time and channel to transmit in a distributed manner. This proposed scheduling algorithm is implemented in NS-3 with the architecture illustrated in Fig. 2. With this architecture implementation, four different scenarios are simulated in NS-3: single cell scenario with one gateway for legacy LoRaWAN and RS-LoRa protocols and multi-cell scenario with seven gateways for legacy LoRaWAN and RS-LoRa protocols.

The author Floris Van den Abeele et al. [9] has presented the analysis of scalability of large scale LoRaWAN networks using NS-3. LoRaWAN. NS-3 module was modeled as presented in Fig. 3. By using one, two or four gateways, different scenarios are simulated and the scalability of LoRa networks was analyzed.

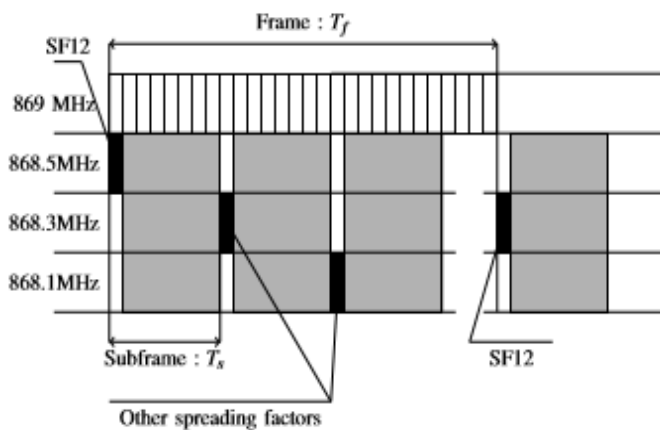


Fig. 1. Channel Assignment (Brecht Reynders et. al.).

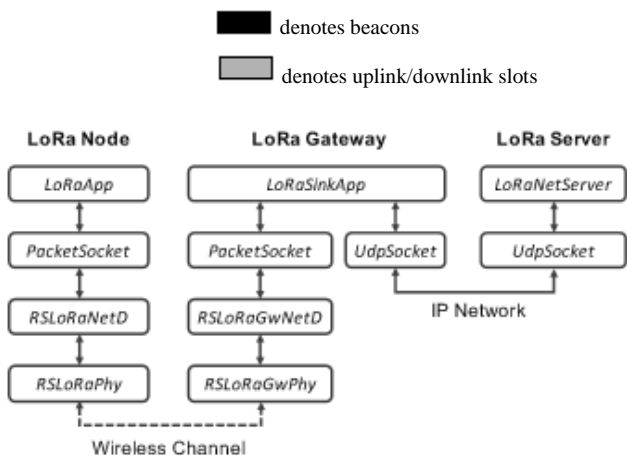


Fig. 2. RS-LoRa Architecture in NS-3 (Brecht Reynders et. al.).

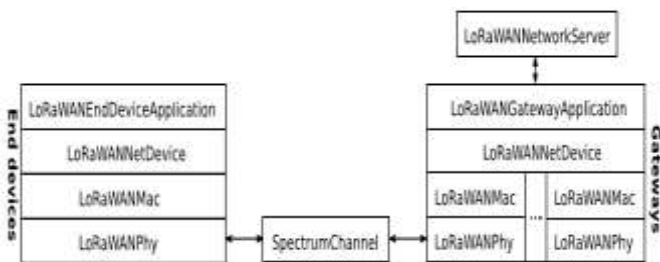


Fig. 3. Class A End Device, Gateway and Network Server in Modeled LoRaWAN NS-3 Module (Floriss Van Den Abeele et. al.).

In [10], the author Davide Magrin et al. implemented NS-3 model in urban environment to cram the performance of LoRa networks. The author Daniele Croce [11] investigated the link level performance of LoRa networks with quasi orthogonal spreading factor chirps. In [12], LoRa network performance with different PHY settings and environmental conditions was analyzed using experimental set-up. For confirmed and unconfirmed UL transmissions, performance investigation of Class A devices in LoRaWAN was presented in [13]. Using NS-3, system level and link level performance of LoRa networks was analysed for IoT and wireless sensor networks in [14-15]. Through a numerical model, LoRa performance was assessed with additive white Gaussian noise channels in

[16]. In [17], authors analyzed the performance metrics of LoRa network by implementing efficient algorithms for spread factor assignment and power allocation. These algorithms have enhanced throughput and fairness with efficient energy consumption. This paper aims to address scalability and reliability issues along with improved throughput and fairness by incorporating lightweight scheduling technique in LoRa networks. In this paper, an energy efficient hybrid algorithm to enhance reliability and scalability with enhanced speed for large computations using pruning algorithm is proposed and implemented in NS-3. Rest of the paper is organized as: Section III presents the lightweight technique and algorithms to be used in LoRa networks. Section IV describes the implementation of LoRa network in NS-3 environment and the performance evaluation of simulated LoRa network in NS-3. Finally, the conclusion of paper is given in Section V.

### III. LORA NETWORK DESIGN

A LoRa network consists of Gateway and end nodes. To address scalability and reliability issues a light-weight scheduling technique is implemented at Gateway as well as at end nodes. In this light weight scheduling technique at Gateway, spreading factor and allowed transmission power for each channel are specified and schedule accordingly. The end nodes have the freedom to select their own transmission parameters by following this scheduling information. These transmission parameters include selected channel, selected spreading factor and transmission power and offset time. Transmission power and spreading factor are the important factors that affect the reliability of a Lora network. The Lora Gateway can control the transmission power of each node and by varying the thresholds of received signal strength the amount of traffic can be increased or decreased thereby resulting in low or high reliability for that specific channel. Spreading factor value usually ranges from 7 to 12. To achieve high channel reliability more spreading factors should be allowed us different spreading factor packets can be received simultaneously. This also reduces the occurrence of collisions in the network. But more energy is consumed when more spreading factors are used in the network by different end nodes. Therefore a balance between energy consumption and reliability should be achieved. At the end notes in Lora network, lightweight scheduling technique is performed by determining the transmission parameters at each node by using the information transmitted by Lora gateway. At each note to determine the transmission parameters the following Algorithm-1 is presented. The input parameters from gateway are:

- $I \leftarrow$  channel set,
  - $P_{RSS} \leftarrow$  Received signal strength at node n,
  - $S_j \leftarrow$  set of allowed spreading factor for channel j,
  - $P_j \leftarrow$  target uplink received signal strength for channel j.
- The transmission parameters determined at each node using light weight scheduling technique are:
- $S_n \leftarrow$  selected spreading factor,

$P_n \leftarrow$  selected transmission power,

$C_n \leftarrow$  selected channel.

In the algorithm,

$P_{SF} \leftarrow$  Probability of selecting spreading factor  $sf$ ,

$P_{PL} \leftarrow$  Path loss power,

$R_{sf}$  and  $R_{sf'}$  are the bit rate achieved with the spreading factor  $sf$  and  $sf'$ ,

---

**Algorithm-1: Determination of transmission parameters**

---

- Step 1:  $P_{Temp} \leftarrow 0$   
 Step 2:  $flag \leftarrow false$   
 Step 3: **for**  $j \in I$  **do**  
 Step 4:     **if**  $P_{Temp} < P_j < P_{RSS}$  **then**  
 Step 5:          $C_n \leftarrow j$   
 Step 6:          $P_{Temp} \leftarrow P_j$   
 Step 7:          $flag \leftarrow true$   
 Step 8:     **end if**  
 Step 9: **end for**  
 Step 10: **if**  $flag = true$  **then**  
 Step 11:  $SF_n \leftarrow$  select spreading factor using  $P_{SF} = \frac{R_{sf}}{\sum_{sf' \in S_j} R_{sf'}}$ ,  $\forall sf \in S_j$      # $P_{SF}$  is the probability of selecting spreading factor of  $sf$   
 Step 12:  $P_n \leftarrow P_{Temp} + P_{PL} - 2.5 SF_n + P_{Offset}$   
 Step 13: **else**  
 Step 14:      $SF_n \leftarrow 7$   
 Step 15:      $P_n \leftarrow 0$  dBm  
 Step 16:      $C_n \leftarrow \arg \max_{j \in I} P_j$   
 Step 17: **end if**
- 

In illustrated Algorithm-1, when a message is to be transmitted by a node, most suited channel is selected in step 3 to step 9. Upon selecting a suitable channel the random spreading factor is selected in step 11. Then in step 12 transmission power is calculated based on the selected spreading factor and target received signal strength. If a channel is not selected in step 3 to 9, indicates that the node is very close to Gateway and in such case the lowest spreading factor with value 7 is selected with minimum power consumption. This is presented in step 13 to 17. Finally transmission time for each node is selected at the end of the algorithm after selecting channel spreading factor and transmission power at each node.

IV. IMPLEMENTATION IN NS-3 AND PERFORMANCE EVALUATION

LoRa network is implemented in Network Simulator NS-3. The implementation architecture consists of four layers: physical, data-link, transport and application layers as described in [8]. For simulation setup initially multi-cell scenario with seven gateways are considered as illustrated in Fig. 4. In simulation setup, gateways are located at 30 meters height and nodes are located 1meter above the ground level. Via gateway, the nodes will transmit data packets to network server for every 2 minutes.

The parameters that are set-up in NS-3 simulation are tabulated in Table I. To investigate the performance of LoRa network, all seven gateways are enabled and the cell radius is increased. The end nodes are randomly distributed and for three values of nodes 1000, 2000 and 3500, simulations results are obtained. The performance parameters packet error ratio, throughput and fairness are analyzed here.

**Packet Error Ratio:** The performance of scalability and reliability can be demonstrated using the parameter Packet Error Ratio (PER). In Fig. 5 packet error ratio for 1000 nodes and 3500 nodes has been illustrated by varying the gateway distance up to 1.4 Km. PER is 6.7% for 1000 nodes at a distance of 1.3 Km from GW1. It can be observed that PER is less when the nodes are at a distance of 700 – 1000 m from GW1 and nearer to GW2, GW3, GW4, GW5, GW6 and GW7. The maximum PER achieved for 3500 nodes is only 20.7% and hence PER is improved. Therefore, scalability and reliability of LoRa network are enhanced using lightweight scheduling technique.

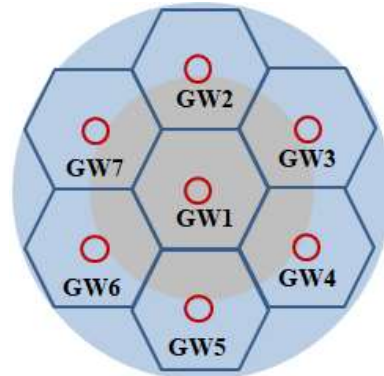


Fig. 4. Location of Seven Gateways in NS-3 Simulation.

TABLE I. PARAMETERS SET-UP IN NS-3 SIMULATION

| Parameter                        | Value       | Unit  |
|----------------------------------|-------------|-------|
| Number of nodes                  | 100 to 4000 | /     |
| Maximal distance to the GW1      | 1000        | m     |
| Tf                               | 10          | min   |
| Ts                               | 1           | min   |
| Average inter-packet interval    | 120         | Sec   |
| Distance between gateways        | 1000        | M     |
| Packet length (excluding header) | 51          | bytes |

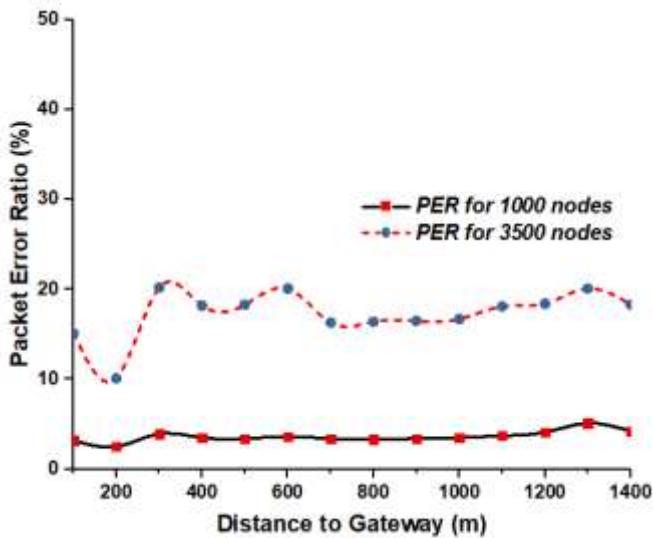


Fig. 5. Packet Error Ratio (PER) for Multi-Cell Scenario with Seven Gateways.

**Throughput:** The throughput of the LoRa network for lightweight scheduling technique under multi-cell scenario is presented in Fig. 6. The results are demonstrated for three different values of nodes: 1000, 2000 and 3500. It is 3.2Kb/s for 1000 nodes, 6.2Kb/s for 2000 nodes and 9.6Kb/s for 3500 nodes. The network throughput increases for increased number of nodes as the network reliability has been increased with lightweight scheduling technique.

**Fairness:** The network performance can be increased by increasing the fairness parameter. The lightweight scheduling technique improves the fairness of network by reducing capture effect [8]. Fairness under multi-cell scenario is investigated for three values of nodes: 1000, 2000 and 3500 and the investigation results are illustrated in Fig. 7. For 1000 nodes the fairness is 99.6%, it is 98.4% for 2000 nodes and 95.1% for 3500 nodes. Hence, with improved reliability and reduced capture effect, fairness of LoRa network is enhanced.

For simulation, only seven gateways are considered as shown in Fig. 4 and the gateway distance is increased till 1.4Km. In real networks, especially in urban areas, network coverage area should be large and so 19 and 37 gateway scenarios are further considered in simulation. The gateways are placed around a central gateway in hexagonal grid structure as found in cellular coverage of large area. The arrangement of 19 and 37 gateways is as presented in Fig. 8. In simulations, each gateway covers a radius of 1500m and the end nodes are placed in a radius of 7500m. Inter-cell interference can be simulated for such large network. Fig. 9 illustrates the simulation coverage of end nodes in NS-3 with different gateway densities and their respective spreading factors.

In a realistic network model, simulation of such large network with 37 gateways and a coverage area of 7500m

radius require some optimizations to speed-up the computations. Pruning algorithm is one such approximation which prunes the end nodes that cause unnecessary interference at the central gateway GW1 and speeds-up the simulation computations. Pruning algorithm is represented in Algorithm-2. Let  $R_s$  be the simulation radius. Let  $I_R$  be the set of end nodes that are centered on central gateway GW1 and inside the radius R. Let  $O_R$  be the set of end nodes that are outside the radius R. Find the smallest R for which the total energy of inside end nodes is less than the energy of outside end nodes. Then prune all the end nodes in the set of  $O_R$  by detaching PHYs from the channel and removing end nodes from simulation. The effect of pruning is well presented in Fig. 10.

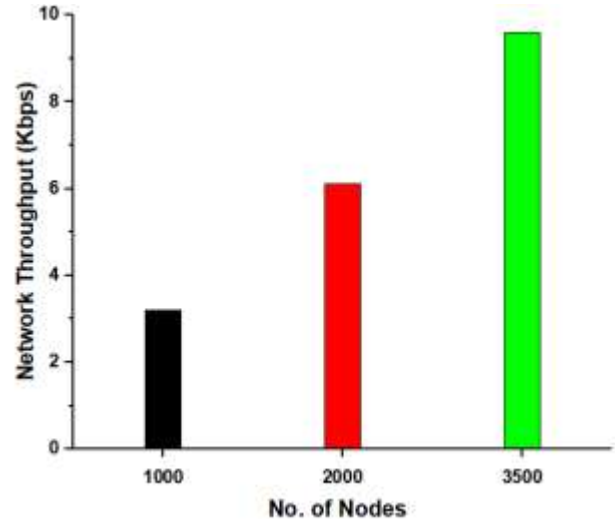


Fig. 6. Network Throughput for Multi-Cell Scenario with Seven Gateways.

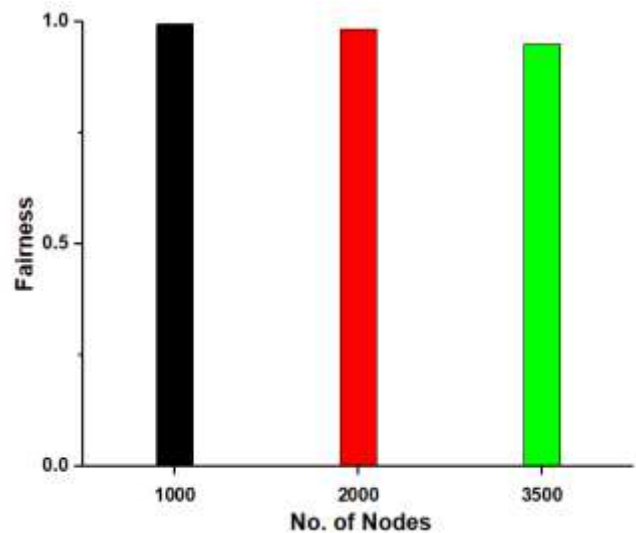
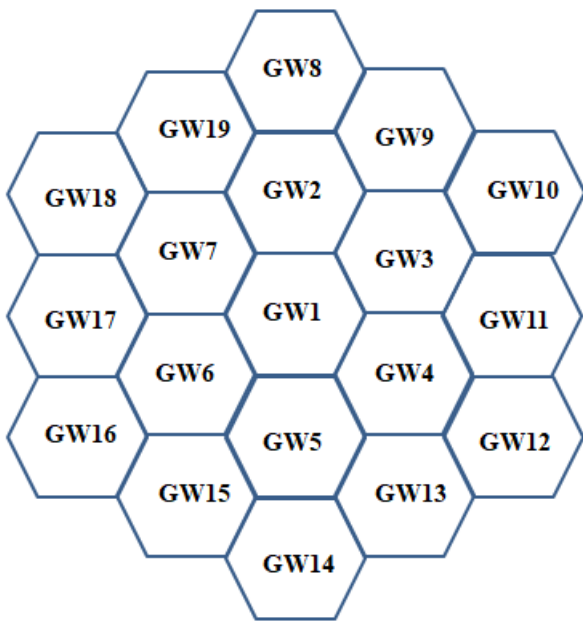
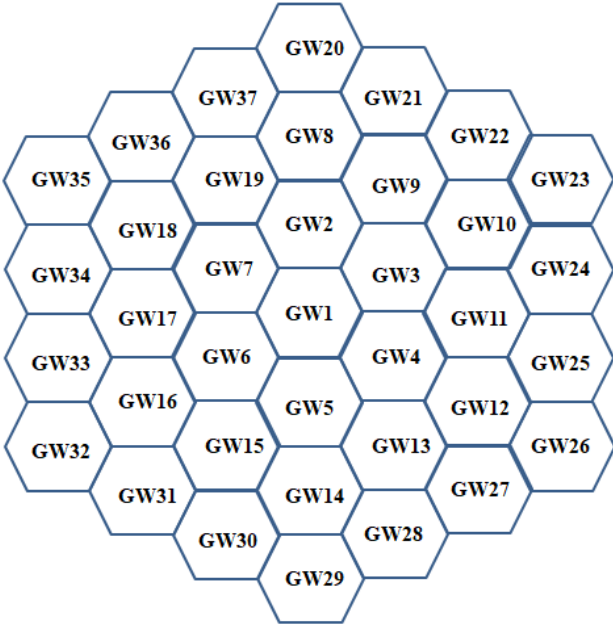


Fig. 7. Fairness for Multi-Cell Scenario with Seven Gateways.

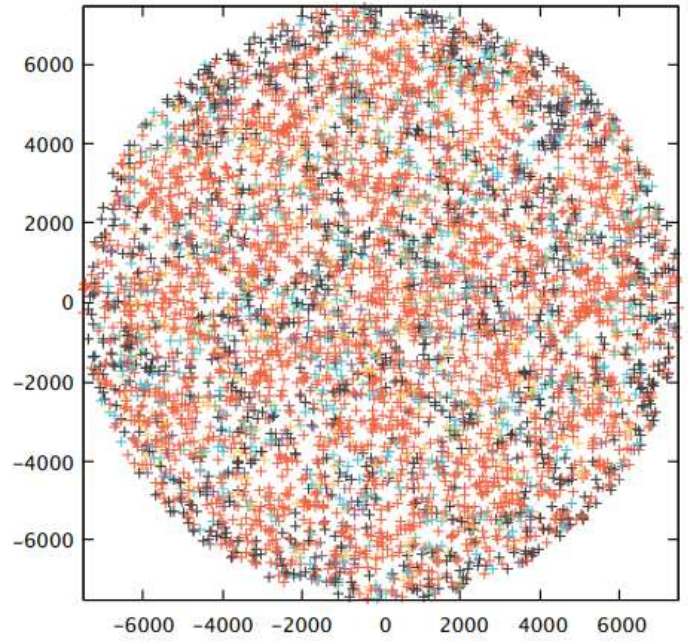


(i) Nineteen Gateways.

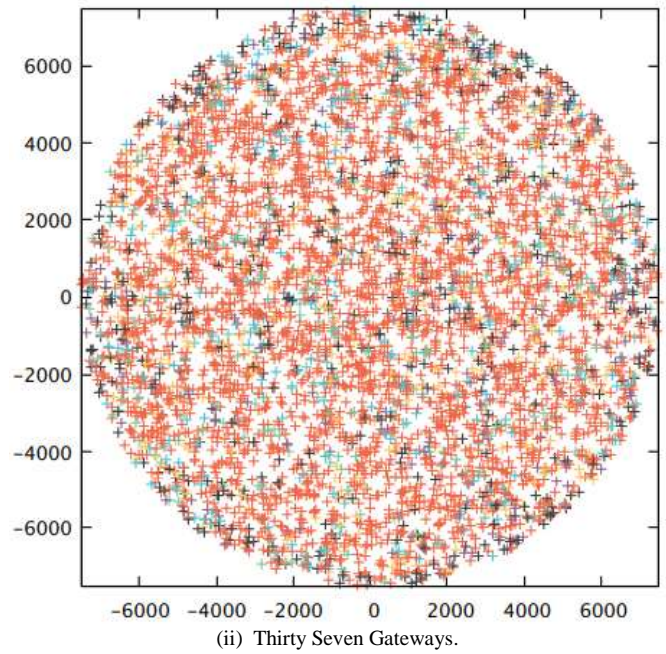


(ii) Thirty Seven Gateways.

Fig. 8. Arrangement of Gateways for Large Coverage Area.



(i) Nineteen Gateways.



(ii) Thirty Seven Gateways.

Fig. 9. End Nodes with different Gateway Densities and their Spreading Factors.

In Algorithm-2:

$E_{inside} \leftarrow$  Energy of inside node

$E_{outside} \leftarrow$  Energy of outside node

Table. II Describes the comparison of present work with previous related works.

**Algorithm-2 End Nodes Pruning.**

Input:  $R_s \leftarrow$  Simulation radius

Step1:  $R \leftarrow 0$

Step2:  $exit \leftarrow false$

Step3: **for**  $SF \in 7, \dots, 12$

Step4: **for**  $R < R_s, exit \leftarrow false$

Step5:  $E_{inside} \leftarrow E_{outside} \leftarrow 0$

Step6: **for** Each end device  $ed$

Step 7: **if**  $ed$ 's  $SF_e \leftarrow SF$

Step 8:  $E \leftarrow$   $ed$ 's energy for a transmission

Step 9: **if**  $ed$ 's distance from center  $< r$

Step 10:  $E_{inside} \leftarrow E_{inside} + E$

Step 11: **else**

Step 12:  $E_{outside} \leftarrow E_{outside} + E$

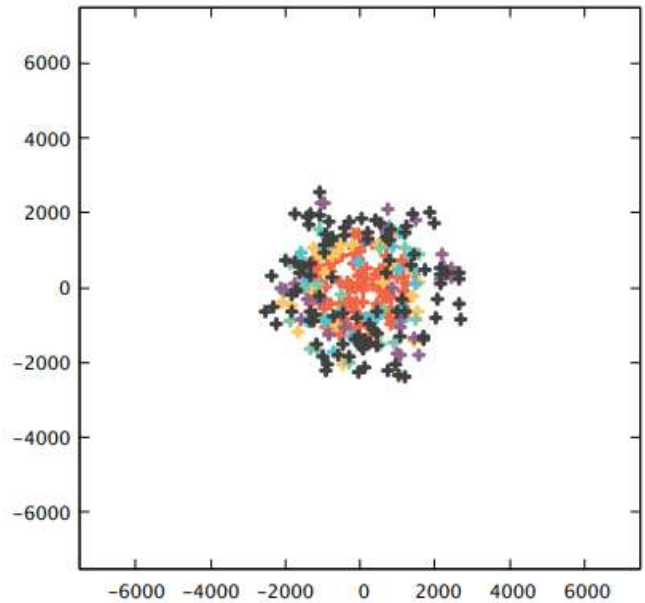
Step 13: **if**  $E_{outside} < E_{inside} / 10$

Step 14: **exit**  $\leftarrow true$

Step 15: **else**

Step 16:  $R \leftarrow R + \epsilon$

Step 17: **return**



(ii) With Pruning Algorithm.

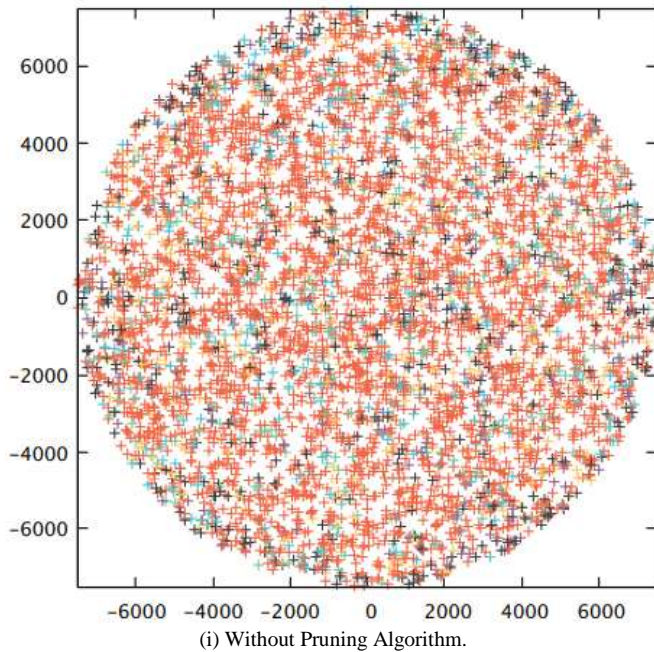
Fig. 10. End Nodes Coverage with and Without Pruning Algorithm.

TABLE II. COMPARISON TABLE

| Reference                 | Simulator | Algorithm/Technique   | Performance parameters    | Simulation Scenario | No. of Gateways |
|---------------------------|-----------|---|---------------------------|---------------------|-----------------|
| Brecht Reynders [8]       | NS-3      | Lightweight Scheduling                                      | PER, Throughput, Fairness | Single, Multiple    | 1,7             |
| Floris Van den Abeele [9] | NS-3      | Interference Model  | Scalability               | Multiple            | 7               |
| Daniele Croce [11]        | Matlab    | Quasi Orthogonal Spreading Factor Chirps                    | BER, SIR threshold        | Single              | 1               |
| Furqan Hameed Khan [13]   | NS-3      | Confirmed and Unconfirmed UL Transmissions, Class A Devices | PDR, UL Throughput        | Single              | 1               |
| This paper                | NS-3      | Lightweight Scheduling, Pruning Algorithm                   | PER, Throughput, Fairness | Multiple            | 7,19,37         |

V. CONCLUSION

In this paper, a hybrid approach to improve scalability, reliability of LoRa network and to speed-up the computations in case of large networks is presented. Lightweight scheduling technique is used to enhance scalability and reliability of LoRa network. Performance parameters packet error ratio (PER), throughput of network and fairness among nodes are evaluated and illustrated. Then by considering large network with 19 and 37 gateways in NS-3 simulation, pruning algorithm is implemented to speed-up the computations by pruning end nodes that are outside the region of interest.



(i) Without Pruning Algorithm.

LoRa networks can be reliable and scalable and for real scenarios, computations can be speed-up. Further, this work can be carried out with real-time experimentation of LoRa network. LoRa network and edge computing can be integrated together to have sustainable solutions for IoT applications.

#### REFERENCES

- [1] Al-Fuqaha, M. Guizani, M. Mohammadi, M. Aledhari, and M. Ayyash, "Internet of things: A survey on enabling technologies, protocols, and applications," *IEEE Communications Surveys & Tutorials*, vol. 17, no. 4, pp. 2347–2376.
- [2] M. Centenaro, L. Vangelista, A. Zanella, and M. Zorzi, "Long-Range Communications in Unlicensed Bands: the Rising Stars in the IoT and Smart City Scenarios," *IEEE Wireless Communications*, vol. 23, Oct. 2016.
- [3] A. J. Wixted et al. Evaluation of lora and lorawan for wireless sensor networks. In 2016 IEEE SENSORS, Oct 2016.
- [4] Aloys Augustin, Jiazi Yi, Thomas Clausen, William Mark Townsley, "A Study of LoRa: Long Range & Low Power Networks for the Internet of Things", *Sensors*, 2016.
- [5] K. Mekki et al. A comparative study of lpwan technologies for largescale iot deployment. *ICT Express*, 2018.
- [6] J. C. Liando et al. Known and unknown facts of lora: Experiences from a large-scale measurement study. *ACM Trans. Sen. Netw.*, 15(2):16:1–16:35, February 2019.
- [7] Qihao Zhou, Kan Zheng, Lu Hou, Jinyu Xinz , Rongtao Xu, "Design and Implementation of Open LoRa for IoT" , *IEEE ACCESS*, Vol. 7, 2019.
- [8] Brecht Reynders, Qing Wang, Pere Tuset-Peiro, Xavier Vilajosana, Sofie Pollin, "Improving Reliability and Scalability of LoRaWANs Through Lightweight Scheduling", *IEEE Internet of Things Journal*, vol. 5, no. 3, 2018.
- [9] Floris Van den Abeele , Jetmir Haxhibeqiri, Ingrid Moerman, Jeroen Hoebeke, "Scalability Analysis of Large-Scale LoRaWAN Networks in NS-3", *IEEE Internet of Things Journal*, vol. 4, no. 6, 2017.
- [10] Davide Magrin, Marco Centenaro, and Lorenzo Vangelista, "Performance Evaluation of LoRa Networks in a Smart City Scenario", *IEEE ICC 2017 SAC Symposium Internet of Things Track*, May 2017.
- [11] Daniele Croce, Michele Gucciardo, Stefano Mangione, Giuseppe Santaromita, Ilenia Tinnirello, "Impact of LoRa Imperfect Orthogonality: Analysis of Link-level Performance", *IEEE Communications Letters*, Vol: 22 , Issue: 4 , 796 – 799, April 2018.
- [12] Marco Cattani, Carlo Alberto Boano and Kay Römer, "An Experimental Evaluation of the Reliability of LoRa Long-Range Low-Power Wireless Communication", *Journal of Sensors and Actuator Networks*, vol. 6, no. 2, 2017.
- [13] Furqan Hameed Khan and Marius Portmann, "Experimental Evaluation of LoRaWAN in NS-3", *International Telecommunication Networks and Applications Conference (ITNAC)*, Nov. 2018.
- [14] Andrew J Wixted ; Peter Kinnaird ; Hadi Larijani ; Alan Tait ; Ali Ahmadi ; Niall Strachan, "Evaluation of LoRa and LoRaWAN for wireless sensor networks", *IEEE SENSORS*, 2016.
- [15] Luca Feltrin, Chiara Buratti, Enrico Vinciarelli, Roberto De Bonis, Roberto Verdone, "LoRaWAN: Evaluation of Link- and System-Level Performance", *IEEE Internet of Things Journal*, vol.5, no. 3, 2018.
- [16] Menno J. Faber, Klaas M. vd Zwaag, Willian G. V. dos Santos, Helder R. O. Rocha, Marcelo E. V. Segatto, Jair A. L. Silva, "A Theoretical and Experimental Evaluation on the Performance of LoRa Technology", *IEEE Sensors Journal*, vol. 20, no. 6, 2020.
- [17] Licia Amichi, Megumi Kaneko, Ellen Hidemi Fukuda, Nancy El Rachkidy, Alexandre Guitton, "Joint Allocation Strategies of Power and Spreading Factors with Imperfect Orthogonality in LoRa Networks", *IEEE Transactions on Communications*, vol.68, no. 6, 2020.



# A New Online Plagiarism Detection System based on Deep Learning

El Mostafa Hambi<sup>1</sup>, Faouzia Benabbou<sup>2</sup>

Information Technology and Modeling Laboratory  
Faculty of Sciences Ben M'sik, University Hassan II, Casablanca, Morocco

**Abstract**—The Plagiarism is an increasingly widespread and growing problem in the academic field. Several plagiarism techniques are used by fraudsters, ranging from a simple synonym replacement, sentence structure modification, to more complex method involving several types of transformation. Human based plagiarism detection is difficult, not accurate, and time-consuming process. In this paper we propose a plagiarism detection framework based on three deep learning models: Doc2vec, Siamese Long Short-term Memory (SLSTM) and Convolutional Neural Network (CNN). Our system uses three layers: Preprocessing Layer including word embedding, Learning Layers and Detection Layer. To evaluate our system, we carried out a study on plagiarism detection tools from the academic field and make a comparison based on a set of features. Compared to other works, our approach performs a good accuracy of 98.33 % and can detect different types of plagiarism, enables to specify another dataset and supports to compare the document from an internet search.

**Keywords**—*Plagiarism detection; plagiarism detection tools; deep learning; Doc2vec; Stacked Long Short-Term Memory (SLSTM); Convolutional Neural Network (CNN); Siamese neural network*

## I. INTRODUCTION

According to Risquez et al. [1] “the Plagiarism is conceptualized as the theft of others’ words or ideas without citing the proper reference and thus without giving the accurate credit to the original author”. Depending of depth of transformation performed on the original text, the plagiarism can be classified into five categories [2]:

- Copy & paste plagiarism (word by word) [3]: it is the act of copying text and passing without reference of original authors.
- Paraphrasing [4]: the content is copied from different source without acknowledging authors.
- Use of false references [5]: There are certain cases where user quotes the original sources, but the information provided in the articles are not match with the source provided at the end of the article.
- Plagiarism with translation [6]: it is the act of translating text from language to another.
- Plagiarism of ideas [7]: it is the most difficult plagiarism to detect where fraudsters steal other authors' ideas and present them in a fully modified version of the original text and own the new version.

Plagiarism is applied in different areas such as literature, music, software, scientific articles, newspapers, advertisements, websites, etc. Despite the sanctions applied in cases of cheating and plagiarism in Bulgarian universities, more than 50% of teachers believe that these procedures are not efficient [8]. As the use of internet increases plagiarism becomes a big challenge in schools, universities to maintain the academic integrity. Thus, the use of efficient plagiarism detection tools has become very urgent in many higher education institutions. However, the effectiveness of these plagiarism detection systems depends on their ability to discover different fraudsters’ strategies to modify the text without changing its semantics [9].

As part of NLP research topic, the plagiarism detection methods are based on natural language techniques to process and analyze the structure of documents. Many solutions have been proposed for plagiarism detection, and most of them are based on concept extraction using corpus such as ontologies (e.g. WordNet) to perform a semantic representation of documents. However, these approaches depend on the quality of corpus and an appropriate annotation to choose the best concept that semantically represents a word. In addition, the problem of ambiguity may arise when choosing the concept that semantically represents the word, so the meaning of the processed sentences may be lost if we choose the wrong concept [10]. Some examples of this classical plagiarism detection methods are [11]: Fingerprinting, String matching, Jaccard similarity, Bag of word analyzing and Shingling.

With the emergence of artificial intelligence, many techniques have been proposed, ranging from supervised, unsupervised machine learning techniques to deep learning, and have been successfully applied in various fields. In-depth learning provides models with multiple processing layers capable of learning data representations with multiple levels of abstraction. Recently many applications of deep learning in NLP domains, has been proposed and their performance was very encouraging as Chatbots programming, sentiment analysis and Question and Answering. In this context, we propose an online plagiarism detection system based on Doc2vec technique for word embedding, and SLSTM and CNN deep learning algorithms. Our system can perform many tasks of plagiarism detection and the results found are very promising.

The rest of this paper is organized as follows. Section 2 presents a review of the most relevant plagiarism software. Section 3 illustrates our plagiarism detection system. In Section 4 we describe the components of our online system.

Section 5 draws some interpretations about the current state of existing discovery tools and compare them with our system. In the section 6, we finish by a global conclusion.

## II. PLAGIARISM TOOLS REVIEW

### A. Software Description

In the context of academic plagiarism, few tools are proposed, and this section is devoted to describing the most recognized in the scientific community and in different universities. In our latest state of the art we focused on the proposed systems for plagiarism detection based on deep learning, unfortunately we did not find any implementation of these systems. Asim M. El Tahir Ali et al. have proposed an interesting comparative study from five plagiarism detection tools [13]: PlagAware, The PlagScan, CheckForPlagiarism.net, iThenticate and PlagiarismDetection.org. Inspired by this research, we conducted an overview of the top plagiarism detection tools based on some important criteria that a good system would have. Firstly, we used the comparison parameters in [13] as:

- Add a new database is the ability to add a new database in comparison and plagiarism detection.
  - Add a new corpus is the ability to add a new corpus for learning to detect other types of plagiarism.
  - Internet Checking is the ability to use internet results in plagiarism detection.
  - Academical Checking is used to check the research publications and compare them to already published papers.
  - Multiple document comparison is the capacity of software to support multiple document comparison.
  - Multiple language support is the ability to support multiple language in document analysis.
  - Sentence Structure/synonymy show that software detection is capable to make sentence structure and synonymy analysis.
  - In our study, we include other parameters to evaluate the relevance of the plagiarism detection tools:
  - Types of plagiarism to detect is a feature which allows the selection of the type of plagiarism to be detected.
  - Machine learning means a machine learning model used in the approach.
  - Similarity based means if the software is based on matching techniques and similarity measurement.
  - Free license or not.
  - Size limitedness describes if the size of the document is limited (e.g. some tools limit the size document to 1000 words).
  - Document file is the file format to be analyzed (e.g. txt, pdf, docx, etc.).
- Classical methods use a corpus to extract the concepts, but recent researches rely on word embeddings techniques as Word2Vec, GloVe, BERT, to preserve the semantic and the syntactic context of the text.
  - Type of plagiarism detected presents whether if the software displays the types of plagiarism encountered or not and gives the rate each type checked.
  - Reports generation describes if the software exports the results as a report.
- 1) *The PlagAware tool*: It is an online tool [12] that uses a classic search engine to detect plagiarism and offers several reports helping the user to decide if the analyzed text is plagiarized or not. It is possible to add new database, to check documents from the internet results, and to compare multiple documents. Verifying sentence structure analysis and synonymy replacement is not supported. The languages supported are German, English and Japanese. It is used in universities to check the originality of the works to be published. PlagAware performs a complete scan of the document, and each sentence is analyzed to subsequently detect whether it contains plagiarism or not.
- 2) *The PlagScan tool*: It is an online tool used for academic plagiarism detection. This tool uses a local database that include millions of documents and includes the results of the internet search for making comparison. It supports adding a new database over the internet. It detects several types of plagiarism such as: copy and paste or words switching [14]. PlagScan supports the UTF-8 encoding languages and all Latin or Arabic languages. It is used in universities to check the originality of the works to be published. Sentence structure analysis and synonymy replacement are not supported. This tool uses a plagiarism detection algorithm that contains three consecutive word matches to subsequently detect plagiarism methods which use the replacement the words by their synonyms. In addition, they apply matching algorithms to detect documents similarity.
- 3) *The CheckForPlagiarism.net tool*: It is a tool for detecting academic plagiarism developed by a professional academic team. It can detect several types of plagiarism. It uses its own database that include millions of documents from several databases with different domains. It performs an internet Checking, Sentence structure analysis and synonymy replacement is detected, and it is possible to compare multiple documents. CheckForPlagiarism.net checks several types of documents, including, newspapers, PDFs, magazines, journals, books, articles etc. It supports several languages: Spanish languages, Portuguese, German, English, Korean, French, Italian, Arabic and Chinese languages. Each document is assigned by a fingerprint and used in document comparison [15].
- 4) *The iThenticate tool*: iThenticate is an online academic plagiarism detection tool for researchers, publishers and authors [16]. It possible in iThenticate to add a new database or use Internet for comparison and in addition it uses its own

database that contains several documents like books, newspapers and articles. Sentence structure analysis and synonymy replacement checking is not supported by iThenticate but it is possible to compare multiple documents. It supports more than 30 languages likes English, Russian, Arabic, etc. Many online scientific journals use it for submitted papers checking. iThenticate performs a matching advanced technics in similarity analysis highlight material within a manuscript that matches documents found in the iThenticate database algorithm to check the contents of a document against an extensive database of published scholarly writing [16].

5) *The PlagiarismDetection.org tool*: This is an online tool that is mostly used by teachers and students [17]. It used its own database that contains millions of documents, but it is possible to add a new database or use internet Checking. Sentence structure analysis and synonymy replacement checking is not supported neither multiple document comparison. It supports all languages using Latin characters. the technique is based on the n-gram method.

6) *The Urkund tool*: URKUND is a web plagiarism prevention system. Today, a vast majority of universities around the world use Urkund to effectively and detect

plagiarism. This system allows to compare the content of a document with several other resources from different sources (Internet, database, internal documents, etc.). Document formats accepted is doc, docx, pdf, etc. Urkund is multilingual, detects plagiarism by paraphrasing and replacements by synonyms, and returns the rate of similarity with the other documents [ 21].

7) *The Turnitin tool*: The Turnitin Plagiarism Detection System allows users to check their documents and compare them with web content and other documents that have already been downloaded by institutions as well as with certain journals [22]. For each submission, a report is produced identifying the sources of its similarities as well as the percentage of correspondence with the submitted document. Turnitin uses a matching algorithm to find strings of words within assignments that are identical to those within its repository.

### B. Comparison and Analysis

In this section we propose a qualitative comparison of plagiarism software detection. We focus on the features and properties of the tools rather than their performance in the first instance. Based on the comparison parameters cited above the results are reported in Table I.

TABLE I. COMPARISON OF THE PLAGIARISM DETECTION TOOLS

| Features                      | PlagAware | PlagScan | iThenticate | CheckForPlagiarism.net | Plagiarismdetecting.org | URKUND | Turnitin |
|-------------------------------|-----------|----------|-------------|------------------------|-------------------------|--------|----------|
| Add new database              | x         | x        | x           | x                      | x                       | ✓      | ✓        |
| Add a corpus for new training | x         | x        | x           | x                      | x                       | x      | x        |
| Internet Checking             | ✓         | ✓        | ✓           | ✓                      | ✓                       | ✓      | ✓        |
| Academical Checking           | ✓         | ✓        | ✓           | ✓                      | ✓                       | ✓      | ✓        |
| Multiple document comparison  | ✓         | ✓        | ✓           | ✓                      | ✓                       | ✓      | ✓        |
| Multiple languages            | ✓         | ✓        | ✓           | ✓                      | ✓                       | ✓      | ✓        |
| Sentence structure/ synonymy  | x         | x        | ✓           | ✓                      | x                       | ✓      | ✓        |
| Types of plagiarism           | x         | x        | x           | x                      | x                       | x      | x        |
| Machine Learning              | x         | ✓        | x           | x                      |                         | ✓      | x        |
| Free License                  | x         | ✓        | ✓           | ✓                      | x                       | x      | ✓        |
| Similarity based              | ✓         | x        | ✓           | x                      | x                       | x      | ✓        |
| Size limitedness              | ✓         | ✓        | ✓           | ✓                      | ✓                       | x      | x        |
| All type of files             | ✓         | ✓        | ✓           | ✓                      | ✓                       | ✓      |          |
| Classical method              | ✓         |          | ✓           | ✓                      | ✓                       | x      | ✓        |
| Type of plagiarism detected   | x         | x        | x           | x                      | x                       | x      | x        |
| Reports generation            | ✓         | ✓        | ✓           | ✓                      | ✓                       | ✓      | ✓        |

From Table I, we can see that all studied plagiarism detection tools can perform Internet Checking to verify if there is any similarity with any resources on internet. Also, the document analyzed can be written in Multilanguage. These systems are almost used in the Academical context to check student reports, thesis, or research papers. Multiple documents comparison is also provided by these tools. But as we see, most of them does not a have the feature of adding a new corpus. This new feature enables adding a corpus to be used as the basic dataset for the plagiarism detection step. It is an opportunity to use more corpus for improving the learning phase. The new corpus contains a source document, suspicious documents and the type of plagiarism. As we can see in Table I, none of the analyzed tools specify the type of plagiarism that has been detected from sources, nor give the user the possibility to specify the type of plagiarism he wants to be detected. Based on this comparison and to benefit from our previous work [18], we propose an implementation with new features to deal with a plagiarism in textual documents. In the next section, we describe the background of the approach and its components and the services that our framework can provide.

### III. GENERAL FRAMEWORK OF OUR APPROACH

The proposed plagiarism detection tool is based on our previous research validated with PAN Dataset where data are labeled with the types of plagiarism [20]. Fig. 1 represents a global architecture of our framework which is based on two Deep learning architectures Siamese Long Short-Term Memory (SLSTM) and Convolutional Neural Networks (CNN).

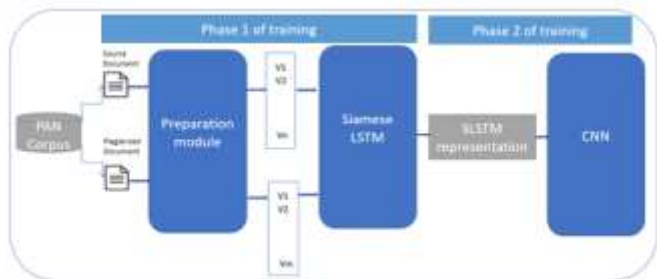


Fig. 1. Global Architecture.

The approach based mainly on three steps as described below:

- **Context representation of documents:** The corpus consists of a set of source documents and a set of suspicious documents plagiarized from each source using a specific kind of plagiarism. Both of sources and the plagiarized document are transformed with doc2vec a list of sentences vectors to be used as input to the SLSTM model.
- **SLSTM Learning phase 1:** The Siamese LSTM is used to learn the different kind of plagiarism in dataset.
- **CNN Learning phase 2:** The output of the first stage is a SLSTM representation of documents. To consolidate our approach, we used CNN deep learning model to detect the types of plagiarism learned in the first part of

the approach. The goal is the classification of the document as plagiarized or not with the type of plagiarism detected in it if yes.

### IV. DEEP LEARNING PLAGIARISM DETECTION SYSTEM

In this section we will present the proposed plagiarism detection framework by illustrating the technical architecture and its different layers. Fig. 2 shows an overview of the proposed system. The system is composed by the following six layers: Front-end Layer, http layer, Controller Layer, preprocessing layer, Learning layer and Detection Layer. Here bellow, we present the description of each layer and its implementation.

#### A. Front End Layer/ Http Layer/ Controller Layer

The Front end is a platform for building mobile and desktop web applications that communicate with the http layer which offers web services to consume. The flask package provides some classes to build a Service layer and exposes an API that interacts with the model. The first idea is to remove all logic of the routes and model of the Flask application and put it in the service layer. The second goal is to provide a common API that can be used to manipulate a model regardless of its storage backend. The controller layer concerns a middleware between the flask layer and the other layers of our system.

#### B. Preprocessing Layer

At first, the corpus is preprocessed as shown in Fig. 3. For each document we realize the cleaning, segmentation and stemming phases [18]. Then the output is given as input to the doc2vec word embedding model layer.

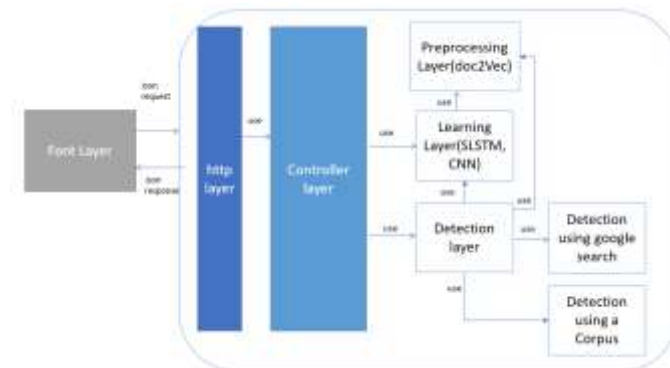


Fig. 2. The Distribution of Application Layers.

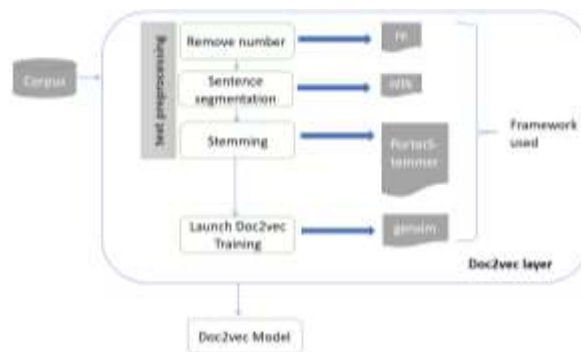


Fig. 3. Data Preprocessing and Vectorization.

Then we launch the training phase to generate a doc2vec model which we will use later to transform each sentence of a document to a vector. We worked with the re framework to build a regular expression that removes numbers, nltk to segment a document by sentence, PorterStemmer to apply the stemming principle that makes a word in the initial form and gensim to start training the doc2vec model.

C. Learning Layer

In this layer, we applied twice the learning process as shown in Fig. 4. In the first step, we used SLSTM algorithm for learning from the output of doc2vec and the output is given to CNN Model to learn again to build our efficient learning model. At the end of this phase we restore the SLSTM model which will be used to test whether a pair of documents are similar or not and we also get the SLSTM representation. In this step we used the keras tensorflow.

To carry out the classification of documents and add the types of plagiarism that have been detected, we used the keras tensorflow to build our CNN model. Hence, the outputs of the SLSTM model are used in the second learning phase which consists of classifying the types of plagiarism already learned in the first part.

D. Detection Layer

For document classification task (whether is plagiarized or not), the users can make choice to use a new corpus, internet search results or a new corpus for comparison. The corpus contains a list of sources documents that will be used in learning step or to search for similarity with the text to be verified. The second option uses python Google package to get the link of the first n search results and compare the text analyzed with the contents of these links. More details will be given in the next subsections.

1) Add a new corpus: To add a new corpus, we respect the process in Fig. 5. Firstly, the user adds the pairs file, which is a text file that contains several lines and each line represents a type of plagiarism. Secondly, the user uploads the corpus containing the source and plagiarized document mentioned in the pair document above.

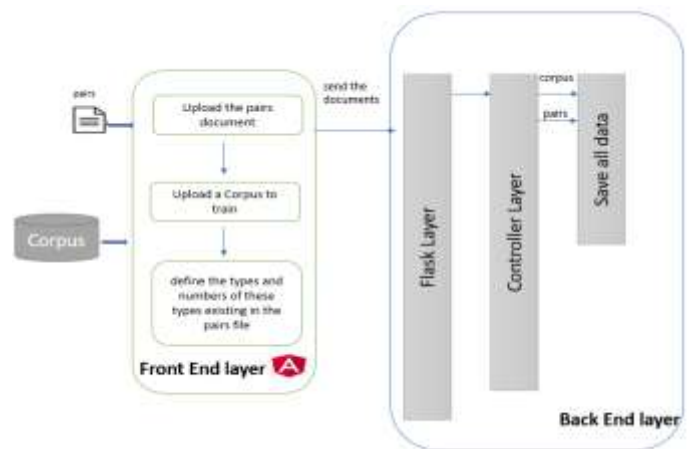


Fig. 5. Add New Corpus.

Finally, the user defines the types and numbers of plagiarism cases. But the number of plagiarism types entered must correspond to the number of lines existing in the pair file. After adding a new corpus, we can launch the training phase which follows the process in Fig. 6.

2) Add a new corpus for comparison: Our framework can also compare a document to a special corpus containing a set of desired source documents to compare with. We must first add corpus which will be the basis of comparison, and the system will compare the document to each document in corpus to detect a kind of plagiarism. Fig. 7 presents the process of this task. The comparison is carried out by using the following steps:

- Select corpus trained and corpus of comparison.
- Segment the analyzed document to a list of paragraphs.
- Retrieve a list of paragraphs for each document in corpus of comparison.
- Using our deep learning system, we compare each paragraph of the analyzed document with all the paragraphs returned via the corpus of document.

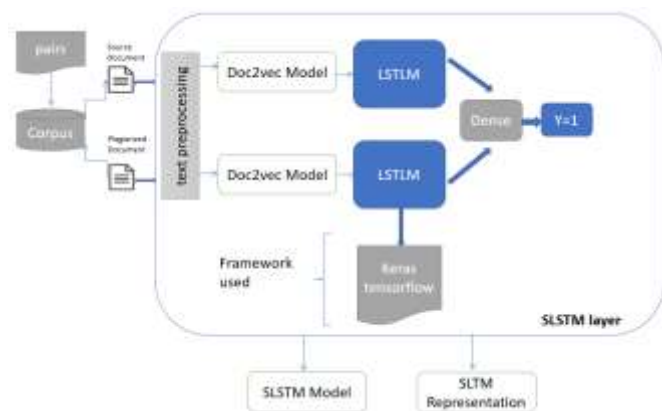


Fig. 4. Learning with SLSTM Model.

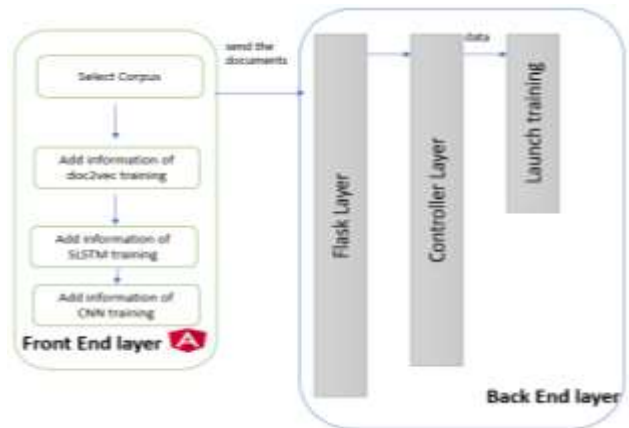


Fig. 6. Launch Training.

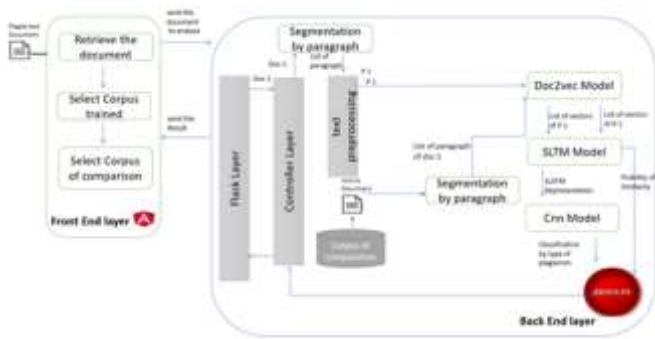


Fig. 7. Plagiarism Detection from Corpus Process.

3) *Using google research engine:* Our system can also detect the plagiarism in documents using google search result as illustrate in the following Fig. 8.

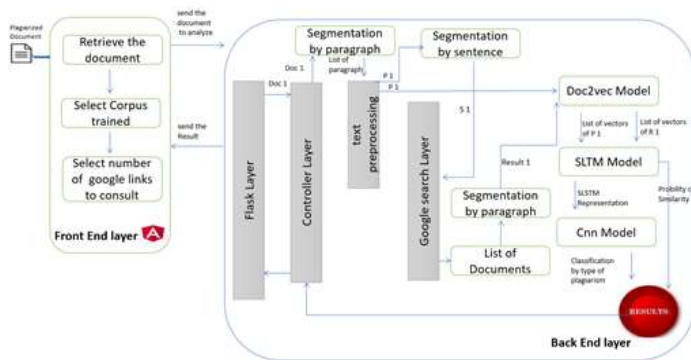


Fig. 8. Check Plagiarism from Internet.

The plagiarism detection consists of the following steps:

- Segment the analyzed document to a list of paragraphs and the list of sentences.
- Use the sentences in this paragraph to retrieve the various links which contain the suspected texts.
- Retrieve a list of paragraphs for each link found.
- Using our deep learning model to compare each paragraph of the analyzed document with all the paragraphs returned via the Google searches.

More precisely we assume that the document contains  $N$  paragraphs, if for example the first paragraph contains  $S$  sentences, so we launch  $S$  internet search to retrieve  $S \times N$  result then we assume once again that each result will offer us  $P$  paragraphs which are considered as suspected initials. So, the first paragraph of the analysis document is compared with  $N \times S \times P$  paragraph.

## V. EXPERIMENTS

In this section, we present different possibilities that our system provides in terms of plagiarism detection. We can proceed three kind of comparison: Two text comparison, online comparison and using an intern corpus for a comparison.

### A. Add New Corpus

For this task we proposed the following IHM in Fig. 9:

Fig. 9. Add New Corpus.

### B. Training a New Corpus

To do that, we must fill some information about corpus, doc2vec training, SLSTM training and finally CNN training. The data requested are used to develop the accuracy rate of our training. For this phase we proposed the following IHM in the Fig. 10.

Fig. 10. IHM to Launch New a New Training.

This part contains hyperparameters used to adjust the three models in learning process, for more information see [19][14].

### C. Comparison of Two Texts

Given two documents, we can make a comparison of two given documents by following the steps in Fig. 11 and Fig. 12. The two documents will be preprocessed and converted to a list of vectors with doc2vec model. The system will detect later if the input documents are similar or not using SLSTM Model and it will report the probabilities of each kind of plagiarism trained in our system when we use CNN Model. Fig. 13 provides an example of two documents comparison.



Fig. 11. Two Documents Comparison.

Result

Value of Similarity: 0.1246744

| Type                           | Similarity |
|--------------------------------|------------|
| NUMBER_NO_PLAGIARISM           | 0.1137468  |
| NUMBER_NO_OBFUSCATION          | 0.1115743  |
| NUMBER_RANDOM_OBFUSCATION      | 0.1124852  |
| NUMBER_TRANSLATION_OBFUSCATION | 0.1117394  |
| NUMBER_SUMMARY_OBFUSCATION     | 0.113384   |
| NUMBER_NO_PLAGIARISM10         | 0.1133344  |
| NUMBER_HUMAN_RETELLING         | 0.1118285  |
| NUMBER_SYNONYM_REPLACEMENT     | 0.1138874  |
| NUMBER_CHARACTER_SUBSTITUTION  | 0.1130215  |

Fig. 12. Results of Comparison of Two Documents.

And we get the result below which result contains the probability of similarity between these two texts, in fact, we also recover the probabilities of each type of plagiarism learned at the training phase.

#### D. Online Checking

For performing plagiarism detection from documents returned from Google research engine, we need to fix several parameters as the learned corpus, number of sites to consult and finally the text to analyze, as mentioned below it proposed an IHM in Fig. 13.

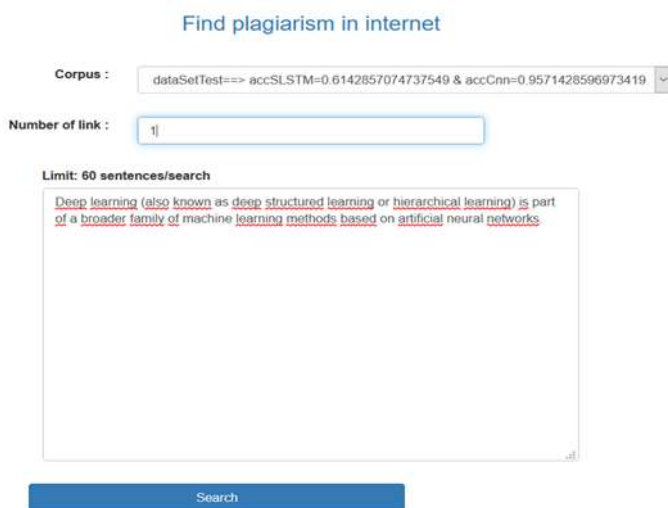


Fig. 13. Plagiarism Detection from Internet IHM.

The results in Fig. 14 shows the source text, link of the source text, the probability of similarity. In the right, the table presents the probabilities of each type of plagiarism learned in the training phase present in the document.

#### E. Using Corpus for a Comparison

Fig. 15 below represents the result of detection of plagiarism using a corpus of source documents instead of consulting the results of the internet. The results consist of a list of blocks containing the following information:

- the paragraphs analyzed.
- the name of the source document.
- probability of similarity.

In addition, we propose a table in Fig. 16 containing the probabilities of each type learned in the training phase.

Result

See also: the source document, learning.html  
[www.dailymail.co.uk/Science/Article-3262276/Deep-learning.html](#)

Value of Similarity: 0.882285

Text to analyze: Deep learning also known as deep structured learning or hierarchical learning is part of a broader family of machine learning methods based on artificial neural networks. Learning can be supervised, semi-supervised or unsupervised.

Original text (HTML) used learned by information processing and distributed communication nodes in biological systems. After how various differences from biological

| Type                           | Similarity |
|--------------------------------|------------|
| NUMBER_NO_PLAGIARISM           | 0.2223422  |
| NUMBER_NO_OBFUSCATION          | 0.2145277  |
| NUMBER_RANDOM_OBFUSCATION      | 0.2223422  |
| NUMBER_TRANSLATION_OBFUSCATION | 0.2223422  |
| NUMBER_SUMMARY_OBFUSCATION     | 0.2223422  |
| NUMBER_NO_PLAGIARISM10         | 0.2223422  |
| NUMBER_HUMAN_RETELLING         | 0.2223422  |
| NUMBER_SYNONYM_REPLACEMENT     | 0.2223422  |
| NUMBER_CHARACTER_SUBSTITUTION  | 0.2223422  |

Fig. 14. Example of Plagiarism Detection from Internet.



Fig. 15. Plagiarism Detection from Corpus IHM.



Fig. 16. Example of Plagiarism Detection from Corpus.

#### F. Proposed System Features

In comparison with existing systems, our plagiarism detection system has all the properties used in the comparison above. We have added new features making it an able to make followed action:

- Upload and Add Any Dataset.
- Add New Corpus for Training Plagiarism.
- Internet Checking.
- Academical Checking: We can add the corpus of publication or get them through Google result.
- Two documents comparison but it could be extended to more than two.
- Multiple languages detection: We can use any language, but you must choose the corpus already trained by this language.
- Check all type of plagiarism.
- Personalize the types of plagiarism to detect: We can define several kinds of plagiarism in our training phase.
- Use the deep learning approaches: our approach uses deep learning algorithms.
- Document size is limited: not limited.
- Show the type of plagiarism detected: Yes.
- Reports generation: Yes.

#### VI. CONCLUSION

In this paper, we proposed a new system for the detection of plagiarism based on the deep learning methods. Its interest is the extraction of characteristics without losing the sense of the document by using doc2vec word embedding technique. The proposed system has the ability to detect not only that there is plagiarism but also the probabilities of the existence of each type of plagiarism. We presented the different functionalities offered by our system, either at the level of the personalized learning phase or the different ways of detecting

plagiarism offered. Compared to the other tools studied in this paper, our proposition offers more functionalities as adding and training new corpus or using a special corpus for comparison. As for our perspectives, we will improve the various interfaces of the application to make it more accessible to the general public and improve the response time due to the learning time. It would also be interesting to compare the performance of different approaches in a quantitative way.

#### REFERENCES

- [1] Risquez, A., Dwyer, M. O.; Ledwith, A. (2011). «Thou shalt not plagiarise»: from self-reported views to recognition and avoidance of plagiarism». *Assessment & Evaluation in Higher Education*, vol. 2, no. 1, p. 34-43. <http://doi.org/10.1080/02602938.2011.596926>. 3 Ruipérez, G.; García-Cabrero, J.C. (2016). «Plagiarism and Academic Integrity in Germany». *Comunicar*, vol. 24, no. 48, p. 9-17. <http://doi.org/10.3916/C48-2016-01>.
- [2] Liddell, J. (2003). *A Comprehensive Definition of Plagiarism*. *Community & Junior College Libraries*, 11(3), 43–52. doi:10.1300/j107v11n03\_07.
- [3] Thomas Lancaster, Fintan Culwin. *A Visual Argument for Plagiarism Detection using Word Pairs*. School of Computing University of Central England Perry Barr Birmingham B42 2SU United Kingdom. Faculty of Business, Computing and Information Management London South Bank University Borough Road London SE1 0AA. *Plagiarism: Prevention, Practice and Policies 2004 Conference*.
- [4] Zubarev D.V. Sochenkov I.V. *Paraphrased plagiarism detection using sentence similarity*. Federal Research Center “Computer Science and Control” of Russian Academy of Sciences, Moscow, Russia.
- [5] Maxim Mozgovo, Tuomo Kakkonen, Georgina Cosma. *Automatic student plagiarism detection: future perspectives*. university of aizutsuruga, ikki-machi, aizu-wakamatsu, fukushima, 965-8580 japan. *articleinjournal of educational computing research* · january 2010.
- [6] Sousa-silva, r. -detecting translanguing plagiarism and the backlash against translation plagiaristslanguage and law / linguagem e direito, vol. 1(1), 2014, p. 70-94.
- [7] Eman s. Al-shamery, Hadeel qasem Ghani. *plagiarism detection using semantic analysis*. published 2016. *computer science indian journal of science and technology*. doi:10.17485/ijst/2016/v9i1/84235 corpus id: 55709933.
- [8] Roumiana Peytcheva-forsyth, Harvey Mellar, Lyubka Alekseiva. *using a student authentication and authorship checking system as a catalyst for developing an academic integrity culture: a bulgarian case study*. *journal of academic ethics* 17, 245-269(2019).
- [9] Amruta Patil, Nikhil Bomanwar. *survey on different plagiarism detection tools and software’s*. computer department, mumbai university. (ijcsit) *international journal of computer science and information technologies*, vol. 7 (5) , 2016,2191-219.
- [10] Hage, J. rademaker, P Vugt, n. *a comparison of plagiarism detection tools*, tech-nical report uu-cs-2010-015,june 2010, department of information and comput-ing sciences utrecht university, utrecht, the netherlands.
- [11] G. Bela, “state-of-the-art in detecting academic plagiarism. *international journal for educational integrity*,” vol. 9no.1 june, 2013 pp. 50-71 issn 1833-2595.
- [12] <https://www.plagaware.com/>(lastaccessfebruary 1, 2020).
- [13] Asim m. El tahir Ali, Hussam m. Dahwa Abdulla, and Vaclav Snasel. *overview and comparison of plagiarismdetection tools*. department of computer science, v̄sb-technical university of ostrava,17. listopadu 15, ostrava - poruba, czech. dateso 2011, pp. 161–172, isbn 978-80-248-2391.
- [14] [http://www.plagscan.com.](http://www.plagscan.com/) (last access february 1, 2020).
- [15] <http://www.checkforplagiarism.net> (last access february 1, 2020).
- [16] <http://www.ithenticate.com/index.html> (last access february 1, 2020).
- [17] <http://www.plagiarismdetection.org> (last access february 1, 2020).



- [18] El mostafa Hambi and Faouzia Benabbou. a multi-level plagiarism detection system based on deep learning algorithms. information technology and modeling laboratory faculty of sciences ben m'sik casablanca, morocco. international journal of computer science and network security. vol. 19 no. 10 pp. 110-117.
- [19] Hambi el Mostafa , Faouzia Benabbou. a system for ideas plagiarism detection: state of art and proposed approach. misc2018. information technology and modeling laboratory science faculty ben m'sik casablanca, morocco.
- [20] Martin Potthast, Matthias Hagen, Anna Beyer, Matthias Busse, Martin Tippmann, Paolo Rosso, and Benno Stein. overview of the 6th international competition on plagiarism detection. web technology & information systems, bauhaus-universität weimar, germany natural language engineering lab, universitat politècnica de valència, spain. at pan 2014.
- [21] <https://www.orkund.com/the-orkund-system/>.
- [22] Dr. Karl o. Jones. practical issues for academics using the turnitin plagiarism detection software. international conference on computer systems and technologies.

# Impact of Project-Based Learning on Networking and Communications Competences

Cristian Castro-Vargas<sup>1</sup>, Maritza Cabana-Caceres<sup>2</sup>, Laberiano Andrade-Arenas<sup>3</sup>

Faculty of Sciences and Engineering  
Univerisdad de Cienciasy Humanidades  
Lima Perú

**Abstract**—The objective of this article is to establish the impact of project-based learning on networking and communication competences I in engineering students from Lima city. The study was of an applied type, quasi-experimental design, and was made up of a population conformed by 39 students of the VI cycle of engineering, an objective test was applied to measure the impact of project-based learning on network and communication competences I. The research results determined the statistically significant relationship of project-based learning and networking and communications competences I in engineering students with pretest values of  $Z = -$ , 498 greater than -1.96 (critical point) and level of significance  $p$ -value = 0.618 greater than  $\alpha = 0,05$  ( $p > \alpha$ ) and then with values in the posttest of  $Z = -$  4,488 less than -1.96 (critical point) and level of significance  $p$ -value = 0.000 less than  $\alpha = 0.05$  ( $p < \alpha$ ), therefore the project-based learning positively and significantly impacts on network and communication competences I, in engineering students, supporting the alternative hypothesis and rejecting the null hypothesis. Consequently, we reached the conclusion that the application of the project-based learning methodology has demonstrated that caused a positive and significant impact on network and communications competences I in engineering students.

**Keywords**—Project-based learning; competencies; networking; communications; network convergence

## I. INTRODUCTION

The rapid changes that occur in the job sector as a result of globalization do not find future engineering students with the skills in networks and communications for an adequate professional development. Higher educational programs do not correctly manage the development of professional competences holistic or comprehensive [1]. The Organization for Economic Co-operation and Development (OECD) indicated that professionals have difficulties in performing properly in the company because they are not aware of the importance of development in network skills and communications, generating not enough contributions in society which worldwide in the coming years will need millions of dozens of jobs that will require that they have competences in the specialty to solve problems in this field [2]. Furthermore, each network technology manufacturer develops new implementations of Network and Communications Convergence, generating the need for more trained professionals in these new technological competences. According to the United Nations Educational, Scientific and Cultural Organization UNESCO [3] mentioned

that all this leads to pressure on academic institutions of higher education, whether in the public or private sector.

In Latin America [4] they were classified as one of the problems, is that not all students have the availability to adapt to work academic-training activities in a collaborative way, although the student team tries to include them, they have rigid behavior when presenting tasks individually, they do not accept changes to new situations and also they cannot communicate their ideas.

In the case of Perú, the concern of the Organization for Economic Co-operation and Development (OECD) [5] found that the student in his university training does not have the competences in the adequate communications network to be able to cover jobs competitively. This is worrisome when considering that overall the penetration of some Information and Communications Technologies (e.g., the Internet) have been gradually increasing in this Latin American nation [6]. Due to this he does not adapt quickly to technological changes, its development is insufficient, and it is vulnerable to inserting itself in the labor population, which is why the National Superintendence of Higher Education SUNEDU [7] mentioned that it is still due to the disorder that exists between the labor markets and higher education, it generates a complex academic problem due to the inequality between the information that is handled in the university that does not contrast with the problems that the market asks to solve, reflecting this in his professional life through inconsequential work jobs. In the university, instead of the research, the teaching of specific and formative knowledge prevails, but it is not enough to achieve an integral formation of the students who agree to achieve competitiveness, whether national or international. Therefore, the research problem is focused on strengthening the student competencies, which in some cases it has not been prioritized in the course of networks and communications.

In Brazil, a study evaluated 20 students for every 4 semesters, with a maximum reference of 10, all the students obtained a score of 8 at the end. As a conclusion of the research, it was found that when carrying out the implementation of the network infrastructure, the experimental group evolved better because they oriented the activities to projects with real problems, they also found a greater commitment to meet the challenges, leading them to search and analyze the best possible solution and understanding of the functionality of the correct design that a data network should have, thus also leading them to have a greater social awareness for the participation in carrying out the projects [8].

In the USA, Rice University, Houston, Texas, researchers carried out a quasi-experimental study, of a student sample equal to 5492 of science and engineering, the participants rated the extent to which they were going to be effective in executing STEM tasks on a Likert-type scale of 6 items, descriptive statistics, correlation  $r = 5.194$ ,  $p < 0.001$  for STEM efficacy competences. It was concluded that the university investment in active learning activities such as PjBL will pay off by increasing student participation and interest in the STEM career [9].

In Spain, a study considered a sample of 50 students and 1 teacher. Its objective was to teach students the importance of the network course, to enhance learning and the acquisition of basic competences of the subject, through project-based learning to achieve better performance of the competences of the subject of networks, to establish the level of participation an evaluation of self-correction mechanisms of the projects was applied to the students. It was concluded that the realization of projects is positive for the learning of communication networks, because it allows them a better acquisition of competences, as well as gaining previous experience for professional performance, verifying with the growth of 70% of good level of the students in the course at the end of the course [10].

In China, a research work was conducted with a population of 80 students. It was concluded that project-based learning not only cultivates student initiative, but also improves their collaborative, practical, and project planning capabilities, as well as the student recognizes how extracurricular time can be used well with this system [11].

The National University of San Agustín de Arequipa, considered a population of 74 students distributed in 4 groups, carried out 9 sessions, aimed to apply the project-based learning methodology to the engineering course to verify better achievement in their professional training competencies. The project-based learning methodology was concluded, which strengthens and increases knowledge, as well as improves the competences, abilities and attitudes of understanding the course, with 72% between the level of good and excellent as a final result [12].

The Universidad Peruana Cayetano Heredia, carried out a study of a quantitative approach and quasi-experimental design, of a population of 76 students, aimed to check the effect of using the virtual platform to improve competences and learn the network course and communications, obtaining as results a significant improvement verifying in the notes at the beginning of 9.24 (65%) in the pretest and at the end of 15.6 (90%) in the posttest of the students, having as a conclusion that employment of the active methodology improves the level of competences in the course of networks and communications [13].

Likewise, at the César Vallejo University, it considered a population equal to 158 students, a non-probabilistic, intentional sample, 57 from the group belong to the experimental group and 60 belong to the control group. I do 8 sessions with the project-based learning methodology for the experimental group. Using a questionnaire instrument of 34 consultations divided into seven dimensions where the level of

investigative competences was evaluated, showing improvements due to the use of project-based learning. The non-parametric Mann-Whitney U test was applied, determining the statistical significance  $p = 0.00$ . Concluding that the application of project-based learning allowed a significant improvement in the level of investigative competences of the experimental group [14].

The San Luis Gonzaga National University of Ica, with a population of 80 students, considered two groups of 40 equal for the experimental and the control, whose objective was to establish the influence of project-based learning on the competencies in systems engineering students. In conclusion, collaborative project-based learning improved the acquisition of competences in Engineering students, as well as a greater ability to: conception, design, development, effective use of engineering techniques and tools, reflected in the notes. The pretest the mean was 13.98 and in the end in the posttest improvements were obtained with a mean 17, achieving a  $p$  value equal to 0.00 [15].

At the National University of Engineering, a census sample of 36 students of the VI cycle of engineering was considered, divided into two groups: 17 control and 19 experimental, and performed the non-parametric U test of Mann Whitney, whose objective was to determine the effect of the application of a sustained program in project-based learning methodology in the development of competencies in students. Concluding the experimental group, he developed the procedural and attitudinal cognitive competences of the course, applying the knowledge acquired in the execution of a project similar to an activity in the professional field, achieving a  $p$ -value = 0.000 [16].

For the variable impact of project-based learning, constructivism was considered as a theoretical approach [17] because it refers to the use of learning resources associated with appropriate scientific methods or learning models, for which this model allows the allocation of time For the design of learning resources, the project-based learning methodology thus generally allows for better development in the institution's learning-teaching processes. Constructivism [18] according to Piaget's theory, indicated that there is a close similarity between project-based learning and constructivism, because there is an active process, permanent participation among students going through various experiences to try to contrast solutions to the problem that normally occurs in a real environment. Project-based learning is a pedagogical approach due to the following characteristics: Students are active during learning through cooperative participation in scientific and engineering practices, students create a set of tangible products and shared artifacts that are accessible to the public to external representations [19].

In this regard, project-based learning is developed by the theory of active learning which is constructivism, and is a learning model that allows to increase learning in students; the habit and the creation of new learning practices, because students have to originally think for the real-life problem, the solution, centered on the student as it allows constant discovery [20]. In the same sense, it encourages the development of transversal competences and promotes autonomy [21]. To

obtain a good learning requires learning environments of higher education conducive to a better development of metacognitive competences [22].

Regarding the theoretical approach of Network and Communications Competences I, it is considered belonging to constructivism. According to Piaget, he maintains that the student builds his own knowledge through the experience of carrying out the assigned activities, obtaining practical results, using the approach of constructionism, which maintains that it is an action-based learning, which allows the student to have a mental structure related to the concrete action, which will allow him to generate a motivation for learning, thus obtaining a correct intellectual construction given by the exchange and the assigned work that is being developed, also generating intellectual and affective autonomy. Another theoretical consideration is the connectivism approach, according to George Siemens [23].

Regarding the definition of the dependent variable network and communications competences I, according to the curriculum EAP Professional Academic School of Systems Engineering, defined it as the curricular experience of networks and communications I belonging to the professional training area and it is theoretical - practical compulsory, with the purpose of creating in the student the competences in design, implementation and administration of computer networks, using both information and communication technologies for the development of the following aspects: Networking, switching, routing and wan technologies [24]. Hence, competence was defined as "an ability to act effectively in a defined type of situation, an ability that relies on knowledge but is not reduced to it" [25].

The computer network was defined as "a communication system that connects two or more computers to each other by means of such optical fibers, radio waves, and electromagnetic waves, allowing them to share information and services." [26]. Regarding the definition of a communication network, it is necessary to interconnect different computer systems throughout the wan, using internet [27]. Academic competencies are the essential knowledge that is learned during general training and are classified into: writing, problem solving, mathematics, reading ability, creative thinking, assertive communication, decision making, assimilation and comprehension, learning and reasoning [28].

The curriculum of the Academic Professional School EAP Systems Engineering dimensioned the competences of networks and communications I: knowledge to solve problems of local area network design, implementation of network convergence and network administration [24]. Theoretically, project-based learning and network and communications competencies are justified. Both have an affinity because they allow the use of characteristics between them such as analysis, planning and design, guiding them to the correct development of similar problem solutions to the professional field.

Thus, there is also an epistemological justification between project-based learning and the competences of the course on networks and communications I, which is given such justification because they refer to the theory of constructivism, which affects having the student as a direct participant in their

learning within the problems that they must understand in order to build or guide concepts to solve said problem [24].

The purpose of this research is to determine the impact of project-based learning on the competences of networks and communications I, in engineering students, which is dictated the sixth cycle of the specialty of Engineering, of Systems of the University of Sciences and Humanities.

In Section II we present the methodology followed to obtain the results shown in this study, in Section III presents our results, in Section IV we present the discussion, in Section V we present the conclusions and finally in Section VI we present the recommendations.

## II. METHODOLOGICAL FRAMEWORK

The present study is based on the positivist paradigm. Thus, the present investigation predicts and controls the phenomena thus verifies theories, being the external observer investigator of the events, which is based on the positivist paradigm [29].

The approach is quantitative because it is the concrete expression of reality, for which it considers as a means of reaching knowledge, explanation, prediction, control of the phenomenon, verification of theories [29]. Because they use criteria such as: validity, objectivity, reliability, statistics, experimentation, used in questionnaires, tests and then data analysis using inferential and descriptive statistics.

The level of investigation is explanatory. The objective is to verify the causal hypotheses, following a sequence of processes or organized steps [30].

The type of research is applied because it is based on the theoretical framework and is practical. The research design is quasi-experimental, quasi-experimental designs come to be strategies for conducting research in order to evaluate the impact of the set of methods carried out, for this purpose it considers two groups, one control and the other experimental through a pretest and at the end of the process the posttest evaluation, in order to demonstrate the hypothesis raised at the beginning [31].

### A. Variable Operationalization

Regarding the conceptual definition of the dependent variable: Network and communications competences I, according to the EAP Professional Academic School of Systems Engineering, networks and communications curricular experience I belongs to the professional training area and is theoretical-practical compulsory, with the purpose of creating in the student the competences in design, implementation and administration of computer networks, using both information and communication technologies for the development of the following aspects: Networking, switching, routing and wan technologies (Internetworking) [24].

To operationalize the impact variable of project-based learning, an organizational matrix was developed, distributed in 10 sessions of 3 academic units (Table I).

For the operationalization of the dependent variable network and communications competences I, a 20-question questionnaire was prepared, on a dichotomous scale containing the indicators of the three dimensions: knowledge to solve

local area network design problems, implementation of the convergence of the network and network administration (Table II).

**B. Population**

The population comes to be a finite set of people or objects that have common characteristics, susceptible to a study and later allow to replicate the findings in other populations [30].

TABLE I. ORGANIZATION MATRIX OF THE INDEPENDENT VARIABLE LEARNING BASED ON PROJECTS

| Units                     | Strategic activities   | Phases                     | Indicators   |
|---------------------------|--|----------------------------|--|
| Local area network design | Session 1: Local area networks<br>Session 2: Structured wiring<br>Session 3: Local area networks wireless<br>Session 4: Wireless Security  | 1. Key question definition | <ul style="list-style-type: none"> <li>Event or occurrence.</li> <li>Spontaneous interest in the student</li> <li>Community proposal</li> <li>Values educational power</li> <li>Causes a commitment</li> <li>Responds to your interests</li> <li>Social relevance</li> <li>Listening and creative attitude</li> <li>Project appointment</li> </ul> |
| Network convergence       | Session 5: IP Telephony. Voice digitization and coding<br>Session 6: IP Telephony Servers.   | 2. Work plan-calendar      | <ul style="list-style-type: none"> <li>Interaction</li> <li>Motivation</li> <li>Collective thinking</li> <li>Research lines</li> <li>Different strategies</li> <li>Different itineraries</li> <li>Decide what to do</li> <li>Feel capable</li> <li>Timeline</li> </ul>   |
| Network administration    | Session 7: Fundamentals of LAN Switches.<br>Session 8: WAN Technologies. Routers Basics.<br>Session 9: Monitoring of end devices.<br>Session 10: Monitoring of intermediary devices. | 3. Follow-up monitoring    | <ul style="list-style-type: none"> <li>Learning moves to action</li> <li>Multiple intelligences</li> <li>Cooperation</li> <li>SWOT tool</li> <li>Creative tools</li> <li>Collect materials</li> <li>Individual and group self-assessment tools</li> </ul>  |
|                           |  | 4. Evaluation              | <ul style="list-style-type: none"> <li>I've learned</li> <li>What is it for</li> <li>What do I do with it?</li> <li>Tools (portfolio, learning journal, rubric, assessment notes)</li> </ul>   |

TABLE II. OPERATIONALIZATION MATRIX OF THE DEPENDENT VARIABLE COMPETENCES OF NETWORKS AND COMMUNICATIONS I

| Dimensions   | Indicators   | Item           | Value scale  | Dimension levels  | Levels    |
|--|--|----------------|--------------|---|-----------|
| 1. Knowledge to solve local area network design problems | 1.1. Configure PCs to share resources                        | 1<br>2         |              | Bad: [00 - 03]<br>Regular: [04 - 06]<br>Good: [07 - 08] | Bad:      |
|  | 1.2. Connect the work area wiring                            | 3<br>4         |              |   |           |
|  | 1.3. Implement a small wireless network                      | 5<br>6         | Dichotomous  |   | [00 - 10] |
|  | 1.4. Configure the security of a wireless local area network | 7<br>8         |              |   |           |
| 2. Implementation of the network convergence             | 2.1. Configure IP phones and softphones                      | 9              |              | Bad: [00 - 01]<br>Regular: [02]<br>Good: [03]           | Regular:  |
|  | 2.2. Configure the services of an IP telephone exchange      | 10             | 0: Incorrect |   | [11 - 15] |
|  | 2.3. Install the softPBX                                     | 11             |              |   |           |
| 3. Network administration                                | 3.1. Basicly configure the Cisco switch                      | 12<br>13<br>14 |              | Bad: [00 - 03]<br>Regular: [04 - 06]<br>Good: [07 - 09] | Good:     |
|  | 3.2. Basically configure the Cisco Router                    | 15<br>16<br>17 | 1: correct   |   | [16 - 20] |
|  | 3.3. Install and operate end device monitoring tool          | 18<br>19<br>20 |              |   |           |
|  | 3.4. Monitor the events of the Switches and Router.          | 21             |              |   |           |

The population considered comes to be 39 Engineering students, from the Private University of the engineering career, who develop the academic cycles of the VI cycle of the network and communications course I (Table III).

TABLE III. POPULATION

| No. | Sections                          | Males | Ladies | Totals |
|-----|-----------------------------------|-------|--------|--------|
| 1   | Sixth Cycle A1 Experimental Group | 14    | 5      | 19     |
| 2   | Sixth Cycle C1 Control Group      | 15    | 5      | 20     |
|     | TOTAL                             |       |        | 39     |

### C. Techniques, Data Collection Instruments, Validity and Reliability

The technique used was the survey, which allowed data to be obtained on the variable under study network and communications competences I. Thus, the instrument used was to collect data on a knowledge test of network and communications competences I.

The factor analysis to measure network and communications competences I through its 20 items, is used to test whether the items that make up each factor can generate correspondence between the dimensions proposed, obtaining that the KMO value is equal to .610, which indicates that there is a relationship between the values reached and the chosen sample (Table IV).

To establish reliability, a pilot test was carried out to verify the reliability of the research instrument, to a population of 30 students, which had to eliminate some items, to improve the reliability of the instruments. Once the correction was made, the instruments were applied again to the study population, obtaining a KR-20 value equal to 0.860 for the variable network and communications competences I, verifying high reliability (Table V).

TABLE IV. ADAPTATION ANALYSIS TO THE FACTORY ANALYSIS

| Statistical                                     | Value               |         |
|---|---------------------|---------|
| Kaiser-Meyer-Olkin measure of sampling adequacy | 0.61                |         |
| Bartlett's sphericity test                      | Approx. Chi squared | 345,434 |
|   | gl                  | 190     |
|   | Sig.                | 0       |

TABLE V. INSTRUMENT RELIABILITY

| Variables                                   | Reliability Statistics | Value | No. of elements |
|---|------------------------|-------|-----------------|
| Networking and communications competences I | Kuder-Richardson       | 0.860 | 30              |

### D. Process

It was carried out through the construction of a questionnaire, to evaluate the use of the project-based learning methodology in the skills of the network and communications course I, consisting of 20 questions which was validated by 5 experts. To do this, he coordinated with the management of the systems engineering school through which he sent a request to the director to authorize the examination of 20 questions. After the authorization of the application was validated, the date for the examination evaluation was coordinated with the students. The next stage was the development of the project-based learning methodology for 10 sessions, culminating the same, we proceeded at the end of the cycle with the examination

agreed upon both for the experimental group and for the control group. The guidelines for the development of the exam were given and the exam was delivered. Finally, with the evaluation carried out, it was extracted and digitized in Excel and then analyzed in SPSS.

### E. Data Analysis Method

For the study, the descriptive analysis was carried out using various contingency tables, as well as bar graphs and for the inferential analysis, the non-parametric Mann-Whitney U test was applied to test hypotheses, analyzing these data in the EXCEL 2016 programs and SPSS 25. Statistical analyzes allow to measure the relationships between variables with a higher [31].

### F. Ethical Aspects

The ethical aspects considered for carrying out this research were the authorization of the EAP Academic Professional School of Systems Engineering to apply the project-based learning method in the sixth cycle of the Private University in the morning and night shifts, as well as the application of the instruments were considered anonymous.

## III. RESULTS

The results obtained from the research carried out are shown:

### A. Description of the Variable Network and Communication Competences I

Table VI and Fig. 1 through the frequency bar diagram show that in the pretest and posttest of the control group they present similar results to 95% without change; while, in the pretest of the experimental group, 100% was found in the bad level and in the posttest of the experimental group: 21.1% was found in the bad level, 26.3% was found in the regular level and 52.6% was found good level, which shows a significant difference.

### B. Description of the Dimension Knowledge to Solve Local Area Network Design Problems

Table VII and Fig. 2 through the frequency bar diagram show that in the pretest and posttest of the control group they present similar results, while in the experimental group, in the pretest 47.4% were found to be at a bad level, 52.6 % was found at a regular level and in the post-test of the experimental group the improvements are observed, of which 10.5% were found at a bad level, 42.1% were found at a regular level and 47.4% were found at a good level, which shows a significant difference.

### C. Description of the Implementation Dimension of the Network Convergence

Table VIII and Fig. 3 through the frequency bar diagram show that in the pretest and posttest corresponding to dimension 2, of which for the control group in the pretest it is observed that 80% of the students are in a bad level, 10% are at a regular level and 10% are at the good level, then in the posttest it increases slightly to 100% of the students are at the bad level, while in the experimental group, the 94.7% was found in a bad level, 5.3% was found in a regular level and in the posttest of the experimental group the improvements are

observed, of which 31.6% were found in a bad level, 15.8% were found in a regular level and 52.6% A good level was found, which has a clear significant difference.

TABLE VI. LEVELS IN THE NETWORK AND COMMUNICATIONS COMPETENCES I OF THE PRETEST AND POSTTEST

|                    |      | Networking and communications competences I |         |       |       | Total |
|--------------------|------|---|---------|-------|-------|-------|
|                    |      | Bad   | Regular | Good  |       |       |
| Control group      | Pre  | fi  | 19      | 1     | 0     | 20    |
|                    |      | % fi  | 95.0%   | 5.0%  | 0.0%  | 100.0 |
|                    | Post | fi  | 19      | 1     | 0     | 20    |
|                    |      | % fi  | 95.0%   | 5.0%  | 0.0%  | 100.0 |
| Experimental group | Pre  | fi  | 19      | 0     | 0     | 19    |
|                    |      | % fi  | 100.0%  | 0.0%  | 0.0%  | 100.0 |
|                    | Post | fi  | 4       | 5     | 10    | 19    |
|                    |      | % fi  | 21.1%   | 26.3% | 52.6% | 100.0 |

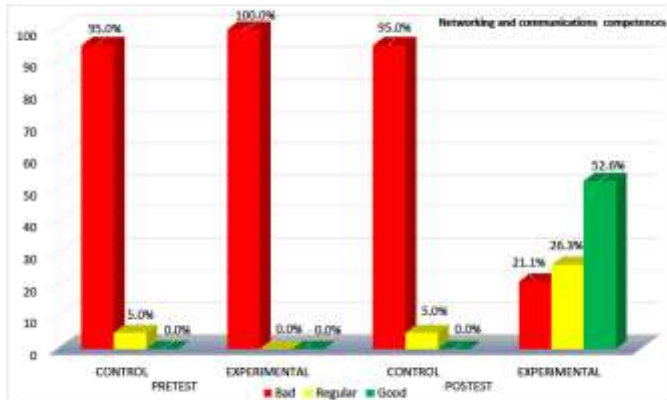


Fig. 1. Levels in Network and Communication Competences I of the Pretest and Posttest.

TABLE VII. LEVELS IN KNOWLEDGE TO SOLVE NETWORK DESIGN PROBLEMS OF LOCAL AREA OF PRETEST AND POSTTEST

|                    |      | Knowledge to solve local area network design problems |         |       |       | Total |
|--------------------|------|---|---------|-------|-------|-------|
|                    |      | Bad   | Regular | Good  |       |       |
| Control group      | Pre  | fi  | 13      | 6     | 1     | 20    |
|                    |      | % fi  | 65.0%   | 30.0% | 5.0%  | 100.0 |
|                    | Post | fi  | 14      | 6     | 0     | 20    |
|                    |      | % fi  | 70.0%   | 30.0% | 0.0%  | 100.0 |
| Experimental group | Pre  | fi  | 9       | 10    | 0     | 19    |
|                    |      | % fi  | 47.4%   | 52.6% | 0.0%  | 100.0 |
|                    | Post | fi  | 2       | 8     | 9     | 19    |
|                    |      | % fi  | 10.5%   | 42.1% | 47.4% | 100.0 |

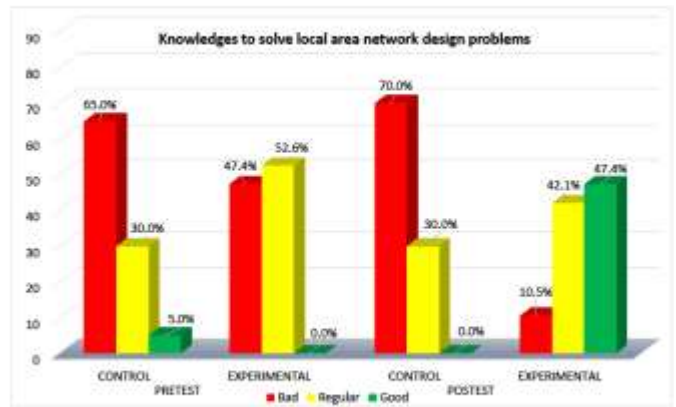


Fig. 2. Levels of Knowledge to Solve Pretest and Posttest Local Area Network Design Problems.

TABLE VIII. LEVELS IN THE IMPLEMENTATION OF THE CONVERGENCE OF THE PRETEST AND POSTTEST NETWORK

|                    |      | Implementation of network convergence |         |       |       | Total |
|--------------------|------|---------------------------------------|---------|-------|-------|-------|
|                    |      | Bad                                   | Regular | Good  |       |       |
| Control group      | Pre  | fi                                    | 16      | 2     | 2     | 20    |
|                    |      | % fi                                  | 80.0%   | 10.0% | 10.0% | 100.0 |
|                    | Post | fi                                    | 20      | 0     | 0     | 20    |
|                    |      | % fi                                  | 100.0%  | 0.0%  | 0.0%  | 100.0 |
| Experimental group | Pre  | fi                                    | 18      | 1     | 0     | 19    |
|                    |      | % fi                                  | 94.7%   | 5.3%  | 0.0%  | 100.0 |
|                    | Post | fi                                    | 6       | 3     | 10    | 19    |
|                    |      | % fi                                  | 31.6%   | 15.8% | 52.6% | 100.0 |

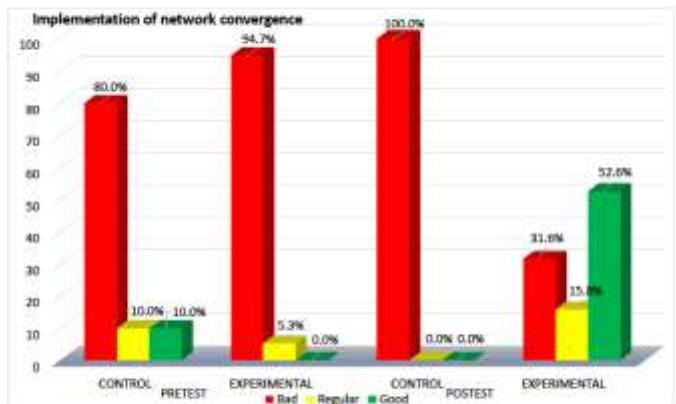


Fig. 3. Levels in the Implementation of the Convergence of the Pretest and Posttest Network.

D. Description of the Implementation Dimension of Network Administration

Table IX and Fig. 4 through the frequency bar diagram shows the pretest and posttest of dimension 3, in the control group in the pretest it is observed that 50% of the students are at a bad level, 45% at the regular level and 5% at the good level, then in the posttest it decreases slightly to 45% of students are at the bad level, 50% at the regular level and 5% at the good level, while the experimental group, in the pretest, found the 63.2% in bad level, 36.8% in regular level and in the

posttest of the experimental group 10.5% of students in bad level, 36.8% in regular level and 47.4% in good level, have significant difference.

TABLE IX. LEVELS IN THE ADMINISTRATION OF PRETEST AND POSTTEST NETWORKS

| Network administration |      |      |       |         |       | Total |
|------------------------|------|------|-------|---------|-------|-------|
| Control group          | Pre  | fi   | Bad   | Regular | Good  |       |
|                        |      |      |       |         |       |       |
| Control group          | Pre  | fi   | 10    | 9       | 1     | 20    |
|                        |      | % fi | 50.0% | 45.0%   | 5.0%  | 100.0 |
|                        | Post | fi   | 9     | 10      | 1     | 20    |
|                        |      | % fi | 45.0% | 50.0%   | 5.0%  | 100.0 |
| Experimental group     | Pre  | fi   | 12    | 7       | 0     | 19    |
|                        |      | % fi | 63.2% | 36.8%   | 0.0%  | 100.0 |
|                        | Post | fi   | 3     | 7       | 9     | 19    |
|                        |      | % fi | 15.8% | 36.8%   | 47.4% | 100.0 |

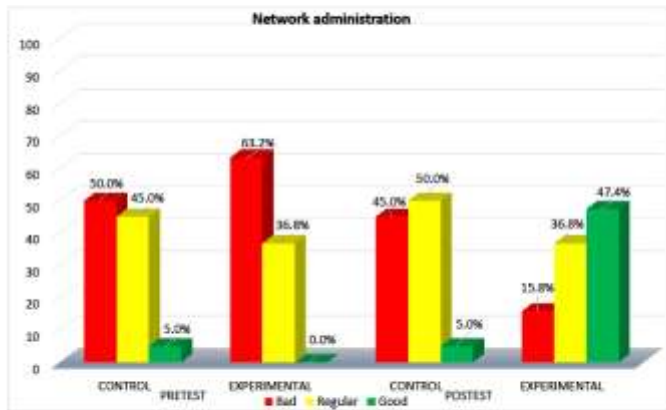


Fig. 4. Levels in Pretest and Posttest Network Administration.

E. Normality Test

The Kolmogorov-Smirnov method was used to check if the results follow a normal distribution. To do this, the Kolmogorov-Smirnov method is used to choose the corresponding test, due to the sample size being equal to or greater than 30.

Accordingly, in Table X, the corresponding significance values are shown to be less than  $\alpha = 0.05$ ; therefore,  $H_0$  is rejected, that is, the data do not follow a normal distribution. However, the 3rd is found. Posttest dimension with data record greater than  $\alpha = 0.05$ , but corresponds to the 6<sup>th</sup> part of the total of the three dimensions analyzed, as larger records with greater significance. Consequently, non-parametric tests will be applied for all inferential studies. Therefore, the Mann-Whitney U test will be used.

F. General Hypothesis Testing

$H_0$ : Project-based learning does not positively and significantly impact on network and communications competences I, in Engineering students, Lima 2020.

$H_a$ : Project-based learning has a positive and significant impact on network and communication competences I, in Engineering students, Lima 2020.

TABLE X. RESULTS OF THE KINDNESS OF ADJUSTMENT TEST FOR THE VARIABLE NETWORK AND COMMUNICATIONS COMPETENCES I

|                            | Kolmogorov-Smirnov |     |       |
|----------------------------|--------------------|-----|-------|
|                            | Statistical        | gl. | Sig.  |
| Pretest_D1_Knowledge       | , 182              | 39  | , 002 |
| Pretest_D2_Implementation  | , 254              | 39  | , 000 |
| Pretest_D3_Administracion  | , 168              | 39  | , 007 |
| Posttest_D1_Knowledge      | , 152              | 39  | , 024 |
| Posttest_D2_Implementation | , 282              | 39  | , 000 |
| Posttest_D3_Administracion | , 140              | 39  | , 051 |

TABLE XI. U MANN-WHITNEY TEST TO TEST THE GENERAL HYPOTHESIS

| Statistical   | Group            |                       | U test       |
|---------------|------------------|-----------------------|--------------|
|               | Control (n = 20) | Experimental (n = 19) | Mann-Whitney |
| Pretest       |                  |                       |              |
| U = 172,500   |                  |                       |              |
| Average range | 20.88            | 19.08                 | Z = -, 498   |
| Sum of ranges | 417.50           | 362.50                | p =, 618     |
| Posttest      |                  |                       |              |
| U = 31,000    |                  |                       |              |
| Average range | 12.05            | 28.37                 | Z = -4,488   |
| Sum of ranges | 241.00           | 539.00                | p = .000     |

According to the result described in Table XI, the network and communications competencies variable I, the control and experimental groups in the posttest show the U-Mann-Whitney = 31,000 and Z = -4,488, evidence of p less than  $\alpha = 0, 05$ ; therefore  $H_0$  is rejected and  $H_a$  is accepted, thus concluding that the variable network and communications competences I of the experimental group shows a significant improvement with respect to the control group, affirming that project-based learning positively and significantly impacts on networking and communications competences I.

G. Specific Hypothesis Test 1

$H_0$ : Project-based learning does not positively and significantly impact knowledge to solve local area network design problems in Engineering students, Lima 2020.

$H_a$ : Project-based learning positively and significantly impacts knowledge to solve local area network design problems in Engineering students, Lima 2020.

TABLE XII. MANN-WHITNEY TEST U TO TEST SPECIFIC HYPOTHESIS 1

| Statistical   | Group            |                       | U test       |
|---------------|------------------|-----------------------|--------------|
|               | Control (n = 20) | Experimental (n = 19) | Mann-Whitney |
| Pretest       |                  |                       |              |
| U = 183,500   |                  |                       |              |
| Average range | 19.68            | 20.34                 | Z = -, 187   |
| Sum of ranges | 393.50           | 386.50                | p =, 852     |
| Posttest      |                  |                       |              |
| U = 23,000    |                  |                       |              |
| Average range | 11.65            | 28.79                 | Z = -4,748   |
| Sum of ranges | 233.00           | 547.00                | p = .000     |



According to the result described in Table XII, the knowledge dimension to solve local area network design problems, the control and experimental groups in the posttest show the U-Mann-Whitney = 23,000 and Z = -4,748, it is evident p lower at  $\alpha = 0.05$ ; therefore Ho is rejected and Ha is accepted.

H. Specific Hypothesis Test 2

Ho: Project-based learning does not positively and significantly impact the implementation of network convergence in Engineering students, Lima 2020.

Ha: Project-based learning has a positive and significant impact on the implementation of network convergence in engineering students, Lima 2020.

According to the result described in Table XIII, the implementation dimension of network convergence, the control and experimental groups in the posttest show the U-Mann-Whitney = 56,000 and Z = -3,975, with evidence of p less than  $\alpha = 0, 05$ ; therefore Ho is rejected and Ha is accepted.

I. Specific Hypothesis Test 3

Ho: Project-based learning does not have a positive and significant impact on the administration of networks in Engineering students, Lima 2020.

Ha: Project-based learning positively and significantly impacts network administration in engineering students, Lima 2020.

TABLE XIII. MANN-WHITNEY TEST U TO TEST SPECIFIC HYPOTHESIS 2

| Statistical   | Group            |                       | U test       |
|---------------|------------------|-----------------------|--------------|
|               | Control (n = 20) | Experimental (n = 19) | Mann-Whitney |
| Pretest       |                  |                       |              |
| U = 185,000   |                  |                       |              |
| Average range | 19.75            | 20.26                 | Z = -, 154   |
| Sum of ranges | 395.00           | 385.00                | p =, 878     |
| Posttest      |                  |                       |              |
| U = 56,000    |                  |                       |              |
| Average range | 13.30            | 27.05                 | Z = -3,975   |
| Sum of ranges | 266.00           | 514.00                | p = .000     |

TABLE XIV. MANN-WHITNEY TEST U TO TEST SPECIFIC HYPOTHESIS 3

| Statistical   | Group            |                       | U test       |
|---------------|------------------|-----------------------|--------------|
|               | Control (n = 20) | Experimental (n = 19) | Mann-Whitney |
| Pretest       |                  |                       |              |
| U = 169,500   |                  |                       |              |
| Average range | 21.03            | 18.92                 | Z = -, 588   |
| Sum of ranges | 420.50           | 359.50                | p =, 557     |
| Posttest      |                  |                       |              |
| U = 65,500    |                  |                       |              |
| Average range | 13.78            | 26.55                 | Z = -3,533   |
| Sum of ranges | 275.50           | 504.50                | p = .000     |

According to the result described in Table XIV, the network administration dimension, the control and experimental groups in the posttest show the U-Mann-Whitney = 65,000 and Z = -3,533, with p-value less than  $\alpha = 0.05$ ; therefore Ho is rejected and Ha is accepted.

IV. DISCUSSION

Regarding the validation of the general hypothesis, the statistical results obtained values of Z = -4.488 and p value of 0.000 (see Table XI), so that project-based learning positively and significantly impacts on network and communication competences I, in engineering students, Lima 2020, for which the results obtained in the post-test, it was determined that the number of students who are at the regular and good level both are approximately 80% of the total population of the experimental group (see Table VI), verifying that the use of the programmed designed project-based learning to improve the learning competences of the communication networks I course, improved the understanding of it in all its dimensions. thus confirming the importance of using the project-based learning methodology [27].

Regarding the validation of the specific hypotheses, the analysis by dimensions was carried out, in which dimension 1 consists of 8 questions, dimension 2 consists of 3 questions and dimension 3 consists of 9 questions, totaling 20 questions. For the confirmation of the specific hypothesis 1; obtained values of Z = -4,748 and p value of 0.000 in the posttest of the experimental group compared to a Z = -, 187 and p value of 0.852 of the control group, thus also at the percentage level it is observed that in the pretest in the bad level they were 47.4% of students and in the post-test the percentage dropped to 10.5% of students, thus it is also observed that in the case of the regular level in the pretest, 52.6% of students are observed and in the posttest the percentage dropped to 42.

To validate the specific hypothesis 2; values of Z = -, 154 were obtained greater than -1.96 (critical point) and a p value = 0.878 in the pretest there are initially no differences and then when obtaining a Z = -3.975 and a p value = 0.000 in the posttest , finding significant differences between the control group and the experimental group, as well as at the percentage level, it is observed that in the pretest the bad level was 94.7% of students and in the posttest the percentage fell to 31.6% of students, it is also observed that in the case of the regular level in the pretest, 5.3% of students are observed and in the posttest the percentage increased to 15.8% of students and finally in the experimental group, a significant growth percentage at the good level, being initially 0.0% to 52.6% in the posttest in the pretest.

To confirm the specific hypothesis 3; values of Z = -3,533 and p value of 0.000 were obtained in the posttest of the experimental group compared to a Z = -, 588 and p value of 0.557 of the control group, as well as percentage results were obtained in order to observe the changes between the pretest and posttest of the control group as well as the experimental group, in which it is observed that in the pretest in the bad level were 63.2% of students and in the posttest the percentage fell to 15.8% of students, Thus, it is also observed that in the case of the regular level in the pretest, 36.8% of students were observed and in the posttest, the percentage was maintained at

36.8% of students, and finally, in the experimental group, a percentage of significant growth in the good level, being initially 0 in the pretest, 0% to 47.4% in the posttest. Thus, according to the statistical results found in specific hypothesis 3, project-based learning positively and significantly impacts network administration in engineering students, Lima (see Fig. 4).

## V. CONCLUSIONS

We showed that project-based learning caused a positive and significant impact on the dependent variable network and communication competences I, in engineering students, Lima 2020, because the significance level was obtained Sig. = 0.000 less than  $\alpha = 0,05$  ( $p < \alpha$ ) and  $Z = -4,488$ , indicating that the proposed model is appropriate.

It was verified that project-based learning caused a positive and significant impact on the knowledge dimension to solve local area network design problems in engineering students, Lima 2020, due to the fact that it reached the significance level Sig. = 0.000 less than  $\alpha = 0.05$  ( $p < \alpha$ ) and  $Z = -4,748$ , demonstrating that the proposed model is acceptable.

Project-based learning was found to have a positive and significant impact on the implementation dimension of network convergence in engineering students, Lima 2020, because the level of significance Sig. = 0.000 is less than  $\alpha = 0,05$  ( $p < \alpha$ ) and  $Z = -3,975$ , indicating that the proposed model is appropriate.

Project-based learning was shown to have a positive and significant impact on the dimension of network administration in Engineering students, Lima 2020, because the level of significance Sig. = 0.000 is less than  $\alpha = 0.05$  ( $p < \alpha$ ) and  $Z = -3,533$  of the experimental group in the posttest calculation, confirming that the proposed model is valid.

## VI. RECOMMENDATIONS

It is recommended to both public and private universities to implement the use of the project-based learning program to improve competences in the course of networks and communications, in order to use it as a model for their best student performance in classes, as well as motivate them in the investigations of all professional academic activities permanently in order to contribute to individual and team-student-teacher-university improvement so that they can interpret, design, plan and implement various solutions to the same problem, both personally, family, social and professionally.

We also suggest to train the teachers of the networks and communications course to learn the project-based learning methodology, to use it as a reference model for class sessions, since activities structured by a sequence of phases are developed to carry out prototype professional projects, managing to generate better research among teachers, thus improving the levels of competences for the knowledge dimension to solve local area network design problems.

It is recommended to develop the implementation dimension of network convergence, through the project-based learning program, in order to improve your competences to

successfully face the new challenges of emerging technological trends. Finally, it is recommended to learn the network administration dimension through project-based learning, in order to improve technological organizations, closing the gap in the job sector.

## REFERENCES

- [1] C. N. Casimiro Urcos, W. H. Casimiro Urcos y J. F. Casimiro Urcos, "Desarrollo de competencias profesionales en estudiantes universitarios". Conrado, vol. 15, n.º 70, pp. 312-319. Epub 02-Dic-2019. ISSN 2519-7320. [http://scielo.sld.cu/scielo.php?script=sci\\_arttext&pid=S1990-86442019000500312&lng=es&nrm=iso](http://scielo.sld.cu/scielo.php?script=sci_arttext&pid=S1990-86442019000500312&lng=es&nrm=iso).
- [2] ITU. Conjunto de Herramientas para las Habilidades Digitales. Suiza-Ginebra, 2018. [https://www.itu.int/en/ITU-D/Digital-Inclusion/Documents/Digital-Competences-Toolkit\\_Spanish.pdf](https://www.itu.int/en/ITU-D/Digital-Inclusion/Documents/Digital-Competences-Toolkit_Spanish.pdf).
- [3] Unesco. Qué hace la UNESCO en materia de educación superior, May 2018. <https://es.unesco.org/themes/educacion-superior/accion>.
- [4] S. Núñez López, J.E. Ávila Palet y S. L. Olivares Olivares, "The development of critical thinking through problem-based learning", RIES, México, UNAM-IISUE/Universia, vol. 8, n.º 23, pp. 84-103, Agosto 2016. <https://www.redalyc.org/jatsRepo/2991/299152904005/html/index.html>
- [5] MTPE. Subsector Telecomunicaciones: Ocupaciones frecuentes y oferta formativa a nivel nacional, 2019. <http://cdn.www.gob.pe/uploads/document/file/458911/telecomunicaciones.pdf>.
- [6] C. Sotomayor-Beltran and L. Andrade-Arenas, "A spatial assessment on internet access in Peru between 2007 and 2016 and its implications in education and innovation," 2019 IEEE 1st Sustainable Cities Latin America Conference (SCLA), Arequipa, Peru, 2019, pp. 1-4.
- [7] SUNEDU, "Informe bienal sobre la realidad universitaria peruana" Sunedu, Lima-Perú, 2018. <https://www.sunedu.gob.pe/informe-bienal-sobre-realidad-universitaria/>.
- [8] V. F. Martins, P. N.M. Sampaio, A. J. A. Cordeiro and B. Ferreira Viana, "Implementing a Data Network Infrastructure Course using a Problem-based Learning Methodology". Journal of Information Systems Engineering & Management, vol. 3, n.º 2, pp. 1-7, April 2018. DOI: <https://doi.org/10.20897/jisem.201810>.
- [9] M. Beier, M. Kim, A. Saterbak, V. Leautaud, S. Bishnoi and J.M. Gilberto, "The effect of authentic project - based learning on attitudes and career aspirations in STEM". Journal of Research in Science Teaching, Mayo 2018. DOI: <https://doi.org/10.1002/tea.21465>.
- [10] M. Fernández Redondo, C. Hernández Espinosa, and J. Sales Gil, "Learning of Computer Networks through the use of Projects in a Video Game Degree", Zaragoza-España, pp. 593-598, 2017. DOI: 10.26754/CINAIC.2017.000001\_124.
- [11] L. Qin, "Design and Realization of Project-based Computer English Learning System". International Journal of Emerging Technologies in Learning, vol. 12, n.º 08, pp. 128-136, June 2017. DOI: <https://doi.org/10.3991/ijet.v12.i08.7147>.
- [12] C. Baluarte Araya, "Project based Learning Application Experience in Engineering Courses: Database Case in the Professional Career of Systems Engineering". International journal of computer science and advanced applications (IJACSA), vol. 11, n.º 3, pp. 128-136, 2020. <http://dx.doi.org/10.14569/IJACSA.2020.0110316>.
- [13] L. Torres, "Plataforma virtual para mejorar el rendimiento en una asignatura del plan curricular de la escuela de tecnologías de la información, Senati", tesis doctoral, escuela de posgrado de la Universidad Peruana Cayetano Heredia, Lima, Perú, 2019. [http://repositorio.upch.edu.pe/bitstream/handle/upch/7726/Plataforma\\_TorresArgomedeo\\_Leonardo.pdf?sequence=1&isAllowed=y](http://repositorio.upch.edu.pe/bitstream/handle/upch/7726/Plataforma_TorresArgomedeo_Leonardo.pdf?sequence=1&isAllowed=y).
- [14] F. Rodríguez, "Aprendizaje basado en proyectos en el nivel de competencias investigativas en estudiantes de Instituto Pedagógico, Trujillo, 2017", tesis doctoral, escuela de posgrado de la Universidad Cesar Vallejo, Lima, Perú, 2018. [http://repositorio.ucv.edu.pe/bitstream/handle/20.500.12692/22688/rodriguez\\_vf.pdf?sequence=1&isAllowed=y](http://repositorio.ucv.edu.pe/bitstream/handle/20.500.12692/22688/rodriguez_vf.pdf?sequence=1&isAllowed=y).
- [15] D. G. Hostia Luque, "Aprendizaje basado en proyectos colaborativos y competencias de los estudiantes de tercer año de Ingeniería de Sistemas de la Universidad Nacional San Luis Gonzaga de Ica", tesis de maestría,

- escuela de posgrado de la Universidad Nacional San Luis Gonzaga de Ica, Lima, Perú, 2018. <http://repositorio.une.edu.pe/bitstream/handle/UNE/2467/TM%20CE-Du%204071%20H1%20-%20Hostia%20Luque.pdf?sequence=1&isAllowed=y>.
- [16] L. M. Zegarra Ramírez, “Efectos de la aplicación de la metodología de aprendizaje basado en proyectos en el desarrollo de competencias en el curso de procesos de manufactura II”. tesis de maestría, escuela de posgrado de la Universidad Peruana Cayetano Heredia, Lima, Perú, 2019. [http://repositorio.upch.edu.pe/bitstream/handle/upch/1443/Efectos\\_ZegarraRamirez\\_Leonor.pdf?sequence=1&isAllowed=y](http://repositorio.upch.edu.pe/bitstream/handle/upch/1443/Efectos_ZegarraRamirez_Leonor.pdf?sequence=1&isAllowed=y).
- [17] A. Rofieq, R. Latifa, E. Susetyarini and P. Purwatiningsih, “Project-based learning: Improving students’ activity and comprehension through lesson study in senior high school”. JPBI (Jurnal Pendidikan Biologi Indonesia), vol. 5, n.º 1, pp. 41-50, March 2019. DOI: <https://doi.org/10.22219/jpbi.v5i1.7456>.
- [18] J. Gutiérrez, G. De la Puente, A. Martínez, and E. Piña. Aprendizaje basado en problemas. Editorial Universidad Nacional Autónoma de México, 2013. [https://portalacademico.cch.unam.mx/materiales/libros/pdfs/librocch\\_abp.pdf](https://portalacademico.cch.unam.mx/materiales/libros/pdfs/librocch_abp.pdf).
- [19] K. Sormunen, K. Juuti and J. Lavonen, “Maker-Centered Project-Based Learning in Inclusive Classes: Supporting Students’ Active Participation with Teacher-Directed Reflective Discussions”. *Int J of Sci and Math Educ* 18, pp. 691–712, April 2020. DOI: <https://doi.org/10.1007/s10763-019-09998-9>.
- [20] I. Saputra, S. Joyoatmojo, and H. Harini, “The implementation of project-based learning model and audio media Visual can increase students’ activities”. *International Journal of Multicultural and Multireligious Understanding*, vol. 5, n.º 4, pp. 166-174, 2018. DOI: <http://dx.doi.org/10.18415/ijmmu.v5i4.224>.
- [21] M. De la Puente, D. Guerra, C. de Oro and C. McGarry, “Undergraduate students’ perceptions of Project-Based Learning (PBL) effectiveness: A case report in the Colombian Caribbean. *Cogent Education*”, *Cogent Education*, vol. 6, n.º 1, 1616364, pp. 1- 17, May 2019. DOI: <https://doi.org/10.1080/2331186X.2019.1616364>.
- [22] G. Geitz and J. de Geus, “Design-based education, sustainable teaching, and learning”, *Cogent Education*, vol. 6, n.º 1, 1647919, pp. 1- 15, July 2019. DOI: <https://doi.org/10.1080/2331186X.2019.1647919>.
- [23] F. D. Mendoza Vargas. “Relación Entre La Actitud Experiencial y La Utilización De Simuladores Como Herramienta Pedagógica”. tesis de maestría, facultad de ciencias administrativas y contables de la Universidad de La Salle, Bogotá, Colombia, Noviembre 2015. [https://ciencia.lasalle.edu.co/cgi/viewcontent.cgi?article=1481&context=maest\\_administracion](https://ciencia.lasalle.edu.co/cgi/viewcontent.cgi?article=1481&context=maest_administracion).
- [24] UCH. Silabo de Redes y Comunicaciones I. Escuela Profesional de Ingeniería de Sistemas 2019, diciembre 2019, Lima, Perú.
- [25] C. M. Amador Ortiz y L. Velarde Peña, “ICT Competences in students of higher education, a case study”. *RIDE*, vol. 10, n.º 19, ago. 2019. DOI: <https://doi.org/10.23913/ride.v10i19.515>.
- [26] C. De Anda, N. Galaviz y R. Santiago. Tecnología de la Información 1. Editorial Gyros S.A. Ciudad Universitaria Culiacán, Sinaloa, México, Agosto 2019. [https://issuu.com/profejrul/docs/libro\\_tiy\\_c\\_1\\_uas](https://issuu.com/profejrul/docs/libro_tiy_c_1_uas).
- [27] J. F. Kurose and K. W. Ross. Redes de computadoras. un enfoque descendente. Editorial Pearson Educación, S. A. Madrid, España. 7ma, 2017. Edición. [https://www.academia.edu/40738627/Redes\\_de\\_computadoras\\_Un\\_enfoque\\_descendente\\_7a\\_Edici%C3%B3n](https://www.academia.edu/40738627/Redes_de_computadoras_Un_enfoque_descendente_7a_Edici%C3%B3n).
- [28] A. López y P. Díaz. “Capítulo 3: Investigación, emprendimiento y TIC, elementos de una propuesta pedagógica en la Universidad Popular del Cesar”, en Congreso Internacional sobre el Enfoque Basado en Competencias, por editorial Corporacion CIMTED, Colombia, vol. 10, n.º 1, pp. 146-162, Marzo 2018. <http://memoriascimted.com/wp-content/uploads/2016/02/Memorias-CIEBC2018.pdf>.
- [29] C. de Pelekais, O.El kadi, C. Seijo y N. Neuman. El ABC de la Investigación. Pauta Pedagógica. Septima Edicion, Maracaibo, Venezuela. Editorial: Astro Data S.A., 2015.
- [30] H. Ñaupas, E. Mejía, E. Novoa y A. Villagómez. Metodología de la investigación: cuantitativa - cualitativa y redacción de la tesis, 4ta. Edicion. Bogotá, Colombia: Ediciones de la U., Abril 2014.
- [31] W. Rengel y M. Giler. Publicar investigación científica Metodología y desarrollo, 1ra. Edicion, Manabi-Ecuador, Editorial: Mar Abierto, 2018.

# Artificial Intelligent Techniques for Palm Date Varieties Classification

Lazhar Khriji<sup>1\*</sup>, Ahmed Chiheb Ammari<sup>2</sup>, Medhat Awadalla<sup>3</sup>

Department of Electrical and Computer Engineering  
College of Engineering, Sultan Qaboos University  
Muscat, Oman

**Abstract**—The demand on high quality palm dates is increasing due to its energy value and nutrient content, which are of great importance in human diet. To meet consumer and market standards with large-scale production, in Oman as among the top date producer, an inline classification system is of great importance. This paper addresses the potentiality of using Machine-Learning (ML) techniques in classifying automatically, without any physical measurement, the six most popular date fruit varieties in Oman. The effect of color, shape, size, and texture features and the critical parameters of the classifiers on the classification efficiency has been endeavored. Three different ML techniques have been used for automatic classification and qualitative comparison: (i) Artificial Neural Networks (ANN), (ii) Support Vector Machine (SVM), and (iii) K-Nearest Neighbor (KNN). Based on the merge of color, shape and size features contributes to achieve the highest accuracy. Experimental results show that the ANN classifier outperforms both SVM and KNN with the highest classification accuracy of 99.2%. This developed vision system in this paper can be successfully integrated in the packaging date factories.

**Keywords**—Palm date; feature extraction; machine learning; computer vision

## I. INTRODUCTION

The great energy and the content of nutrients present in the date fruits are the reasons behind the great importance of them in human diet. Referring to Food Agricultural Organization (FAO), the largest date producer countries in the world are located in the Middle East and North Africa. Oman produces about 260-270 Million-tons and ranks among the largest dates producing countries in the world [1-2]. However, only around 7000 tons of date fruits is reported under export [3]. This low output could be related to the required higher international export standards such as color, size, softness, freshness, etc. In order to diversify the sources of income, Oman has considered palm dates as a priority to its economy. It encouraged to plant palm dates that reached more than 250 varieties. Texture, size, shape and color are the main used features to differentiate between varieties [4-6]. Khalas, Fardh, Khunaizi, Qash, Naghal, and Maan are known as the six most popular date fruit varieties in Oman. The sweetest variety is Khunaizi and the delicious variety is Khalas [2]. In date's industries, the classification of dates into diverse classes is an essential task. Using intelligent computerized systems this classification task can be automated to produce an accurate and fast classification of date's varieties compared to traditional way. Therefore, the

related industries are improving their products in quality and quantity [7, 8].

This paper aims to propose a computerized vision system that automatically classifies date's varieties based on image processing techniques combined with Artificial Intelligence algorithms. Traditionally, palm date's classification is performed based on grade [9]. Starting from 2012, computer vision and pattern recognition have been introduced for automatic date's fruits classification. The authors have tested seven categories of dates and fifteen features have been extracted. To compare the results, they used multiple classifiers such as Neural Networks (NN), Linear Discriminant Analysis (LDA) and Nearest Neighbor [10]. A sorting system based on ANN was presented in the context of date fruits in 2012. Two neural networks models have been used. The first model is using a multi-layer perceptron (MLP), the second model is using Radial Basis Function (RBF) networks with a backpropagation learning algorithm. The performance accuracy of 87.5% and 91.1% have been achieved using MLP and RBF, respectively [11]. An automatic system for classification, which uses different images of dates, is used to classify different types of dates [12, 13]. In this work, different features are extracted from the images of the dates such as the shape, texture, and the color. Fisher discrimination Ratio (FDR) has been used to reduce the dimension of features vector where SVM was used as a classifier [12, 13]. For date's classification relying on hardness, a system equipped with a monochrome camera was presented in 2016. This study used histogram and texture features in their system and LDA and ANN were implemented as classifiers [14-16]. In 2018, an automated system that identifies different date fruit maturity status and classifies their categories is developed. Color, size and skin texture features are extracted. The system counts the number of dates, classifies them into different classes and identifies the defects [17].

Our aim is to classify automatically, without any physical measurement, the top six date palm varieties in Oman. We will work to extract color, shape, size, and texture features of various date images. Artificial Neural Network (ANN), Support Vector Machine (SVM), and K-Nearest Neighbor (KNN) classifiers are proposed and comparative performance analysis are conducted.

This paper is organized as follows: Section II describes the materials and methods. Experimental results and discussions are given in Section III. Section IV concludes the paper.

\*Corresponding Author

This project was funded partially by Sultan Qaboos University, Deanship of Scientific Research, under grant number "IG/ENG/ECED/19/01", and partially by OMANTEL under grant number "EG/SQU-OT/18/01".

## II. MATERIALS AND METHODS

The proposed and developed system flowchart is given in Fig. 1. First, a dataset of colored images of dates must be achieved, where each image contains only a single date. To prepare the images free of noisy segmented pixels, different operations such as segmentation and mathematical morphological are then applied. Then different features are extracted from these segments. In the training phase, to address the importance of each feature, different classifiers are trained. Then the classifiers are trained again with a combination of two or more features together. Based on the achieved results in each training process. The most effective features have been determined and used to update the learning parameters of the classifiers. In the testing phase, new and unseen data are presented to the classifiers for testing.

### A. Samples Collection

The most valuable and common dates in Oman such as Khalas, Fardh, Khunaizi, Qash, Naghal and Maan have been used in this study as shown in Fig. 2. AL-Dhahirah Governorate is the main source for the collected samples. The dataset comprises six types, 100 samples for each class.

### B. Image Acquisition System

RGB color camera, fluorescent lights and, EOS 1100D, Canon, Taiwan, resolution of  $4272 \times 2848$  pixel Personal computer have been used in the developed system [18, 19] are. For background, a white paper with A4 size is used. Each date is 15 cm away from the used camera and is manually placed. The mode of self-timer embedded in the camera helps to take images for each sample. To reduce expected noise, more images have been take.

### C. Preprocessing and Segmentation

MATLAB Toolbox (Version R2014a, The Mathworks Inc., Natick, MA, USA) helped to develop algorithms to accomplish different operations on the acquired images. The procedure for image processing is illustrated in Fig. 3. For simple and prompt processing, the image samples are resized. After that, grayscale images are obtained by the conversion of the colored images meanwhile another samples of colored images are maintained for more processes. The foreground and background regions of the images are separated from the grayscale images based on Otsu's method [20] and then followed by morphological operations.

### D. Extraction of Features

The identification of the effective features is considered as the most challenging process of classification of dates. The most vital features that can be used for date's classification are colors features, size-shape features, and skin texture features as depicted in Fig. 4.

1) *Color features*: Since the color feature is the most important feature dates varieties that provide the remarkable information for the classification of dates. At the starting, the channels of Red, Green, and Blue are separated from the cropped RGB images. Then, the mean and the standard

deviation were calculated from each channel. However, the minimum, maximum and mean intensities are determined from the gray images. The pixel that has the smallest intensity and the pixel that has the biggest intensity represent the minimum and maximum intensity, respectively. The mean intensity is represented by the mean values of all pixels [10].

2) *Shape and size features*: In addition to the color feature, the size and shape are important features for the classification of the dates. These features enhance the classification accuracy. Different features of the shape and size can be achieved from the segmented images such as Area, Major axis length, Minor axis length, Ellipse eccentricity, solidity and perimeter (see Fig. 5). Equation 1 is used to calculate the Eccentricity and Equation 2 computes the Solidity [7, 10, 21].

Eccentricity ( $e$ ) is computed as,

$$e = \frac{c}{a} = \frac{\sqrt{a^2 - b^2}}{a} \quad (1)$$

$$\text{Solidity} = \frac{\text{Area}}{\text{Convex Area}} \quad (2)$$

3) *Texture features*: Skin texture can be used to separate some sort of palm dates. Therefore, it is essential to consider dates texture as features. Statistical texture features can be calculated by the Gray Level Co-occurrence Matrix (GLCM) method [22]. GLCM indicates how often the particular gray level pixel pair  $(i, j)$  with a distance  $l$  and relative orientation  $\theta$  has applied in a given matrix and is represented by  $G_{\theta,l}(i, j)$ .

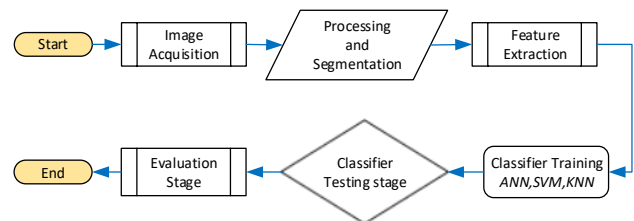


Fig. 1. System Flowchart.

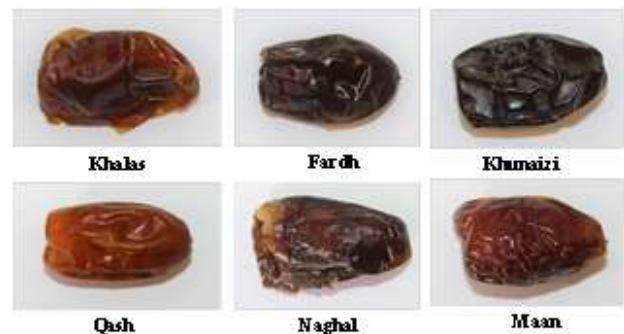


Fig. 2. Samples of the Date's Dataset used.

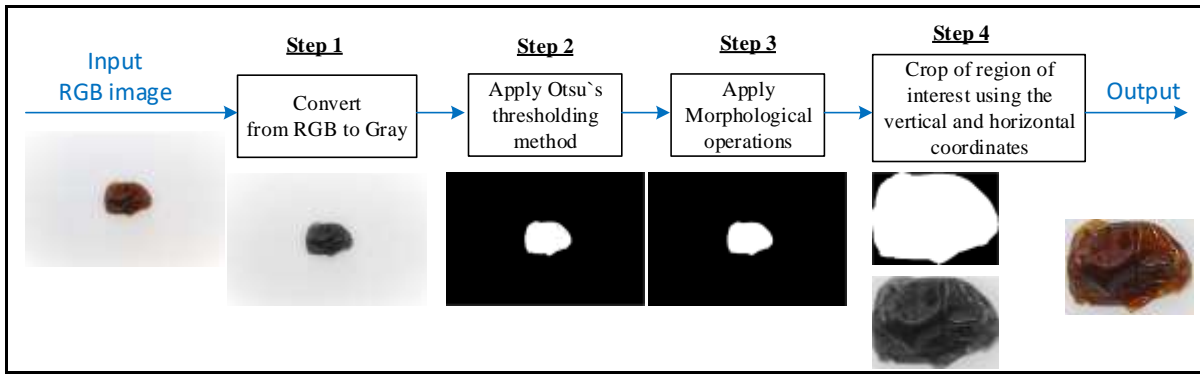


Fig. 3. Segmentation Steps Applied to the Date Samples.

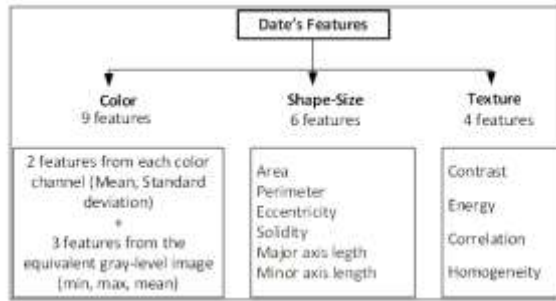


Fig. 4. Categorization of Date's Features.

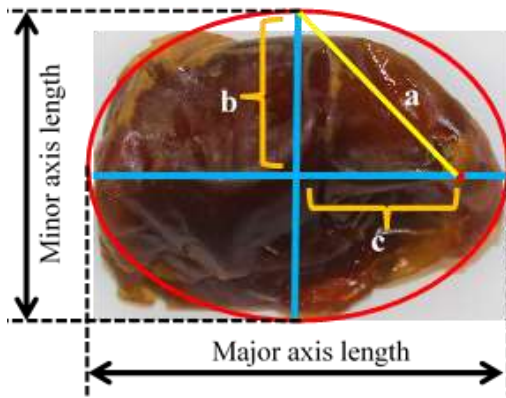


Fig. 5. Major Axis, Minor Axis and Eccentricity Parameters of the Ellipse.

The 4 directions are the orientations of  $\theta$  in  $0^\circ$ ,  $45^\circ$ ,  $90^\circ$  and  $135^\circ$ . GLCM is a function of distance  $d$  and at different angles. Only one single direction of  $\theta = 0^\circ$  is considered in our work leading to four features. The corresponding GLCM is given by;

$$G_{0^\circ}(i, j) = \left\{ \begin{array}{l} [(q, r), (m, n)] \in D \\ q - m = 0, |r - n| = l \\ f(q, r) = i, f(m, n) = j \end{array} \right\} \quad (3)$$

The value of the intensity in  $(m \times n)$  images is determined by the function  $f(q, r)$ . Assuming that  $P_{\vec{O}}(i, j)$  is the GLCM of an image  $I(q, r)$  within the region  $I_B$  for an offset vector  $\vec{O}$  showing co-occurrence count of intensity pair  $(i, j)$  the four features are defined as:

a) Contrast represents the sudden change in value of intensity in a given image contrast between a pixel and its neighbor over the whole image.

$$C = \sum_{i=1}^{N_g} \sum_{j=1}^{N_g} (i-j)^2 P_{\vec{O}}(i, j) \quad (4)$$

b) Energy (E) is the sum of squared elements in GCLM and is given by equation 5.

$$E = \sum_{i=1}^{N_g} \sum_{j=1}^{N_g} P_{\vec{O}}(i, j)^2 \quad (5)$$

c) Correlation is the measure of the similarity (S) and is given by equation 6.

$$S = \sum_{i=1}^{N_g} \sum_{j=1}^{N_g} \frac{P_{\vec{O}}(i, j)(i - \mu_i)(j - \mu_j)}{\sigma_i \sigma_j} \quad (6)$$

d) Homogeneity (H) is the closeness measure and is given by equation 7.

$$H = \sum_{i=1}^{N_g} \sum_{j=1}^{N_g} \frac{P_{\vec{O}}(i, j)}{1 + |i - j|} \quad (7)$$

Where,  $N_g$  is the number of distinct gray levels and  $p(i, j)$  represents the  $(i, j)^{th}$  entry in the GLCM. The means in the row and the column directions are  $\mu_i$  and  $\mu_j$ , respectively. The standard deviations in the row and the column directions are  $\sigma_i$  and  $\sigma_j$ , respectively.

Only one single direction of  $\theta = 0^\circ$  is considered in our work leading to four features.

### III. EXPERIMENTAL RESULTS AND DISCUSSION

#### A. Setup of the Artificial Neural Network

In this work, feed-forward (two layers) neural networks trained with Levenberg-Marquardt backpropagation algorithm. The function of tansig and logsig is used as an activation function in the hidden neurons and the function of softmax is used in the output neurons [22]. Four hundreds and eight images from the collected samples are used for the neural network training, and seventy-two images are used for the network validation while one hundred and twenty images are used for the network testing. Different neural network

architectures in terms of different numbers of neurons per each hidden layer are trained with different numbers of features to achieve the most remarkable architecture and the most valuable features. The training process has been repeated for thirty times and the results have been averaged for the sake of reliability.

Fig. 6 and 7 show the plots of the accuracies of the ANN classifier using different features vs. number of hidden neurons for logsig and tansig activation functions respectively. It is clear that as the number of hidden neurons increases from 1 to 3, the increase in the accuracy is very clear. It is noticed that, when the number of used neurons in the range between four and ten, the improvement in the classification accuracy is subtle.

In addition, the classification accuracy based on the four texture features was low compared with the other features. The performance obtained were 36.08% and 53.72% when the function of logsig is used. However, the performance obtained were 35.36% and 54.36% when the function of tansig is used. It concludes that the texture feature cannot contribute thoroughly in the classification process of the dates. While, the color and shape features can participate to a high extent in the classification of the dates. The achieved performance accuracy using the features of the color and shape was in the range of 58.03% to 80.06% and from 62.58% to 79.67% for using both of logsig function and of tansig function, respectively. As shown in Fig. 6 and Fig. 7, the contribution of both color and shape are comparable. The classification accuracy using shape accuracy is in the range of 65.88% to 81.11% for using logsig function and in the range of 63.67% to 81.19% for using tansig function. There was an improvement in the classification process when all features are used together. The accuracy was 96.22% for using of logsig function and 96.21% for using tansig function. However, there was an improvement in the accuracy (97% and 97.26%) when the texture features were excluded and only color and shape features were used. In our research work, a remarkable performance accuracy of 97.26% was achieved, using a hidden layer that includes seven-neurons and the tansig function as used as depicted in Table I. In case of using logsig function and nine neurons in the hidden layer, a higher accuracy was obtained.

**B. Setup of the Support Vector Machine**

In this approach, the training subset is portioned to 10 parts. Nine portions in each iteration is used for the training process, and one portion only is used in the validation process. The rest of the dataset (20%) are used in the testing process. The optimal value of the kernel scale of Radial Basis Function, RBF, was set automatically but the Super Vector Machine, SVM, optimization parameter “C” changed for five different values in the range [1e-1, 1e-3].

The achieved results that represent the relation between accuracy of different features and the box constraint parameter C is illustrated in Fig. 8. It is clear from the achieved results that the performance of SVM and ANN are similar. Again, the achieved performance in terms of the features of texture were the smallest, in the ranges of 47.06% to 57.14%.

However, the combination of color and shape-size features reaches the highest accuracy of 97.1386% when the box

constraint is 10. When the box constraint increases more than 10, the accuracy of different features either decreases or increases with very small amounts. When all features are used, the accuracy was reduced little bit as compared to color and shape-size combination which shows that texture feature can be ignored.

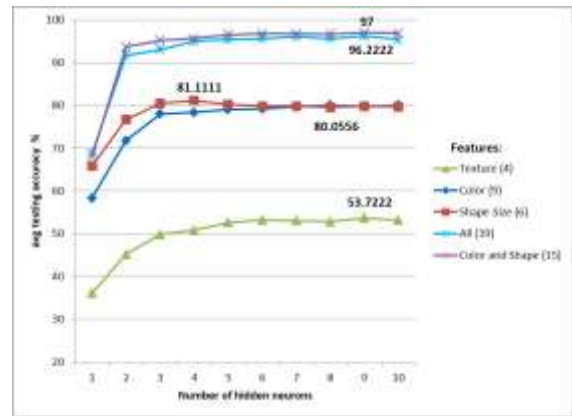


Fig. 6. Plots of the Accuracies of the ANN Classifier using different Features vs Number of Hidden Neurons (Logsig-Softmax).

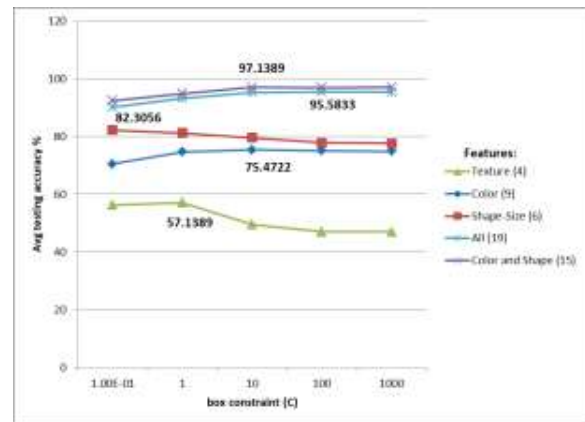


Fig. 7. Plots of the Accuracies of the ANN Classifier using different Features vs Number of Hidden Neurons (Tansig-Softmax).

TABLE I. THE HIGHEST ACHIEVED PERFORMANCES (%) OF THE CLASSIFIERS ARE GIVEN BY THEIR CONFUSION MATRIX USING COLOUR AND SHAPE-SIZE FEATURES

| Classifier Class | ANN-logsig   |             | ANN-tansig   |              | SVM          |              | KNN         |              |
|------------------|--------------|-------------|--------------|--------------|--------------|--------------|-------------|--------------|
|                  | Recall       | Precision   | Recall       | Precision    | Recall       | Precision    | Recall      | Precision    |
| Khalas           | 100          | 100         | 100          | 100          | 100          | 100          | 100         | 100          |
| Fardh            | 90.5         | 90.5        | 100          | 100          | 100          | 91.3         | 95.2        | 90.9         |
| Khunaizi         | 95.2         | 90.9        | 100          | 95.5         | 90.5         | 95           | 90.5        | 90.5         |
| Qash             | 100          | 100         | 100          | 100          | 100          | 100          | 100         | 95.7         |
| Naghah           | 100          | 100         | 94.7         | 100          | 94.7         | 100          | 94.7        | 100          |
| Maan             | 100          | 100         | 100          | 100          | 100          | 100          | 93.8        | 100          |
| Average          | <b>97.62</b> | <b>96.9</b> | <b>99.12</b> | <b>99.25</b> | <b>97.53</b> | <b>97.72</b> | <b>95.7</b> | <b>96.18</b> |
| Accuracy         | <b>97.5</b>  |             | <b>99.2</b>  |              | <b>97.5</b>  |              | <b>95.8</b> |              |

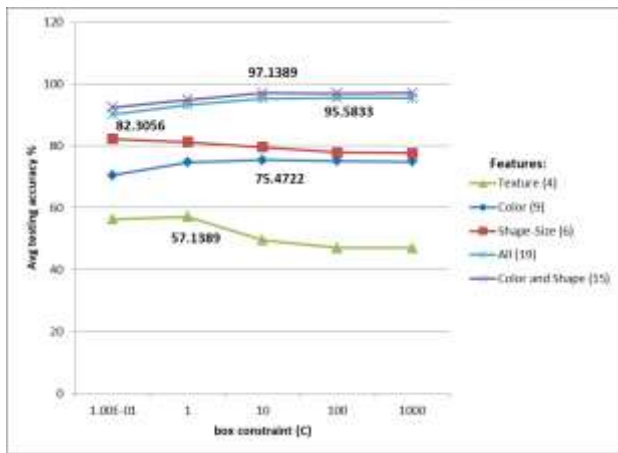


Fig. 8. Plots of the Relation between Average Testing Accuracy of different Features Versus the Box Constraint Parameter C.

### C. Setup of the K- Nearest Neighbor Classifier

The most basic data classification and pattern recognition classifier is the K- Nearest Neighbor, KNN [23]. In this approach, even though selecting the value of the constant K and distance metric is critical, this process does not need tuning many parameters, its efficiency is high. This is the highly recommended advantage of using KNN in object classification. As illustrated in Fig. 9, the lowest performance was achieved using KNN when we compared the obtained results with that of both ANN and SVM. The obtained performance was 53.33% when the texture features were used and K was seven. It was found that the value of K affects the performance in the case of using both color and shape-size features. When k=10, the performance was 70% when the color features were used while the achieved performance was 82.5% when the shape-sized features were used. The performance improved when both features were mixed together, color and shape-sized, and the value of K was five only.

### D. Classifiers Performance using Confusion Matrix

Confusion matrix has been used as a metric to measure the performance of different classification algorithms. It evaluates the accuracy of the networks classification system for training, validation and testing dataset. The column indicating the desired output represents the target class. The class of the output is in the rows of the matrix indicating the output of the system. The results of ANN confusion matrix for the features of both color and shape-size features when the function tansig and logsig are used in the hidden neurons respectively as shown in Fig. 10 and Fig. 11. Table I summarizes the highest achieved performances (%) of the classifiers (ANN with tansig hidden neurons, SVM, KNN with K=5) using color and shape-size features. We see clearly that when the function of tansig is used the neurons activation function in ANN, the performance of ANN is perfect (recall of 100). When C=10, SVM performance is less than the performance achieved using ANN. It could classify 4 out of 6 classes perfectly (recall of 100%).

At the recall of 90.5% and 94.7%, Khunaizi and Naghal are classified. SVM with a precision of 100% managed to classify Khalas, Qash, Naghal, and Maan. However, with a precision of 91.3% and 95%, respectively Fardh and Khunaizi are

classified. The lowest performance is given by KNN. Two classes are only classified perfectly (recall of 100%). At a recall with the range of 90.5% to 95.2%, the rest of the classes are classified.

### E. Time Complexity Analysis

As shown in Table II, the testing time (seconds) of different classifiers are presented. It is found that the time used for classification for both ANN and SVM is almost the same. ANN with tan-sigmoid hidden neurons was able to classify the testing samples in about 0.92 seconds/sample, which is considered as the lowest classification time in this paper. Logsig neural network and tansig neural network reach the highest accuracy in very close time. From Table II, we can judge that the times taken by both classifiers (ANN and SVM) are comparable. However, simulation results show that SVM takes much more time to achieve the classification. The KNN classifier achieves the highest classification time of 2.56 seconds /sample (i.e. the slowest algorithm) since it needs to calculate the distance from each testing sample to all the training samples when a classification is required.

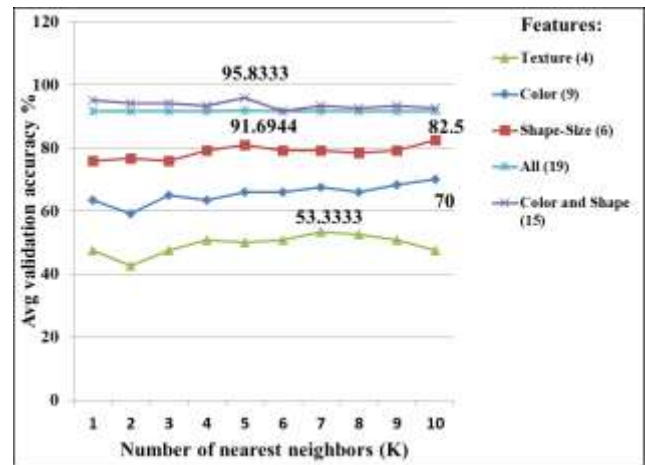


Fig. 9. Plots of the Relation between the Average Testing Accuracy of the KNN Classifier of different Features and the Number of Nearest Neighbors.

|   | 1            | 2            | 3            | 4            | 5             | 6            |               |
|---|--------------|--------------|--------------|--------------|---------------|--------------|---------------|
| 1 | 21<br>17.5%  | 0<br>0.0%    | 0<br>0.0%    | 0<br>0.0%    | 0<br>0.0%     | 0<br>0.0%    | 100%<br>0.0%  |
| 2 | 0<br>0.0%    | 21<br>17.5%  | 0<br>0.0%    | 0<br>0.0%    | 0<br>0.0%     | 0<br>0.0%    | 100%<br>0.0%  |
| 3 | 0<br>0.0%    | 0<br>0.0%    | 21<br>17.5%  | 0<br>0.0%    | 1<br>0.8%     | 0<br>0.0%    | 95.5%<br>4.5% |
| 4 | 0<br>0.0%    | 0<br>0.0%    | 0<br>0.0%    | 22<br>18.3%  | 0<br>0.0%     | 0<br>0.0%    | 100%<br>0.0%  |
| 5 | 0<br>0.0%    | 0<br>0.0%    | 0<br>0.0%    | 0<br>0.0%    | 18<br>15.0%   | 0<br>0.0%    | 100%<br>0.0%  |
| 6 | 0<br>0.0%    | 0<br>0.0%    | 0<br>0.0%    | 0<br>0.0%    | 0<br>0.0%     | 16<br>13.3%  | 100%<br>0.0%  |
|   | 100%<br>0.0% | 100%<br>0.0% | 100%<br>0.0% | 100%<br>0.0% | 94.7%<br>5.3% | 100%<br>0.0% | 99.2%<br>0.8% |
|   | 1            | 2            | 3            | 4            | 5             | 6            |               |

Fig. 10. The Results of ANN Confusion Matrix for Colour and Shape-Size Features using Tansig Hidden Neurons.



**Confusion Matrix**

|              |              |             |                |                |                |                |                |                |
|--------------|--------------|-------------|----------------|----------------|----------------|----------------|----------------|----------------|
| Output Class | 1            | 16<br>13.3% | 0<br>0.0%      | 1<br>0.8%      | 4<br>3.3%      | 0<br>0.0%      | 1<br>0.8%      | 72.7%<br>27.3% |
|              | 2            | 1<br>0.8%   | 14<br>11.7%    | 7<br>5.8%      | 0<br>0.0%      | 2<br>1.7%      | 1<br>0.8%      | 56.0%<br>44.0% |
|              | 3            | 1<br>0.8%   | 6<br>5.0%      | 9<br>7.5%      | 1<br>0.8%      | 1<br>0.8%      | 0<br>0.0%      | 50.0%<br>50.0% |
|              | 4            | 3<br>2.5%   | 0<br>0.0%      | 0<br>0.0%      | 15<br>12.5%    | 1<br>0.8%      | 0<br>0.0%      | 78.9%<br>21.1% |
|              | 5            | 0<br>0.0%   | 0<br>0.0%      | 1<br>0.8%      | 0<br>0.0%      | 8<br>6.7%      | 6<br>5.0%      | 53.3%<br>46.7% |
|              | 6            | 0<br>0.0%   | 1<br>0.8%      | 3<br>2.5%      | 2<br>1.7%      | 7<br>5.8%      | 8<br>6.7%      | 38.1%<br>61.9% |
|              |              |             | 76.2%<br>23.8% | 66.7%<br>33.3% | 42.9%<br>57.1% | 68.2%<br>31.8% | 42.1%<br>57.9% | 50.0%<br>50.0% |
|              | Target Class | 1           | 2              | 3              | 4              | 5              | 6              |                |

Fig. 11. The Results of ANN Confusion Matrix for Colour and Shape-Size Features using Logsig Hidden Neurons.

TABLE II. PROCESSING TIME OF DIFFERENT CLASSIFIERS FOR THE BEST ACHIEVED ACCURACY

| Classifier | Time (sec/sample ) |
|------------|--------------------|
| ANN-logsig | 0.939              |
| ANN-tansig | 0.917              |
| SVM        | 0.997              |
| KNN        | 2.56               |

Even though there is an increase in the achieved performance when a combination of features is used in the classification process, 19 features, as shown in Fig. 6 to 9, the computational overhead will be increased. Using the irrelevant date's features leads to the expletive of dimensionality and decreases the performance of the classification system as shown in Fig. 6 to 9, where the performance of the system in terms of classification accuracy is decreased after including texture features (19 features). By choosing an appropriate feature dimension (15 features as a combination of color features with shape features), balanced performance is achieved.

#### IV. CONCLUSION

The potential of CV systems (combination between color image processing and Machine-Learning techniques) in classifying automatically date fruit varieties in Oman has been investigated. Three ML techniques (ANN, SVM, and KNN) have been used and compared to each other in achieving the classification tasks. Intensive experiments and qualitative comparison are conducted among the developed approaches. Based on the combination of both color and shape-sized features give the highest performance accuracy in all approaches. This implies that date fruits have significant differences in colors and shape-size rather than textures. Meanwhile, the former combination represents an optimum solution of maximum accuracy with less number of features as

well as better processing time is achieved. The highest classification accuracy obtained by ANN, SVM, and KNN classifiers are 97.2581%, 97.1386%, and 95.83%, respectively. Thus, CV systems can be effectively used to classify date fruits and hence could be successfully used as an automatic date separator in the packaging date factories.

#### REFERENCES

- [1] FAO Statistics, "http://www.fao.org/faostat/en/#data/QC", Retrieved Sept. 19, 2019.
- [2] T. Gabriel, A. Manickavasagan & R. Al-Yahyai, "Classification of dates varieties and effect of motion blurring on standardized moment features", Journal of Food Measurement and Characterization, Springer, 6,1-4,2012.
- [3] R. Al-Yahyai & M.M. Khan, "Date palm status and perspective in Oman", in Date palm genetic resources and utilization, Springer, 207-240, 2015.
- [4] S.Ghnmimi, S.Umer, A.Karim, & A.Kamal-Eldin. "Date fruit (Phoenix dactylifera L.): An underutilized food seeking industrial valorization", NFS journal, 6, pp. 1-10, 2017.
- [5] A. Bhargava & A. Bansal, "Fruits and vegetables quality evaluation using computer vision: A review", Journal of King Saud University-Computer and Information Sciences, 2018.
- [6] A.Janecek, W.Gansterer, M.Demel & G.Ecker, "On the relationship between feature selection and classification accuracy", New challenges for feature selection in data mining and knowledge discovery, 90-105, 2008.
- [7] Ministry Of Information, "Popular varieties of dates grown in Oman", 2019, Retrieved from https://omaninfo.om/english/module.php?module=topics-howtopic&CatID=35&ID=3793.
- [8] K.Hameed, D.Chai, & A. Rassau, "Comprehensive review of fruit and vegetable classification techniques", Image and Vision Computing, 80, 24-44, 2018.
- [9] S.Naik & B.Patel. "Machine Vision based Fruit Classification and Grading-A Review", International Journal of Computer Applications, 170, 22-34, 2017.
- [10] H.A.aidar, H.Dong & N. Mavridis, "Image-based date fruit classification", IV International Congress on Ultra Modern Telecomm. and Control Systems. IEEE. St. Petersburg, Russia, 2012.
- [11] K.M.Alrajeh & T.A.A. Alzohairy, "Date fruits classification using MLP and RBF neural networks", International Journal of Computer Applications, 41(10), 36-41, 2012.
- [12] G. Muhammad, "Automatic Date Fruit Classification by Using Local Texture Descriptors and Shape-Size Features", European Modelling Symposium (EMS), IEEE, Pisa, Italy, 2014.
- [13] M. Ghulam. "Date fruits classification using texture descriptors and shape-size features. Elsevier, Engineering Applications of Artificial Intelligence, 37, 361-367, 2015.
- [14] A. Manickavasagan, N.H.Al-Shekaili, N.K. Al-Mezeini, M.S. Rahman & A. Guizani, "Computer vision technique to classify dates based on hardness", Journal of Agricultural and Marine Sciences, 22(1), 36-41, 2017.
- [15] G. Thomas, A. Manickavasagan, R. Al-Yahyai & L. Khriji, "Contrast Enhancement using Brightness Preserving Histogram Equalization Technique for Classification of Date Varieties", The Journal of Engineering Research, 11(1), 55-63, 2014.
- [16] T. Najeeb & M. Safar, "Dates Maturity Status and Classification Using Image Processing", International Conference on Computing Sciences and Engineering (ICCSE), Malaysia, 2018.
- [17] R.C. Gonzalez, R.E.Woods & S.L.Eddins, "Digital image processing using MATLAB", Tata McGraw Hill Education Private Limited, 2011.
- [18] V.K.Mishra, S.Kumar & N.Shukla, "Image Acquisition and Techniques to Perform Image Acquisition", Journal of Physical Sciences, Engineering and Technology. 9(1), 2229-7111, 2010.
- [19] Gongal A., Amatya S., Karkee M., Zhang Q., & Lewis K., "Sensors and systems for fruit detection and localization: A review", Computers and Electronics in Agriculture, 116, 8-19, 2015.

- [20] N.Otsu, "A threshold selection method from gray-level histograms", IEEE Transactions on systems, man, and cybernetics, 9, 62-66, 1979.
- [21] M.J.Zdilla, et al., "Circularity, solidity, axes of a best fit ellipse, aspect ratio, and roundness of the foramen ovale: a morphometric analysis with neurosurgical considerations", The Journal of craniofacial surgery, 27, 222, 2016.
- [22] M. M. Aly, "Survey on multiclass classification methods", Neural Netw, 19, 1-9, 2005.
- [23] R. Durgabai & Y.R. Bhushan, "Feature selection using ReliefF algorithm", International Journal of Advanced Research in Computer and Communication Engineering, 3, 8215-8218, 2014.

# Computational Analysis of Arabic Cursive Steganography using Complex Edge Detection Techniques

Anwar H. Ibrahim<sup>1</sup>, Abdulrahman S. Alturki<sup>2</sup>  
College of Engineering, Qassim University  
Mulaidah, Qassim Province  
Saudi Arabia

**Abstract**—Arabic language contains a multiple set of features which complete the process of embedding and extracting text from the compressed image. In specific, the Arabic language covers numerous textual styles and shapes of the Arabic letter. This paper investigated Arabic cursive steganography using complex edge detection techniques via compressed image, which comprises several characteristics (short, medium and Long sentence) as per the research interest. Sample of images from the Berkeley Segmentation Database (BSD) was utilized and compressed with a diverse number of bits per pixel through Least Significant Bit (LSB) technology. The method presented in this paper was evaluated based on five complex edge detectors (Roberts, Prewitt, Sobel, LoG, and Canny) via MATLAB. Canny edge detector has been demonstrated to be the most excellent solution when it is vital to perform superior edge discovery over-compressed image with little several facts, but Sobel appears to be better in term of the execution time for Long sentence contents.

**Keywords**—Arabic language; Berkeley Segmentation Database (BSD); Least Significant Bit (LSB); Roberts; Prewitt; Sobel; LoG; Canny

## I. INTRODUCTION

Protected of correspondence information between two nodes through a communication system ought to be secured from attack, consequently, numerous ways are utilized for that reason. Data covering up is utilized for forestalling an interloper to recognize them. Steganography is a method used to shroud the data and send them to the sender with changed over an arrangement to secure the data [1]. Another method provides high security to the data is the cryptography. It keeps data over the organization through changing over the plaintext into figure text. A few kinds of cryptography are utilized which are symmetric, topsy-turvy, and hashing. Cryptography calculation utilizes a similar key for encryption and unscrambling measures is called symmetric cryptography, while unbalanced cryptography utilizes various keys for encryption and decoding.

The transmission of a huge amount of information over the channel in a communications network involves high protection to secure the information. Consequently, steganography has a crucial function in communication to encapsulate such data throughout the edge and cover of an image. Steganography is practised by using those wishing to

deliver a mystery message or code through the image. While many valid methods make use of the steganography, such malware builders have additionally been located to use steganography to obscure the transmission of malicious code. Steganographic methods categorized into two classes: transform domain names and spatial methods [2]. Virtual Image for Steganography is one frequently [3]; which required two documents: The message to be embedded into the images for secretly hidden [4]. Steganography based data protection is crucial for confidential facts transfer. There are 3 fundamental requirements within the subject of digital steganography, each significant of mystery information is represented by the way of 8 bits and those bits are embedded inside the edge of the photo once creating the arithmetic processes on it. The first fundamental condition is capacity, which depends on the number of secret bits to be embedded in each cover pixel. The second constraint is robustness that avoids hidden information from attack. The third obligation is imperceptibility, typically intended by peak signal to noise ratio (SNR). Thus, Steganography technology is truly important in terms of information destiny of internet protection and privacy on open systems inclusive of the network which considered respectable when the faintness is high during secret data transmission while needing communication robustness [5]. Most of the existing methods using the Least Significant Bit (LSB) due to the redundant bits on the cover of the images embeds in the spatial area of the image with less effective in which it occurring clear misrepresentation [6,7].

## II. BACKGROUND OF STEGANOGRAPHY

As stated formerly, photographs are taken into consideration because of the maximum famous record formats used in steganography. They are acknowledged for constituting a non-causal medium, because of the possibility to get entry to any pixel of the photograph at random. Further, the hidden data should continue to be invisible to the attention. Fig. 1 represents the general data protection scheme Classification tree.

Steganography is another way of having messages secured during data communication. The end goal of steganography and cryptography is the same but they have different methods. Steganography does not change data or message format and keeps its actual data present while cryptography keeps the data secret by converting it into an unreadable form. The drawback

of the cryptographic approach lies in the existence of original data as the original data was encrypted. Steganography techniques, therefore, provide additive protection to cryptographic techniques. This offers an additional layer of protection for the message during data communication, with the combination of both.



Fig. 1. Classification Tree of Security Systems [8].

### III. STEGANOGRAPHY FEATURES

#### A. Why Steganography is Important?

Nowadays, Steganography can be utilized to cover up hidden information interior to other files so that the parties expecting to induce the message indeed knows a mystery message exists. Steganography gets to be the foremost basic approach utilized to secure the information. The word Steganography implies, hide the secret information just like e-content or advanced arrange. It points to conceal the mystery information eventually between two parties and make it not visual to the third party and without any doubts around the existing of any covered up data. There are a few sorts of Steganography have been isolated into two mediums, which are advanced Steganography and normal dialect Steganography. Computerized Steganography is the craftsmanship that bargains with the computerized medium, for illustration, picture, video, and sound, whereas characteristic dialect Steganography bargains with the content records. Indeed even though computerized Steganography has the most considerations by the analysts, in any case, the content is the foremost basic information that has to be secured since most of the documentation will be within the content such as sending basic data or doling out pressing appointments [9,10]. Also, Steganography in the natural language is divided into two groups, which are linguistic Steganography and Steganographic text. Linguistic Steganography is about the text (a secret message) concealed in a text medium [8]. In the meantime, auto-Steganography adjusts the document format or a specific character, without altering the context of the sentences [11, 12].

Hiding the data involves certain strategies using the natural language Steganography. The sort has its techniques which are, word-rule-based and feature-based techniques used by the researchers in text Steganography [13]. Meanwhile, linguistic Steganography uses five techniques, such as synonymous

substitution, syntactic substitution, semantic substitution, PCFG, and hybrid technique.

#### B. Steganography Features based Technique

So far, numbers of image Steganography methods have been implemented, the simplest approach implemented is the LSB substitution technique. The least important bits of the picture pixel are used in this technique for embedding hidden message bits [14,15].

The feature-based approach works, for example, with the shape or style of the text, by modifying the size or type of font. This strategy will make readers believe that no improvements are made in Text so that the reader cannot notice the hidden message embedded in the cover [16].

### IV. EMBEDDING AND EXTRACTION TECHNIQUES

A Steganography embedding and extraction technique refer to all items with redundancy in the data. People frequently transfer digital images through email and other Internet communication and JPEG is one of the most popular image formats. Also, Steganography systems seem more appropriate for the JPEG format because the systems run in a transformed space and are not affected by visual attacks [8].

An image's edge representation greatly decreases the amount of data to be processed, but it preserves important knowledge about the shapes of the objects in the picture. This description of an image is easily implemented in a large number of object recognition algorithms used in computer vision along with other applications for image processing. The edge detection technique's main property is its ability to determine the exact edge line with reasonable orientation, as well as more literature on edge detection has been available in the last three decades. On the other hand, the efficiency of the edge detection techniques is not yet measured by any typical performance index. The efficiency of an edge detection technique is often judged individually and independently based on its application. The literature includes several edge detection techniques for image segmentation. This section looks at the most widely used discontinuity-based edge detection techniques. These are Roberts edge detection methods, Sobel Edge detection, Prewitt edge detection, Kirsh edge detection, Robinson edge detection, Marr-Hildreth edge detection, LoG edge detection and Canny Edge detection [17].

#### A. Roberts Edge Detection

Lawrence Roberts (1965) implements Roberts Edge Detection. It performs a simple, easy to calculate, 2-D measurement of the spatial gradient on an image. This approach emphasizes high-spatial-frequency regions that often correspond to edges. The operator input is a grayscale image the same as the output is the most common use for this technique, pixel values at each point in the output reflect the approximate complete magnitude of the input image's spatial gradient at that point [18]. The Roberts Edge filter is used to detect edges that are based on sequentially applying a horizontal and vertical filter as shown in table A and B. Both filters refer to the image and are summed up to create the final picture.

| A.Horizontal Filter |   |
|---------------------|---|
| 0                   | 1 |
| 1                   | 0 |

| B. Vertical Filter |   |
|--------------------|---|
| 1                  | 0 |
| 0                  | 1 |

**B. Prewitt Edge Detection**

Prewitt in (1970) proposed by Rafael C for edge detection technique. It was found that as a correct way to estimate the magnitude and orientation of the edge. Although different gradient edge detection requires a time-consuming calculation to estimate the direction from the magnitudes in the x and y directions, the compass edge detection obtains the direction directly from the kernel with a high grade of reacting. This gradient-based edge detector is estimated for eight directions in the 3x3 area. All 8 convolution masks are calculated. A complication mask is then selected, i.e. for the largest module [19].

| A.Horizontal Magnitudes |    |    |
|-------------------------|----|----|
| -1                      | -1 | -1 |
| 0                       | 0  | 0  |
| 1                       | 1  | 1  |

| B. Vertical Magnitudes |   |   |
|------------------------|---|---|
| -1                     | 0 | 1 |
| -1                     | 0 | 1 |
| -1                     | 0 | 1 |

**C. Sobel Edge Detection**

The Sobel method introduced by Rafael C (1970) for the segmentation of the image finds the edges using the Sobel approximation to the derivative. It precedes the edges at the points where the gradient is the highest. The Sobel technique performs a 2-D spatial gradient quantity on the image, thus highlighting regions with a high spatial frequency corresponding to the edges. In general, it is used to find the estimated absolute gradient magnitude for each gradient [18].

| A.Horizontal Magnitudes |    |    |
|-------------------------|----|----|
| -1                      | -2 | -1 |
| 0                       | 0  | 0  |
| 1                       | 2  | 1  |

| B. Vertical Magnitudes |   |    |
|------------------------|---|----|
| -1                     | 0 | -1 |
| -2                     | 0 | 2  |
| -1                     | 0 | 1  |

**D. LoG Edge Detection**

The Laplacian of Gaussian (LoG) was introduced by Marr (1982) for edge detection. Laplacian filters are derivative filters that are used to detect areas of rapid change (edges) in images as shown at equation (1) for the relation of the second derivative of image  $f(x,y)$  [20].

$$\nabla^2 f = \frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2} \tag{1}$$

Since the derivative filters are very sensitive to noise, it is common to smooth the image using a Gaussian filter before applying the Laplacian filter. Subsequently, the operator of Laplace can detect both edges and noise (isolated, out-of-range), it may be desirable first to smooth the image with a Gaussian kernel of width.

**E. Canny Edge Detection**

The Canny edge detector is an operator of the edge detection, using a multi-stage algorithm to detect a wide range of edges in images. It was founded in 1986, by John F. Canny. Canny also developed a computational edge detection theory which explains why the technique works. In industry one of the popular edge detection techniques is the Canny edge detection technique. It was first created by John Canny for his Master's thesis at MIT in 1983, and it still outperforms many

of the newer algorithms that have been developed. To find the edges by separating the noise from the image before finding the edges of the image, Canny is a very important method [17].

**V. A PROPOSED METHOD FOR EMBEDDING AND EXTRACTION STEGANOGRAPHY**

The proposed method of image steganography intends to improve/increase the cover image's hiding capacity. The suggested approach uses the inclusion of the edge region in the cover picture to add more hidden information than embedding it into the non-edge region. The method of embedding and extraction in the proposed work is Widespread introduced in two steps. Transmitter and receiver with high secure user name and password. Fig. 2 shows the steps of embedded and extracted Arabic text steganography.

An important aspect of the techniques used in this thesis is that it used to embed a text in colour images. Fig. 2 shows the button to select an algorithm to be considered from six types (Log, Robert, Prewitt, Canny, Demirel, and Sobel) to perform the proposed method. The user has an option to choose only one algorithm to embed the text as shown in Fig. 2.

The method of modification of the Least Significant Bit (LSB) is used very effectively in the Image Steganography technique. To enhance its wide-ranging application, this research paper proposes that, as a first step, this method can also be applied to images that have undergone edge detection techniques the reason why the edge detected image is different from the original image so that any edge detection changes can be made.

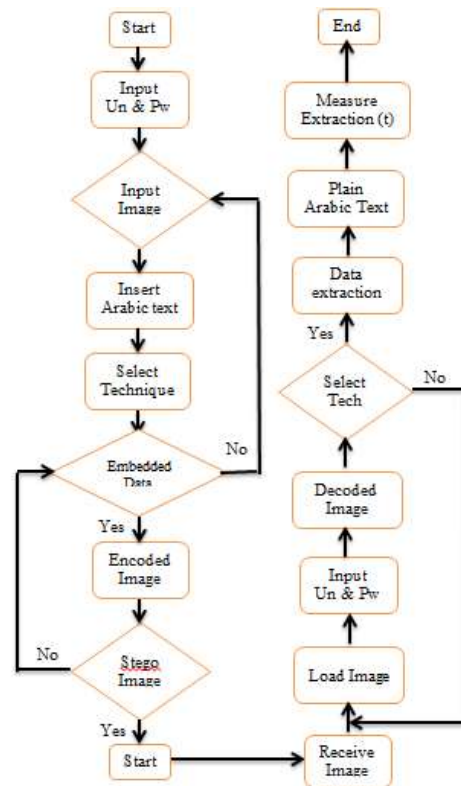


Fig. 2. State Transition Diagram of Proposed Steganography.

A. The Unicode Standard

Unicode is a universal standard that was adopted for the production, encoding, and handling of digital texts represented in most of the writing systems in the world from 1987 until now [21]. The Unicode standard, in other words, is an encoding scheme designed to facilitate the worldwide display, processing and exchange of texts with different languages and technical discipline.

B. Arabic Text Hiding Criteria

When programmers build a text hiding algorithm there are many things to remember. In recently implemented algorithms, however, the fundamental requirements can be easily found: invisibility, embedding power, robustness, and security [22]. For an active or passive warden, respectively, the contact medium through which the stego-image is being transmitted can be noisy or noiseless.

C. Arabic Text Embedding and Extraction

The information stream of image format (MPEG, JPG and SVG) were mostly made out of head data, image encoded information on vector stream utilizing movement remuneration forecast method with a least significant bit (LSB) created movement vector information stream. The design-based strategies include changing a few highlights of the spread content of text embedding, for example, text dimension, style, shading, and so forth that could be modified to cover mystery image. In the extraction process, the inverse method should be applied to extract the data we less time according to the image format and capacity.

D. Algorithm Selection

A significant part of the procedures utilized in this theory is that it used to insert a book in shading pictures. Fig. 2 shows the catch to choose calculation to be considered from six sorts (Robert, Prewitt, Canny, Log, and Sobel) to play out the proposed technique. The client has an alternative to pick just a single calculation to install the content as appeared in Fig. 2.

VI. RESULTS AND DISCUSSION

The research proposed is designing robust algorithms to perform the Results:

Increase robustness by embedding random bit in the edge of the image, employing value shift technique based Matlab algorithms. The concept is only Embedding bits in the consecutive pixels of the samples in the selected area.

Table I displays the steganography overall description techniques and provides a clear understanding that each technique has its advantages and inconveniences. Each is unique to the application and the program requirement justifies the use of such a system with given parameters. Sobel is one of the most successful techniques for the systems requiring fast computation without having to maintain data.

Fig. 3 represented the level of the embedding and extraction time of Arabic text steganography which are limit intangibility, accessibility, and reading time.

Fig. 4 shows the inexact implanting time and the all-out limit character for every calculation through three

configurations; it was discovered that the best calculation for stage installing is finished by the Canny method.

The primary edge assessment was limit according to the image format utilization, the representation capacity measured based on three image format, based on the classification, accessibility, and capacity that utilized the total measure.

TABLE I. OVERALL DESCRIPTION TECHNIQUES OF STEGANOGRAPHY

| Parameters            | Pewit          | Sobel          | Canny          | Roberts        | LoG            |
|-----------------------|----------------|----------------|----------------|----------------|----------------|
| Computational         | Complex        | Simple         | Complex        | simple         | Complex        |
| Signal to Noise ratio | Low            | Low            | High           | High           | Low            |
| Texture based image   | Less efficient | High efficient | Less efficient | Less efficient | High efficient |
| Embedding time        | less time      | efficient time | time-consuming | time-consuming | efficient time |
| Extraction time       | less time      | efficient time | time-consuming | time-consuming | time-consuming |
| Security              | more           | less           | more           | less           | more           |
| Capacity              | more           | less           | more           | more           | less           |

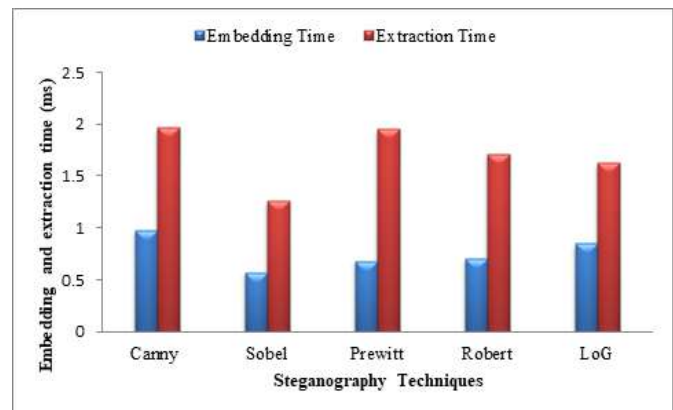


Fig. 3. Embedding and Extraction Time.

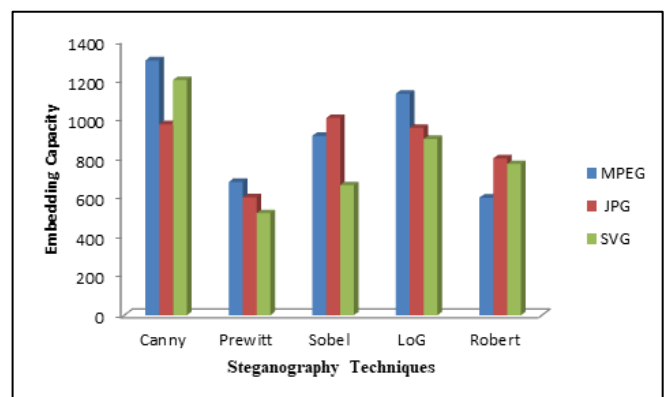


Fig. 4. Embedding Capacity.

## VII. CONCLUSION

The proposed strategy applies the edge discovery method on the spread picture and disorderly guide. We utilized the edge location method with the end goal that Sobel channel, it is utilized to give various pieces utilized in installing. Likewise, we utilized tent guide, it gives the area of pixels which used to inserting pieces. The trials and result directed to affirm that stego-picture inserts the two mystery bits if pixel present edges or installs the one mystery bit if pixel, not the current edge. We utilized the arrangement of the measures to discover the proportion of clamour between pictures. This contextual investigation presents an examination of existing content concealing methods, particularly on those concentrated on adjusting the basic attributes of advanced instant message for Arabic text. The results outlined a scope of crucial rules, applications, and assaults covering the content concealing territory to clarify the current embedding and extraction time challenges in the image steganography. Additionally, the study concludes the three significant evaluation process (Extraction time, Embedding time and effect of the image format) of Arabic text concealing procedures dependent on the best way to deal with spread instant messages to decode the mystery bits, in particular, the maximum capacity with less time for embedding and extraction time. Based on the critical condition for the best sketched out the restrictions and qualities of every classification to show their effectiveness in different image format. Also, we assessed the as of late proposed approaches concerning the key measures to feature their advantages and disadvantages. Finally, we have suggested some of the rules and bearings that legitimacy further consideration in future works.

## REFERENCES

- [1] Hassanain Raheem Kareem, Hadi Hussein Madhi, Keyan Abdul-Aziz Mutlaq. Hiding encrypted text in image steganography. Periodicals of Engineering and Natural Sciences. Vol. 8, No. 2, June 2020, pp.703-707.
- [2] A. Cheddad, J. Condell, K. Curran, P. Mc Kevitt, "Digital image steganography: Survey and analysis of current methods", Signal Processing, Volume 90, Issue 3, March 2010, Pages 727-752.
- [3] A. Cheddad, J. Condell, K. Curran, and P. McKevitt, "Digital image steganography: Survey and analysis of current methods", Signal Processing, vol. 90, no. 3, pp. 727-752, 2010.
- [4] B. Li, J. He, J. Huang, and Y. Q. Shi, "A survey on image steganography and steganalysis", Journal of Information Hiding and Multimedia Signal Processing, vol. 2, no. 2, pp. 2073-4212, 2011.
- [5] Kehui Sun, "Chaotic secure communication: principles and technologies," Berlin Boston De Gruyter, 2016.
- [6] Ramadhan J. Mstafa and Khaled ElleithyKhaled Elleithy, "A Novel Video Steganography Algorithm in the Wavelet Domain Based on the KLT Tracking Algorithm and BCH Codes. 015 IEEE Long Island Systems, Applications, and Technology Conference At: NYC, May 2015.
- [7] G. Sugandhi and C. P. Subha . Efficient steganography using least significant bit and encryption technique . 2016 10th International Conference on Intelligent Systems and Control (ISCO). 7-8 Jan. 2016.Wavelet Domain Based on the KLT Tracking Algorithm and BCH Codes", 2015.
- [8] Karrar Abdallah Mohammed, Int. Journal of Computer Science & Mobile Computing, Vol.7 Issue.10, October- 2018, pp. 25-32.
- [9] R. Din and S. Utama, "Critical Review of Verification and Validation Process in Feature-Based Method Steganography," in Int. Conf. E-Commerce, 2017, pp. 15-19.
- [10] S. S. Iyer and K. Lakhtaria, "New robust and secure alphabet pairing Text Steganography Algorithm," Int. J. Curr. Trends Eng. Res., vol. 2, no. 7, pp. 15-21, 2016.
- [11] H. T. Ciptaningtyas, R. Anggoro, and M. B. A. Fadhillah, "Text Steganography on Sundanese Script using Improved Line Shift Coding," in 2018 International Electronics Symposium on Knowledge Creation and Intelligent Computing (IES-KCIC), 2018, pp. 254-261.
- [12] S. Utama, R. Din, and M. Mahmuddin, "The Performance Evaluation of Feature-Based Technique in Text Steganography," J. Eng. Sci. Technol., vol. 12, pp. 169-180, 2017.
- [13] R. Din, R. Bakar, S. Utama, J. Jasmis, and S. J. Elias, "The evaluation performance of a letter-based technique on text steganography system," Bulletin of Electrical Engineering and Informatics, vol. 8, no. 1, pp. 291-297, 2019.
- [14] Deepali Singla and Mamta Juneja. New Information Hiding Technique using Features of Image. Journal of Emerging Technologies in Web Intelligence 6(2). 237-242. 2014.
- [15] Farah Qasim Ahmed Alyousuf, and Roshidi Din. Analysis review on feature-based and word-rule based techniques in text steganography. Bulletin of Electrical Engineering and Informatics.Vol. 9, No. 2, April 2020, pp. 764~770. (Main).
- [16] A. Westfeld and A. Pfitzmann, "Attacks on Steganographic Systems," Proc. Information Hiding 3rd Int'l Workshop, Springer Verlag, pp. 61–76, 1999.
- [17] Muthukrishnan.R and M.Radha. International Journal of Computer Science & Information Technology (IJCSIT) Vol 3, No 6, Pp: 259 – 267. Dec 2011.
- [18] Rafael C. Gonzalez, Richard E. Woods & Steven L. Eddins (2004) Digital Image Processing Using MATLAB, Pearson Education Ltd. Ltd, Singapore.
- [19] Inas Jawad Kadhim, Peter James Vial and Brendan Halloran. A comprehensive survey of image steganography: Techniques, Evaluations, and trends in future research. Neurocomputing. Volume 335, 28 March 2019, Pages 299-326.
- [20] Srinivas B.L, Hemalatha and Jeevan K.A. Edge Detection Techniques for Image Segmentation. International Journal of Innovative Research in Computer and Communication Engineering. Vol. 3, Special Issue 7, October 2015.
- [21] Robert Lockwood and Kevin CurranKevin Curran. Text based steganography. International Journal of Information Privacy Security and Integrity 3(2):134. January 2017.
- [22] Milad Taleby Ahvanooy, Qianmu Li, Jun Hou, Ahmed Raza Rajput and Chen Yini. Modern Text Hiding, Text Steganalysis, and Applications: A Comparative Analysis. Entropy 2019, 21, 355.

# The Most Efficient Classifiers for the Students' Academic Dataset

Ebtehal Ibrahim Al-Fairouz<sup>1</sup>, Mohammed Abdullah Al-Hagery<sup>2</sup>

Department of Computer Science, College of Computer  
Qassim University, Buraydah, Saudi Arabia

**Abstract**—Educational institutions contain a vast collection of data accumulated for years, so it is difficult to use this data to solve problems related to the progress of the educational process and also contribute to achieving quality. For this reason, the use of data mining techniques helps to extract hidden knowledge that helps in making the decisions necessary to develop education and achieve quality requirements. The data of this study obtained from the College of Business and Economics at Qassim University. Three of the classifiers were compared in this study Decision Tree, Random Forest and Naïve Bayes. The results showed that Random Forest outperforms other algorithms with 71.5% of Precision, 71.2% F1-score, and also it got 71.3% of Recall and Classification Accuracy (CA). This study helps reduce failure by providing an academic advisor to students who have weaknesses in achieving a high-Grade Point Average (GPA). It also helps in developing the educational process by discovering and overcoming weaknesses.

**Keywords**—Data mining; student performance; classification algorithms; evaluation

## I. INTRODUCTION

Over the past few years, the world has seen rapid progress in technology. This progress led to the accumulation of information and data and its availability in all sectors systems, such as educational, health, social, and others. This data can be used to discover and analyze the obstacles and problems facing these sectors by employ data mining techniques [1], [2].

Data mining is considered as an interdisciplinary approach and an essential step in knowledge discovery [1]. It is used to extract useful and hidden information from large databases. Through the use of data mining tasks, it is also possible to answer questions that cannot be known through other techniques such as queries or reports[3]. The essential function of data mining is to use various algorithms such as classification, clustering, regression and association rules to discover hidden patterns that help in many important decisions [2]. The classification technique is a supervised learning task, the data in this method is classified into pre-defined classes. It is a frequently used method for creating models that are used to predict futuristic patterns. Examples of classification algorithms are Decision Tree, Naïve Bayes, Logistic Regression, K-nearest neighbor, Neural Networks, etc.[4]. The process of discovering knowledge involves several steps: collect data, clean data, pre-processing data, and then using data mining techniques.

This paper aims to use data mining techniques to examine student performance, the classification algorithms will be used

to classify the student's GPA through the use of historical data from the College of Business and Economics at Qassim University from 2014 to 2018.

Moreover, the importance of these results was situated in their practical application in educational institutions, as recommended, to use the best classification model on any academic data such as data of university students, institutes, schools, etc., for students' performance prediction. Furthermore, classify students in many aspects that assist the institution to enhance the educational process.

The remaining of the paper is divided into the following sections: the second section reviews the work related to this study, the third section the methodology of the study, the fourth section includes the results and the discussion, whereas the last section presents the conclusions.

## II. RELATED WORKS

The research by Ramaphosa et al. [5] was about primary schools students from four cities in South Africa. The goal of their study was to recognize a predictive algorithm to detect learners' performance and make appropriate decisions for improvement. They analyzed the data using WEKA tool by employing classification algorithms namely Naive Bayes, BayersNet, J48 and JRip. They proved by their results that the J48 algorithm is the best model of prediction when compared to other algorithms by 99.13% classification accuracy. Ultimately, they reported that their study assists the schools in early discovering the academic performance of learners and enable stakeholders to improve the results of weak students.

According to the study by Abu Amrieh et al. [6], they found there is a relationship between the academic performance of the student and the behavior of the student (student interaction with the e-learning system). In their study, they used the dataset from Kalboard 360 which contains 500 records and 16 attributes. The student performance was predicted by applying classification algorithms which are Decision Tree, Artificial Neural Network and Naive Bayes by using WEKA tool. Besides, for improving classifier performance, they implemented ensemble methods namely Boosting, Random Forests and Bagging. Their results showed a robust relationship among academic performance and the behavior of the student, where the predictive model with behavioral attributes achieved higher accuracy than the predictive model without behavioural attributes. Furthermore, they observed an improvement in accuracy when they used ensemble techniques. Finally, they explained that this model



supports stakeholders in understanding students and identify weaknesses and develop their learning process in addition to reducing failure. Another study by Al-Noshan et al. concentrated on a set of important factors affecting the students' performance in the first year of the university [7], likewise Al-Rofiyee et al. [8]. Also, in [9], [10], [11] the authors compared the classifiers accuracy but using a medical dataset.

On the other hand, Rahman and Islam [12] applied four traditional classification algorithms which are K-NN algorithms, Naïve Bayes, Decision Tree and Artificial Neural Network algorithms. Besides, they used bagging ensemble method, boosting ensemble method and at last ensemble filtering technique, which helps extract hidden knowledge from student data that makes it easier for educational establishments to improve their quality of education. Their results indicated that the ensemble filtering technique obtained the best accuracy among all the algorithms.

Roy and Garg [13] applied the J48, Naïve Bayes, and MLP. The results showed that J48 obtained 73.92% accuracy which was the highest accuracy among the used algorithms. The objective of their study was to identify and predict the factors that affect student academic performance, where the performance of students can be affected by different attributes such as related to school, social and demographic.

The aim of the study by Guerra et al. [14] was to predict the performance of students in specific courses. In their study, they applied Decision Tree techniques on the dataset of IFMS in Brazil from 2012 to 2015 by utilizing WEKA tool. Their results showed that the J4.8 classification algorithm achieves the best results with cross-validation and pruning by 75.8%.

Ahmed and Elaraby [15] applied classification techniques to predict students' performance in the final assessment. They collected the dataset from the information system department from the year 2005 to 2010. The tool used in this work is WEKA by applying a Decision Tree algorithm (ID3). Through their results, they explained that their study help improves student performance as well as identify students who need the advice to guide them and make the appropriate decision.

The study by Tsiakmaki et al. [16] aimed to predict students' marks in the final exams of the courses of the second semester based on the first semester grades. They used a dataset from the Business Administration department of the TEI of Western Greece from 2013 to 2017 which contains 592 students. They only applied methods of regression using the WEKA tool, namely Linear Regression, Bagging, M5 algorithm, Gaussian processes (GPs), M5-Rules, Sequential Minimal Optimization (SMO), Random Forest and 5NN. The evaluation measure used in their study was MAE. After all the experiments they had done on the data set, they concluded that all the algorithms had achieved fair accuracy.

Pérez et al. [17] presented their initial results that prediction of attrition of students from a large dataset of Systems Engineering (SE) undergraduate students after six years of registration at a Colombian university, the dataset includes 762 students. In their study, they applied four algorithms which are Decision Tree, Random Forest, Naive

Bayes and Logistic Regression. Then, they found that SE courses performance is linked to mathematics and physics courses performance where they obtained the best AUC from Random Forest by 97% in the 3rd semester. These results showed them plainly that the courses which related to Systems Engineering have a dominant effect in predict dropout.

The objective of the study by Adekitan and Salau [18] was to perform predictive analysis to determine the final CGPA of graduation using their GPA of the previous three years as well as to define the class to which the student belongs at graduation. They applied six algorithms namely Decision Tree, Tree Ensemble, Random Forest, Naive Bayes, Logistic Regression and the Probabilistic Neural Network on the dataset which was gathered from Covenant University at Nigeria for the engineering students. The tools used in their study were KNIME and MATLAB. Their results demonstrated that the logistic regression obtained the best accuracy by 89.15%. Hence, they pointed out that students' results can be predicted the last year of their study using their performance in the previous three years. On the other hand, a few of studies included large datasets with records ranging from 14,333 records to 21,314 records.

Yulianto et al. [19] applied classification algorithms to student data to identify features that affect student achievement. Besides, they expected that the results of the analysis would be able to find the reasons that led to the delay of some students in the study period. They used two models of algorithms k-Nearest Neighbor and Decision Tree C4.5. They concluded that the k-Nearest Neighbor got better accuracy than the other. Quinn and Gray [20] used data from the Moodle to predict students' grades whether they will succeed or fail in the course. They applied the classification algorithms Random Forest, Gradient Boosting, k-Nearest Neighbours and Linear Discriminant Analysis using R. They summarized that the use of data from Moodle gives the ability to early detection of students at risk. The aim of the study by Walia et al. [21] was to build classification models to predict academic performance for students through the use of classification algorithms Naive-Bayes, Decision Tree, Random-Forest, JRip, and ZeroR. The results indicated that the school and study time were influential factors in the students' final grades.

A comparison of classification algorithms has been applied in several fields like emotion classification, precipitation, Spatial modelling of storm dust provenance etc. Fauziastuti et al. [22], classified students' graduation on time or overtime, by used two classification algorithms to compare their performance: Naive Bayes Classifier and K-Nearest Neighbor, but using a small dataset.

A similar research paper was in emotion classification to find the best classifiers amongst a set of classifiers [23], whereas in this paper we are concentrating on the extraction of the hidden knowledge embedded in the academic data of undergraduate students by a set of classifiers to find the best classifier for getting the hidden knowledge from this kind of data. The study achieved by Lazri et al. [24], focused on estimating precipitation from Meteosat Second Generation images, by combining six models of classification. They also used a linear regression model. Likewise, a study conducted by

Gholam et al. [25], was applied eight classification algorithms, for spatial maps to predict the source of dust in Khuzestan.

### III. METHODOLOGY

The methodology used focuses on the use of classification algorithms in analyzing student performance to discover hidden patterns that help officials make the necessary decisions in the educational process.

The knowledge discovery process consists of four phases: data collection, pre-processing, data mining technique (classification), and interpretation of results, as in Fig. 1. The tool used in this study is the Orange data mining platform.

#### A. Data Collection

Data was collected from the College of Business and Economics from 2014-2018, which contains 72259 records for male and female students from several majors. The dataset contains the following attributes: Semester, Course code, Course name, CRD hours, Gender name, Entry date, Confirmed mark, Grade, Cumulative GPA (CGPA), Semester GPA (SGPA), Student status, Major name and Student level.

#### B. Data Preprocessing

Real data is usually incomplete and inconsistent due to individual errors or computer errors. Therefore, before starting to use data mining techniques, pre-processing of data is required. The process of data pre-processing includes first, clean the data from missing values, as the data was cleaned by the Orange program through the use of Impute widget. After cleaning the data, the number of records became 52,430. Second, data transformation, where the students' GPA was classifying into five categories as follows:

- 1) Excellent (GPA >=4.5)
- 2) Very Good (GPA >=3.75)
- 3) Good (GPA >=2.75)
- 4) Average (GPA >=2.00)
- 5) Fail (<2.00)

This classification was done using the Feature Constructor Widget on the Orange platform based on the CUM\_GPA attribute and the new attribute was named as Class\_GPA. The first class is Excellent, second class is Very Good, third class is Good, the fourth class is Average and the fifth class is Fail.

#### C. Data Mining Techniques (Classification)

The classification method is known as supervised classification, where, the data are organized into given known

classes. The dataset in classification is divided into a training dataset and test dataset. The classification algorithm is trained through a training dataset to build a model and test the model by test dataset since this model is used later to classify new data [26]. For example, predict students' performance by using the classification of GPA to good or bad. Algorithms used in classification such as Decision Trees, Random Forests or Bayes models. The classification techniques that were used in this study include Decision Tree, Naïve Bayes and Random Forest.

The classification algorithms are connected to Predictions widget to shows models predictions on the data. Hence, to evaluate the performance of models, we have focused on four different metrics which are CA, F1-score, Precision and Recall, given in Equations (1-4), where true positives (TP), true negatives (TN), false positives (FP) and false negatives (FN). Moreover, the dataset was divided into a training set, and a test set by using a fixed proportion of data in Data Sampling widget, 75% of the data were used for the training set and 25% for the test set. The target variable is Class\_GPA. Fig. 2 shows the model workflow of the classification task.

$$CA = \frac{TP+TN}{TP+FP+FN+TN} \quad (1)$$

$$Precision = \frac{TP}{TP+FP} \quad (2)$$

$$F1\text{-score} = \frac{2 \times Precision \times Recall}{Precision+Recall} \quad (3)$$

$$Recall = \frac{TP}{TP+FN} \quad (4)$$

Where,

TP: the model predicts correctly the positive class.

TN: the model predicts correctly the negative class.

FP: the model predicts incorrectly the positive class.

FN: the model predicts incorrectly the negative class.

Predictions widget was used to recognize model predictions on the data. Data Sampler widget is used to sample the data by using a fixed proportion of data. The dataset was divided into training data by 75% and test data by 25%. The data sample was sent to three algorithms widgets by Data Sampler widget so that they can produce the corresponding model; after that, the models were sent into Predictions widget while the remaining data was directly sent to Predictions widget.

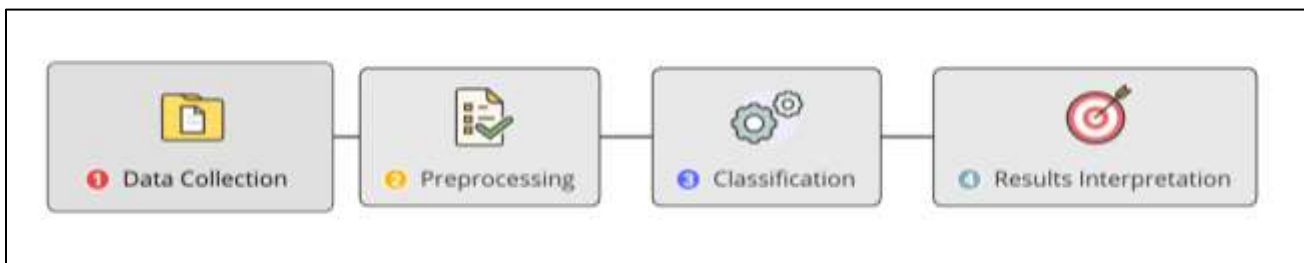


Fig. 1. Knowledge Discovery Process.

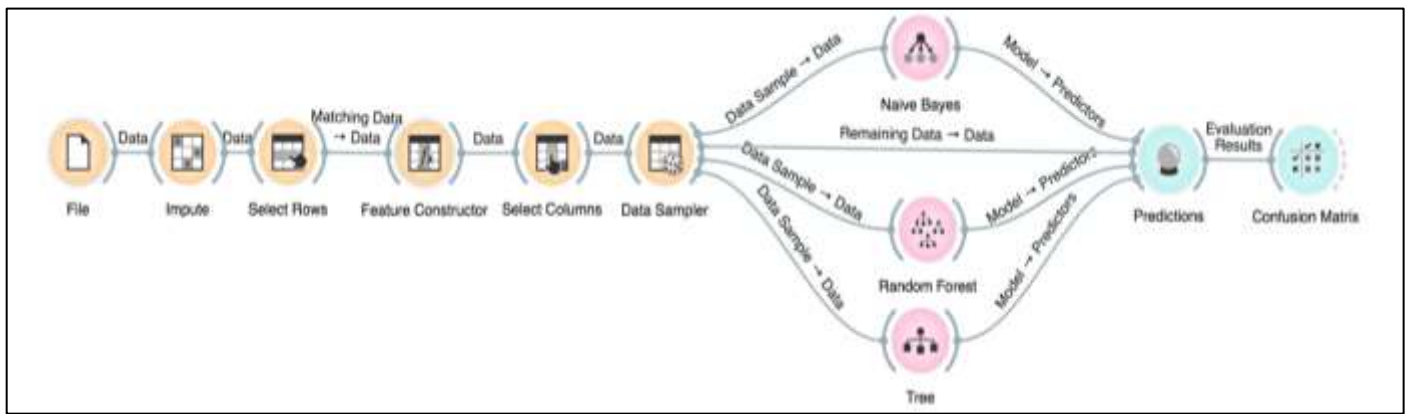


Fig. 2. The Model Workflow of Classification.

#### IV. RESULTS AND DISCUSSION

In this section, we presented the evaluation results of Decision Tree, Random Forest and Naïve Bayes. Fig. 3 shows the Predictions widget, which presented data with added predictions and the results of testing classification algorithms.

The widget received the dataset and then constructed a predictive model with Decision Tree, Random Forest and Naïve Bayes widgets, and it found the probabilities in predictions.

Table I presents the evaluation results of the classification. As we can see from the table, the Random Forest was the best classifier with 71.5% of Precision, 71.3% Recall and CA, and also it got 71.2% of F1-score. While the worst algorithm was the Naïve Bayes with 60.5% of Precision, 59.4% of CA and Recall, 59.5% of F1-score.

Confusion Matrix will be displayed that aims to assess the predictive performance of the models for each class to recognizing prediction of TP, FP, TN, and FN. The class labels are Excellent, Good, Acceptable and Fail. Table II, Table III, and Table IV illustrate the confusion matrices for Naïve Bayes, Random Forest and Decision Tree, respectively.

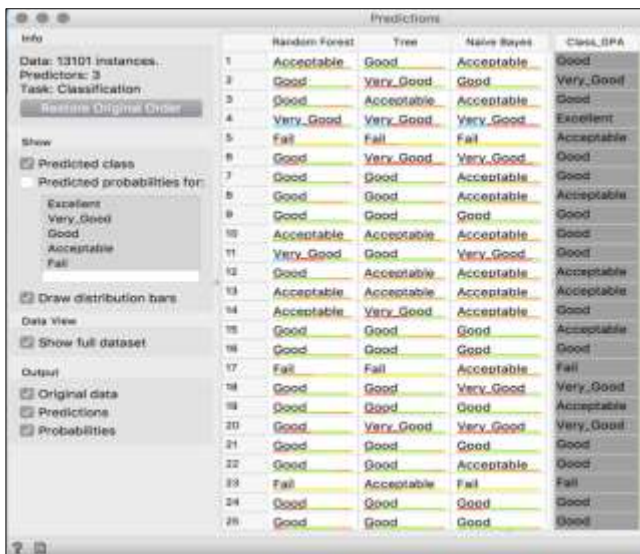


Fig. 3. Predictions Widget.

TABLE I. THE EVALUATION RESULTS OF PREDICTION

| Model         | CA    | F1-score | Precision | Recall |
|---------------|-------|----------|-----------|--------|
| Random Forest | 0.713 | 0.712    | 0.715     | 0.713  |
| Decision Tree | 0.698 | 0.697    | 0.699     | 0.698  |
| Naïve Bayes   | 0.594 | 0.595    | 0.605     | 0.594  |

TABLE II. CONFUSION MATRIX OF THE NAÏVE BAYES

|        |            | Predicted |           |           |      |            | Σ     |
|--------|------------|-----------|-----------|-----------|------|------------|-------|
|        |            | Model     | Excellent | Very Good | Good | Acceptable |       |
| Actual | Excellent  | 700       | 276       | 13        | 0    | 0          | 989   |
|        | Very Good  | 601       | 1254      | 615       | 28   | 3          | 2501  |
|        | Good       | 113       | 875       | 3200      | 916  | 136        | 5240  |
|        | Acceptable | 4         | 62        | 1021      | 2021 | 491        | 3599  |
|        | Fail       | 0         | 0         | 3         | 159  | 610        | 772   |
|        | Σ          | 1418      | 2467      | 4852      | 3124 | 1240       | 13101 |

TABLE III. CONFUSION MATRIX OF THE RANDOM FOREST

|        |            | Predicted |           |           |      |            | Σ     |
|--------|------------|-----------|-----------|-----------|------|------------|-------|
|        |            | Model     | Excellent | Very Good | Good | Acceptable |       |
| Actual | Excellent  | 617       | 344       | 27        | 1    | 0          | 989   |
|        | Very Good  | 162       | 1573      | 756       | 10   | 0          | 2501  |
|        | Good       | 10        | 412       | 4140      | 671  | 7          | 5240  |
|        | Acceptable | 0         | 18        | 986       | 2468 | 127        | 3599  |
|        | Fail       | 0         | 0         | 7         | 208  | 557        | 772   |
|        | Σ          | 789       | 2347      | 5916      | 3358 | 691        | 13101 |

TABLE IV. CONFUSION MATRIX OF THE DECISION TREE

|        |            | Predicted |           |           |      |            | Σ     |
|--------|------------|-----------|-----------|-----------|------|------------|-------|
|        |            | Model     | Excellent | Very Good | Good | Acceptable |       |
| Actual | Excellent  | 711       | 265       | 13        | 0    | 0          | 989   |
|        | Very Good  | 254       | 1593      | 641       | 13   | 0          | 2501  |
|        | Good       | 31        | 601       | 3989      | 608  | 11         | 5240  |
|        | Acceptable | 2         | 41        | 1045      | 2316 | 195        | 3599  |
|        | Fail       | 0         | 0         | 18        | 218  | 536        | 772   |
|        | Σ          | 998       | 2500      | 5706      | 3155 | 742        | 13101 |

Fig. 4 shows the evaluation of the three models Naïve Bayes, Random Forest and Decision Tree. The models were evaluated using four measures CA, F1-score, Precision and Recall Through the figure, we notice that the Random Forest outperform in all measures than other algorithms, followed by decision tree algorithm. As for the worst model was Naïve Bayes.



Fig. 4. Evaluation of the Models.

## V. CONCLUSIONS

Educational institutions often require an analysis of student data to obtain useful knowledge that contributes to enhancing the learning process in addition to achieving quality in education. For this reason, data mining techniques were used to extract hidden knowledge from student data, and a comparison was made between three classifiers, Naïve Bayes, Random Forest and Decision Tree. Experimental results showed that Random Forest exceeded other classifiers with the accuracy of 71.3%, followed by the Decision Tree by 69.8%, then the last classifier was the Naïve Bayes by 59.4%. This study helps to know students' performance in advance by relying on previous results to improve their achievement in the future. Also, educational institutions must provide an academic adviser to failed students to enhance their academic performance.

## ACKNOWLEDGMENTS

The authors would like to thank the College of Business and Economics at Qassim University that provided the data required for this research.

## REFERENCES

- [1] J. Han, M. Kamber, and J. Pei, *Data Mining: Concepts and Techniques*, 3rd ed. Elsevier, 2012.
- [2] E. C. Abana, "A decision tree approach for predicting student grades in Research Project using Weka," *Int. J. Adv. Comput. Sci. Appl.*, vol. 10, no. 7, pp. 285–289, 2019, doi: 10.14569/ijacsa.2019.0100739.
- [3] N. Rehman, "Data Mining Techniques Methods Algorithms and Tools," *Int. J. Comput. Sci. Mob. Comput.*, vol. 6, no. 7, pp. 227–231, 2017, [Online]. Available: <http://www.ijcsmc.com>.
- [4] A. H. Awlla, "Performance Analysis and Prediction Student Performance to build effective student Using Data Mining Techniques," *UHD J. Sci.*

- Technol., vol. 3, no. 2, p. 10, Jun. 2019, doi: 10.21928/uhdjt.v3n2y2019.pp10-15.
- [5] K. I. M. Ramaphosa, T. Zuva, and R. Kwuimi, "Educational Data Mining to Improve Learner Performance in Gauteng Primary Schools," in 2018 International Conference on Advances in Big Data, Computing and Data Communication Systems, *icABCD 2018*, Aug. 2018, pp. 1–6, doi: 10.1109/ICABCD.2018.8465478.
- [6] E. A. Amrieh, T. Hamtini, and I. Aljarah, "Mining Educational Data to Predict Student's academic Performance using Ensemble Methods," *Int. J. Database Theory Appl.*, vol. 9, no. 8, pp. 119–136, 2016, doi: 10.14257/ijtda.2016.9.8.13.
- [7] A. Abdulrahman Al-Noshan, M. Abdullah Al-Hagery, H. Abdulaziz Al-Hodathi, and M. Sulaiman Al-Quraishi, "Performance Evaluation and Comparison of Classification Algorithms for Students at Qassim University," *Int. J. Sci. Res.*, vol. 8, no. 11, pp. 1277–1282, 2018, doi: 10.21275/ART20202907.
- [8] N. Al-Mufadi and M. A. Al-Hagery, "Using prediction methods in data mining for diabetes diagnosis," *Using Predict. Methods Data Min. Diabetes Diagnosis*, vol. 1, no. 4, p. 2014, 2014.
- [9] M. Abdullah Al-Hagery, A. Saleh Alfaiz, F. Suliman Alorini, and M. Saleh Althunayan, "Knowledge Discovery in the Data Sets of Hepatitis Disease for Diagnosis and Prediction to Support and Serve Community," *Int. J. Comput. Electron. Res.*, vol. 4, no. 6, pp. 118–125, 2015, [Online]. Available: <http://ijcer.org>.
- [10] S. Al-qarzaie, S. Al-odhaibi, B. Al-saeed, and M. Al-hagery, "Using the Data Mining Techniques for Breast Cancer Early Prediction," *Symp. Data Min. Appl.*, vol. 1, no. May, p. 2014, 2014.
- [11] M. Abdullah and H. Al-Hagery, "Classifiers' Accuracy Based on Breast Cancer Medical Data and Data Mining Techniques," *Int. J. Adv. Biotechnol. Res.*, vol. 7, no. 2, pp. 976–2612, 2016, [Online]. Available: <http://www.bipublication.com>.
- [12] M. Hasibur Rahman and M. Rabiul Islam, "Predict Student's Academic Performance and Evaluate the Impact of Different Attributes on the Performance Using Data Mining Techniques," in 2nd International Conference on Electrical and Electronic Engineering, *ICEEE 2017*, 2018, no. December, pp. 1–4, doi: 10.1109/ICEEE.2017.8412892.
- [13] S. Roy and A. Garg, "Predicting academic performance of student using classification techniques," in 2017 4th IEEE Uttar Pradesh Section International Conference on Electrical, Computer and Electronics, *UPCON 2017*, 2017, vol. 2018-Janua, pp. 568–572, doi: 10.1109/UPCON.2017.8251112.
- [14] M. S. Guerra, H. A. Neto, and S. A. Oliveira, "A Case Study of Applying the Classification Task for Students Performance Prediction," *IEEE Lat. Am. Trans.*, vol. 16, no. 1, pp. 172–177, Jan. 2018, doi: 10.1109/TLA.2018.8291470.
- [15] A. Badr, E. Din, and I. S. Elaraby, "Data Mining : A prediction for Student' s Performance Using Classification Method," *World J. Comput. Appl. Technol.*, vol. 2, no. 2, pp. 43–47, 2014, doi: 10.13189/wjcat.2014.020203.
- [16] M. Tsiakmaki, C. Pierrakeas, G. Kostopoulos, S. Kotsiantis, G. Koutsonikos, and O. Ragos, "Predicting university students' grades based on previous academic achievements," in 2018 9th International Conference on Information, Intelligence, Systems and Applications, *IISA 2018*, Jul. 2019, pp. 1–6, doi: 10.1109/IISA.2018.8633618.
- [17] B. Pérez, C. Castellanos, and D. Correal, "Predicting student drop-out rates using data mining techniques: A case study," in *Communications in Computer and Information Science*, 2018, vol. 833, pp. 111–125, doi: 10.1007/978-3-030-03023-0\_10.
- [18] A. I. Adekitan and O. Salau, "The impact of engineering students' performance in the first three years on their graduation result using educational data mining," *Heliyon*, vol. 5, no. 2, p. e01250, Feb. 2019, doi: 10.1016/j.heliyon.2019.e01250.
- [19] L. D. Yulianto, A. Triayudi, and I. D. Sholihati, "Implementation Educational Data Mining For Analysis of Student Performance Prediction with Comparison of K-Nearest Neighbor Data Mining Method and Decision Tree C4.5," *J. Mantik*, vol. 4, no. 1, pp. 441–451, May 2020.
- [20] R. J. Quinn and G. Gray, "Prediction of student academic performance using Moodle data from a Further Education setting," *Irish J. Technol.*

Enhanc. Learn., vol. 5, no. 1, pp. 1–19, Oct. 2019, doi: 10.22554/ijtel.v5i1.57.

- [21] N. Walia, M. Kumar, N. Nayar, and G. Mehta, “Student’s Academic Performance Prediction in Academic using Data Mining Techniques,” *SSRN Electron. J.*, pp. 1–5, Apr. 2020, doi: 10.2139/ssrn.3565874.
- [22] V. T. Fauziastuti and L. A. Rakhman, “A Review of Students’ Graduation Classification: A Comparison of Naive Bayes Classifier and K-Nearest Neighbour,” in *1st International Multidisciplinary Conference on Education, Technology, and Engineering (IMCETE 2019)*, Mar. 2020, pp. 219–221, doi: 10.2991/assehr.k.200303.052.
- [23] M. AbdullahAl-Hagery, M. AbdullahAl-Assaf, and F. MohammadAl-Kharboush, “Exploration of the best performance method of emotions classification for arabic tweets,” *Indones. J. Electr. Eng. Comput. Sci.*, vol. 19, no. 2, pp. 1010–1020, 2020, doi: 10.11591/ijeecs.v19.i2.pp1010-1020.
- [24] M. Lazri, K. Labadi, J. M. Brucker, and S. Ameer, “Improving satellite rainfall estimation from MSG data in Northern Algeria by using a multi-classifier model based on machine learning,” *J. Hydrol.*, vol. 584, p. 124705, May 2020, doi: 10.1016/j.jhydrol.2020.124705.
- [25] H. Gholami, A. Mohamadifar, and A. L. Collins, “Spatial mapping of the provenance of storm dust: Application of data mining and ensemble modelling,” *Atmos. Res.*, vol. 233, p. 104716, Mar. 2020, doi: 10.1016/j.atmosres.2019.104716.
- [26] R. Lawrance and V. Shanmugarajeshwari, “An assay of teachers’ attainment using decision tree based classification techniques,” in *Proceedings of IEEE International Conference on Circuit, Power and Computing Technologies, ICCPCT 2017*, Apr. 2017, pp. 1–6, doi: 10.1109/ICCPCT.2017.8074382.

#### AUTHORS’ PROFILE

**Ebtehal Ibrahim Al-Fairouz:** received her BSc in Computer Science from the Qassim University, Buraydah, KSA. She is a teaching assistant in the Department of Management Information System (MIS) at the College of Business and Economics (CBE) and a Master’s student in Computer Science Department, Qassim University, KSA. Her research interests include data mining, data analytics, data visualization and machine learning.



**Mohammed Abdullah Al-Hagery:** received his BSc in Computer Science from the University of Technology in Baghdad Iraq-1994. He got his MSc. in Computer Science from the University of Science and Technology Yemen-1998. Al-Hagery finished his Ph.D. in Computer Science and Information Technology, (Software Engineering) from the Faculty of Computer Science and IT, University of Putra Malaysia (UPM), November 2004. He was the head of the Computer Science Department at the College of Science and Engineering, USTY, Sana’a from 2004 to 2007. From 2007 to this date, he is a staff member at the College of Computer, Department of Computer Science, Qassim University, Buraydah, KSA. He published more than 31 papers in various international journals. Dr. Al-Hagery was appointed the head of the Research Centre at the Computer College, and a council member of the Scientific Research Deanship Qassim University, KSA from September 2012 to October 2018. Currently, he is teaching the master degree students and a supervisor of four master thesis. He is a jury member of several PhD and master thesis, as an internal and external examiner in his field of his specialist.

# An Empirical Study of e-Learning Interface Design Elements for Generation Z

Hazwani Nordin<sup>1</sup>, Dalbir Singh<sup>2</sup>, Zulkefli Mansor<sup>3</sup>

Center for Software Technology and Management, Faculty of Technology and Information Science  
National University of Malaysia, 43600 Bangi, Selangor, Malaysia

**Abstract**—E-learning is the latest evolution of electronic-based learning that creates, fosters, delivers and facilitates the learning process anytime and anywhere with the use of interactive network technology. The use of e-learning as a learning platform makes users want a high quality of interface design to interact with the e-learning system. Interface design that meets students' needs and expectations may increase their involvement and satisfaction towards e-learning, especially generation Z students. However, interface design is always being criticized and has become a part of issues that cause the failure of e-learning. Lack of understanding about students' cultural background and preferences towards e-learning interface design are the major factors that contribute to this phenomenon. To ensure the success of e-learning, a model of interface design specifically for generation Z students' culture that consists of interface elements and design characteristics need to be developed. So, this study aimed to address this subject by identifying e-learning interface elements and design characteristics from existing literature, confirming the elements and design characteristics and discovering related elements for e-learning interface design from generation Z students' perspective. This study used semi-structured for a focus group interview that included seven students. The focus group interview involved three main steps which were sampling, protocol and research instruments. This study validated several interface elements and design characteristics that contribute to the model of e-learning interface design. The findings could guide the interface designer in designing e-learning interface for generation Z students.

**Keywords**—e-Learning; interface design; generation Z; culture; focus group

## I. INTRODUCTION

The education system in Malaysia has undergone many changes as information and communication technology (ICT) is widely used in educational institutions consistent with the modernization of the globalization era. The development of ICT affects the education system where learning in the classroom is integrated with electronic learning or e-learning as an effective teaching and learning method [1][2]. The integration of e-learning has opened up the opportunities for an active communication between lecturers and students especially generation Z students who dominate educational institutions and known as a generation of technology literate or digital natives [3][4].

Generation Z students were born between 1995 and 2010 [5][6]. Generation Z has been shaped by technological advanced since they were born and known a world with

internet and mobile devices. Technology has become their nature as they always stay connected to social media such as Facebook, Twitter, YouTube and WhatsApp as part of their interaction and learning [6]. It becomes their culture as they always depend on technology in their daily life. Therefore, generation Z students should enjoy e-learning as a learning platform because they can learn and stay connected with their friends at the same time without disengaged from technology and social media. In addition, generation Z students love to do something fast because they want to get information and answers immediately from various sources, even from unreliable sources [7]. However, their involvement and engagement towards e-learning are still low and unsatisfying [9][11]. Besides, the rate of drop out courses in e-learning is higher than classroom learning despite the growing popularity of e-learning [8][12][13]. One of the reasons for low student involvement towards e-learning is due to poor e-learning interface design [10][8][14][15]. For example, users are not familiar with the design features would have a tendency to disconnect with learning experience because interface designers had failed to develop engaging e-learning environment [16]. Previous studies discovered that adapting culture into interface design could improve the look and feel of the e-learning system [17][18].

Previous researchers have studied cultural differences in interface design to determine how culture can affect the interface design [19][20][21]. For example, the minimalist interface design of Google is accepted in most Western countries but has failed to attract users in South Korea (SK) where Google is far behind the local search engine Naver.com [22][23]. Google failed to understand users need in SK with the minimalist concept because the interface design of Naver.com is more complex and colourful with interesting animations, images, links to other web pages and more features on the main page [22][23]. Complex interface design is common features on SK web pages, but in Western countries, users find it hard to accept the complexity of interface design. Interface designed with appropriate culture looks more attractive and functional to the targeted users. However, many developers usually develop interface design based on the culture of the country and not individual [24][25], in addition, less prescriptive about interface design elements in the existing models.

Less attention to the importance of interface elements and design characteristics complicate the process of developing cultural interface design [26][27][28]. Hofstede's cultural model is commonly used to explain culture based on six

dimensions which are power distance (PDI), masculine (MAS), individualistic (IDV), uncertainty avoidance (UAI), long term orientation (LTO) and indulgence (IVR). Hofstede has conducted research on more than 60 countries, including Malaysia. The result of Hofstede study in Malaysia shows that PDI scores highest compared to other dimensions and highest in the world compared to other countries [29][30]. Various interface developers refer to the result of Hofstede's study [29] whereby interface design is more highlighted on PDI dimension. Examples of high PDI in the interface design are using the corporate colour of institutions, the image of leaders, organization charts that show hierarchy and special titles in the interface design.

The emergence of a new generation causes cultural interface design needs to be revised constantly in order to make sure the interface meets the users' preferences. To refer the result of Hofstede study in Malaysia in designing the interface for the new generation is questionable whether or not the result is still valid. This is because the Hofstede study was more than 40 years ago and obsolete. Thus, this paper aims to address the issue of e-learning interface design based on generation Z students' preferences and needs by identifying and collating the interface elements and design characteristics from theoretical and empirical perspectives. Interface elements such as colour, graphics and typography are known as cultural markers. Each marker has its own design characteristics which would suit the needs of generation Z students [31][19][32]. These identified elements or cultural markers and design characteristics of the interface can help developers to develop e-learning system in the future. Besides, there are several limitations in this study, such as non-comparative study because it focuses on e-learning interface design for generation Z only. This study also focused on a few interface elements only such as colour, graphic, layout and navigation. It is because of the limitation of past research about the other elements. Moreover, the existing interface design model fewer details in describing the design characteristics of each element.

This paper is organized as follow. Section II presents interface elements and design characteristics that were obtained from the literature. Section III presents the methodology used in the empirical study, particularly the focus group interview. Section IV presents the results from the focus group interview, and Section V concludes the paper by summarising the whole paper and future work.

## II. LITERATURE REVIEW

### A. Cultural Markers of Interface Design

Interface element is a key factor in bridging the gap between users and the system [21]. Cultural adaptation in interface design is important where interface elements carry a different meaning for different cultures [10][21]. There are many interface elements, and it is impossible to put all of them in the e-learning interface. Past research on e-learning interface design is insufficient to determine the elements that usually used in e-learning interface design. Therefore, past researches on various interface design such as in educational and banking websites are carefully studied in order to determine commonly used interface elements and to provide

clear functionality to generation Z students. Table I show frequently used elements of interface design in past researches.

Table I shows several interface elements that are categorized to form one element. The main reason for classification is because they have the same functionality. For example, image, symbol, logo and metaphor element are classified as graphical elements because they can function as a visual element and make interface design more attractive and easy to understand [53][34]. Structure, information organization [46][54] and appearance [50][51][55] are categorized as layout elements because their definition about the arrangement of elements and information such as the position of the navigation bar, symbol, logo, date & time and image are same with the definition of the layout. Meanwhile, link element usually in text, button or graphic are used to navigate from page to page so link is classified under navigation element. Typography is textual style of appearance including the arrangement of written language to make the text readable and appealing [38][56]. Past researchers classify font types and writing direction as part of the language, which is also typography features, including size and textual distance [51][56]. Therefore, the language element is categorized with typography element. In addition, language can be recognized easily, directly accessible and less defined in cultural context [38][51]. Table I also shows the frequency of use for interface elements in previous studies. Five elements that always been used in interface are colour, graphic, navigation, layout and typography are proposed as generation Z students' preferences for e-learning interface design. The following is a brief description of how these elements or cultural markers interpreted in various cultures.

Colour is a crucial design element that provides a visual representation of interaction [21]. The right colour combination can highlight the interface design layout, facilitate the discrimination of screen components, highlight the differences and can make the interface design more attractive. Furthermore, colours can be used to communicate specific meaning in different cultures, whereby they use relevant colours in their interface design [57]. Colour can affect users' expectation, and influence user's perceived towards navigation, link or content of interface design [57][25]. This is because colour can give a different meaning to users from different cultures. For example, red colour is said to bring luck and happiness in China, but in Japan red is considered a colour of anger and symbolizes danger in United State (US) [21][25].

Graphical elements are always used in interface design. There are several elements categorized as graphical elements such as image, icon, symbol and logo. Graphical elements can help students to understand the content without reading text. Cultural differences cause users to interpret graphics from different perspectives [57]. The system may receive more attention if interface design uses culturally appropriate images [58]. For example, images of high ranking people from institutions' management on web pages were interpreted as fraud by students from Western countries but appreciated in Eastern countries [53].

TABLE I. LITERATURE REVIEW ON INTERFACE ELEMENTS

| Elements        | Sources  | Frequency  |    |
|-----------------|--|--|----|
| Colour          | [33], [34], [35], [36], [32], [37], [31], [38], [39], [40], [23], [41], [42], [43], [44] | 15   |    |
| Graphic         | Image  | [34], [36], [24], [32], [38], [39], [40], [45], [41], [46], [43], [44]             | 20 |
|                 | Symbol   | [47], [40], [45], [43]   |    |
|                 | Logo   | [48], [45]   |    |
|                 | Metaphor   | [49], [50]   |    |
| Layout          | Layout   | [35], [47], [36], [32], [37], [38], [39], [45], [41], [42], [44]                   | 16 |
|                 | Structure  | [23], [46]   |    |
|                 | Information Organization   | [51]   |    |
|                 | Appearance   | [49], [50]   |    |
| Navigation      | Navigation   | [34], [48], [33], [24], [32], [37], [49], [38], [50], [40], [45], [23], [46], [44] | 20 |
|                 | Links  | [48], [37], [45], [42], [46], [43]   |    |
| Typography      | Typography   | [34], [35], [33], [37], [38], [40], [42]   | 15 |
|                 | Language   | [24], [51], [38], [39], [52], [41], [42], [46]                                     |    |
| Audio           | [39], [45], [46], [43], [44]   | 5  |    |
| Video           | [45], [46], [44]   | 3  |    |
| Animation       | [52], [45], [41], [43], [44]   | 5  |    |
| Games           | [46]   | 1  |    |
| Menu            | [48], [33]   | 2  |    |
| Search button   | [48], [33], [52]   | 3  |    |
| Help & support  | [48], [32]   | 2  |    |
| Error           | [24]   | 1  |    |
| Map             | [52], [46]   | 2  |    |
| Attention       | [48]   | 1  |    |
| Calendar        | [48]   | 1  |    |
| Feedback        | [48]   | 1  |    |
| Personalization | [52]   | 1  |    |

Layout involves the process of placing and arranging interface elements such as graphics, links colours and information architecture on the interface [10][17][25]. The layout of visual objects on the screenplays a role in user's perception towards e-learning interface design and important for gaining trust from users when using e-learning [16]. Research conducted by [16] to find the best layout to increase student's learning experience in e-learning. In their study, there were two types of layout which focus on functional and chronological layout. They discovered that a combination of the functional and chronological layout was the ideal layout in e-learning interface design.

Good navigation allows smooth movement in handling task and exploring e-learning [59] while poorly developed navigation would cause the student trapped and move into the same space [60]. This cause student to take longer time and feel frustrated to continue using e-learning. A research conducted by [24] shows that many Germany websites offered flexible paths for navigation while in Vietnam navigation paths are less flexible. Besides, several cultures use text as links to other pages instead of images or icons as links [21]. Generally, smooth navigation needs to be well designed with navigational aid such as arrows, icons or buttons for users to understand and use it.

Past researches show that every culture has different preferences for interface design. Culture coordination on interface design is based on cultural dimensions in order to understand the design characteristics of each interface elements. In this study, the Hofstede cultural model is referred to adapt culture into interface design.

### B. The Relationship between Hofstede Cultural Dimensions and Elements of Interface Design in Malaysia

Hofstede cultural model is popular among researchers because it is more comprehensive than other cultural models [61][19][57]. Hofstede studied culture involving more than 60 countries, including Malaysia. The result of Hofstede study is shown in Fig. 1.

According to the past researchers, Malaysia has the highest score of PDI compared to other countries [29][30] whereby a perfect score is obtained from this study. As a result, interface design is more prominent with PDI features compared to the other dimensions [29]. For example, the website of National University of Malaysia (UKM) uses national colours and focus on university expertise and ranking top 150 universities in the world [62]. Table II shows relations between Hofstede cultural dimensions and interface design characteristics in Malaysia based on previous researches.

Table II shows interface design characteristics in Malaysia based on past research. Previous research depends on the score of Hofstede study for developing interface design [67][25] including in Malaysia. Hofstede study in Malaysia happened more than 40 years ago, and users' preferences towards e-learning interface design had also changed. Thus, the interface element and design characteristics would be verified among generation Z students in the empirical study.

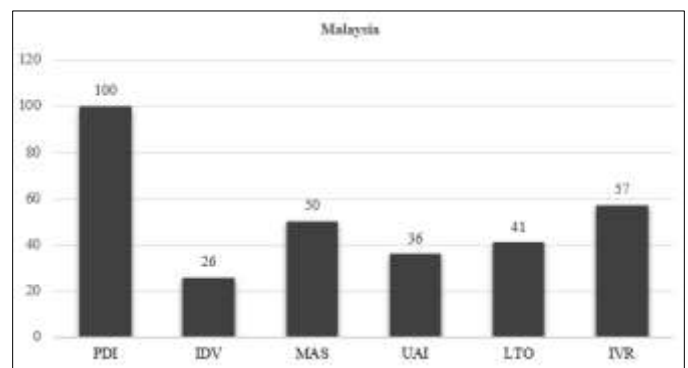


Fig. 1. Result of Hofstede study in Malaysia.



TABLE II. RELATIONS BETWEEN HOFSTEDE CULTURAL DIMENSIONS AND INTERFACE DESIGN IN MALAYSIA

| Cultural dimensions  | Design characteristics   | Sources                            |
|----------------------|--|------------------------------------|
| High PDI             | <ul style="list-style-type: none"> <li>Using images of important people or institutional building to represent history.</li> <li>Complete information about the board of directors and institutional management in an organizational chart.</li> <li>Special titles must be displayed.</li> <li>Displaying awards, hierarchical information and people with rank or authority are key features in the media.</li> </ul>  | [18], [63], [64], [54], [41], [69] |
|                      | <ul style="list-style-type: none"> <li>Use white or black colour as the background to focus on the name and image of the institution.</li> <li>Use the official colour of the institution.</li> <li>Highly structured information.</li> </ul>  |                                    |
| High IDV             | <ul style="list-style-type: none"> <li>Information provided on the website is complete and usable.</li> </ul>  | [64]                               |
| Low IDV (Collective) | <ul style="list-style-type: none"> <li>Use group images.</li> <li>Using many different colours (colourful).</li> <li>Provide a solution for the error message.</li> </ul>  | [41], [48], [49]                   |
| Low UAI              | <ul style="list-style-type: none"> <li>Complex menu with various paths and need to scroll.</li> <li>Vertical page scrolling.</li> <li>A large amount of information is placed on the first or main page.</li> <li>Using animation through text movement.</li> <li>Using a text link to the main page from any pages.</li> <li>Provide feedback on user location.</li> <li>Main menu on the left panel of the page.</li> <li>Use colours for classification.</li> <li>Non-linear navigation.</li> </ul> | [41], [54], [65], [48]             |
| High UAI             | <ul style="list-style-type: none"> <li>Using a site map.</li> <li>Place important items at the top and middle page.</li> </ul>   | [65]                               |
| High MAS             | <ul style="list-style-type: none"> <li>The main text colour is black while the blue and red colour to highlight important information.</li> </ul>  | [41]                               |
| Low MAS (Feminine)   | <ul style="list-style-type: none"> <li>Using a site map.</li> <li>Highlight critical data.</li> </ul>  | [66]                               |
| High LTO             | <ul style="list-style-type: none"> <li>Long path while navigating.</li> <li>Left text alignment.</li> </ul>  | [41] [69]                          |

### III. METHOD

The empirical study began with a qualitative method by conducting a focus group interview among generation Z students. The main purpose of the focus group interview was to verify the identified interface design elements that were derived from the literature as well as discovering other interface elements. Besides, the focus group interview was also aimed to find out about design characteristics preferred by generation Z students which would be used later in the survey. There were several steps taken into consideration before the focus group interview.

#### A. Sampling

The selection of focus group informants was based on several criteria, such as whether they had prior experience in using e-learning and undergraduate student who study in UKM besides being born in 1995 and above. Therefore, purposive sampling technique was adopted. Seven informants who had actively experience in e-learning from various faculties were invited to participate in this focus group. The profiles of seven informants are shown in Table III.

TABLE III. INFORMANTS' PROFILE FOR FOCUS GROUP INTERVIEW

| Informant Code | Faculty | Experience in E-Learning |
|----------------|---------|--------------------------|
| Inf_1          | A       | 2                        |
| Inf_2          | B       | 2                        |
| Inf_3          | B       | 2                        |
| Inf_4          | C       | 2                        |
| Inf_5          | D       | 2                        |
| Inf_6          | E       | 2                        |
| Inf_7          | F       | 2                        |

Based on Table III, all informants who joined the focus group interview were second-year students. They were selected because they have more experience using e-learning compared to the first-year student who is still new to e-learning. The selection of informants from various faculties was to gather opinions from different perspectives and not focussed on one faculty only.

#### B. Instrument

Interview questions were used as the instrument for the focus group interview. To ensure the suitability of the questions, all questions drafted were reviewed by two experts. Both experts were academicians and had experience in e-learning and interface design. The questions were divided into two parts, A and B. Part A covers about informants' experience and difficulty while handling e-learning. A semi-structured interview was conducted in this study; thus, extended questions were asked based on informant's answer in order to understand clearly what has been discussed. Part B revolves around five proposed interface elements from literature which are colour, graphic, typography, layout and navigation. Every informant had to give their opinion about each interface element in order to find out about design characteristics preferred by generation Z students. Table IV summarised and described five elements of the interface that were included in interview questions.

#### C. Protocol

A protocol is the rules provided to conduct a session of focus group interview in a more organized manner. Before the interview session, informants were asked via the online instant messaging application (*WhatsApp*) about the date and time that suits with all of the participants. A week before the focus group interview, the invitation letter containing a brief description of interview objectives, date, time and venue were sent to each informant. The focus group interview session was recorded. Before starting the interview, informants were given a few minutes to read and understand the consent letter before signing it. Generally, the consent letter was about the

confidentiality of the information was given during the interview that indicates the name of the informants would not be revealed if the report is made public.

TABLE IV. INTERVIEW QUESTIONS DESCRIPTION

| Elements   | Description   |
|------------|---|
| Colour     | To verify whether colours are important for generation Z students in e-learning and colour combinations to improve e-learning interface design. |
| Graphic    | To confirm the types of graphic that should be included and its position in e-learning interface design.  |
| Typography | To identify the types of font to use as heading and common text and font size.  |
| Layout     | To validate the position of each element that is preferred by generation Z students to make captivating e-learning interface design.            |
| Navigation | To affirm whether navigation could affect generation Z students while navigating e-learning and links in the form of text, icon or button.      |

#### IV. RESULT AND DISCUSSION

The results from the focus group interview are presented in the following section. The design characteristics pertaining to respective elements are shown in bold.

##### A. Colour

Past researchers stated that colour is important in e-learning interface design in order to make it look interactive. All informants agreed that colour should be given extra attention in e-learning interface design. Informants also suggested a colour combination that should be used in e-learning interface design. Below are some of the comments from informants:

- “Interface design looks more attractive if the combination of striking and pastel colour that suitable for both male and female are used in e-learning.” – Inf\_1.
- “Use more than one colour and combine with animated icons would look more attractive, and background colour must be consistent on each page to look formal.” – Inf\_3.
- “Use colour corporate of the university as theme colour on e-learning interface design to make e-learning more professional looking and must be consistent on every pages” – Inf\_4.
- “The background colour on each page should be different. No need to be consistent with one colour only. It is easier for the students to remember the position of the information on each page based on different colours”- Inf\_5.

Based on the comments, there were several characteristics of colour that were pointed out by informants such as a combination of striking and pastel colours, use more than one colour, corporate colour, background colour and consistency of colour background. There was disagreement among informants about the consistency of background colour. Based on the design principle, interface design should be consistent,

including the background colour. Besides, an informant suggested to use the corporate colour of the university as a sign of professionalism and to show that e-learning belongs to a respective university. However, other informants disagreed with the suggestion because of the limited choice of colour regarding the university’s corporate colour. This is because more colours on interface design would make e-learning more attractive and engaging to the students. Thus, contradict answers between informants would be asked in another empirical study that employs a questionnaire to find out the appropriate interface design of e-learning for generation Z students.

##### B. Graphic

The graphic is the main element in e-learning interface design. Informants agreed that graphic element was important to the e-learning users. Past researchers stated that users in Malaysia prefer the image of leaders or historical building in the interface design. However, generation Z students have different opinions about the types of images used in e-learning. Below are some comments from informants:

- “Use many images or change the image in e-learning frequently. Do not use the same image. If the image is always changing, interface design looks more attractive” – Inf\_6.
- “Images that can inspire students are more appropriate to display in e-learning. For example, the image of students studying or graduating as a motivation for students to complete their studies” – Inf\_5.
- “Image of successful university’s alumni with encouragement words as motivation to students” – Inf\_2.
- “No need to use a complicated image. Just use inspirational words as an image to inspire students who are using e-learning” – Inf\_7.
- “Encouragement words were spoken by high ranking university members. For example, the Vice Chancellor’s recent speech was full of motivational words, and this can inspire students to succeed. Therefore, a graphic that contains the image of people who deliver the speech along with motivational words should be placed in e-learning” – Inf\_4.
- “Use graphic that can motivate students not only related to learning but related to university and students” - Inf\_1.

Based on the comments, informants suggested changing the images frequently to avoid students from feeling boring with the same images and use motivational words as images to encourage students in their learning. Besides, majority informants suggested using graphics that are related to students. This is because the majority students are using e-learning so it should be student-oriented. However, an informant suggested using the image of a high ranking person in the university. So, this issue would be addressed through a large scale questionnaire. The informants also informed a few types of graphics that should be avoided used in e-learning

because it can distract the students from using the e-learning. The feedback obtained from the informants are:

- “Blurred graphic is not suitable in e-learning because it is hard to understand those graphics” – Inf\_3.
- “Graphics that have no motive or have an implicit meaning can make the students feel distracted while using e-learning. Not everyone can understand the message of these graphics is trying to convey. Therefore, graphic with clear meaning must be used in e-learning” – Inf\_5 & Inf\_6.
- “Do not use complicated graphic because students would misunderstand the graphic as for advertisement” – Inf\_2 & Inf\_7.

Furthermore, all informants agreed that institutional logo must be placed in e-learning because it symbolizes the identity of the institution. Below are some comments about the logo in the e-learning interface design:

- “The use of the university logo in e-learning is a must. This is because the logo symbolizes the institution” – Inf\_2, Inf\_5 & Inf\_6.
- “In Malaysia, the position of logo usually is at the top left. So, such a position is strategic, and it does not interfere with the student’s line of view when using e-learning” – Inf\_2 & Inf\_7.

### C. Typography

Past studies have stated that the selection of appropriate typography must be focussed in e-learning interface design. The result of the interview found that generation Z students are less interested in typography. They only agreed that the fonts used in e-learning should be easy to read and understand. Below are the comments about typography element in interface design:

- “Text symbolizes the content of e-learning page, whether it is formal or informal. If the e-learning page is informal, fancy fonts can be used while formal e-learning page use fonts that are easy to understand” – Inf\_6.
- “Use text that has been used prior in order to save time to process and understand the text” – Inf\_5.
- “Types and size of text should be consistent so that it is easy to read and text and background colour should be contra for it to be easy to see” – Inf\_7.
- “For important information, the text should be bold or highlight” – Inf\_2.
- “Or can use other colours to point out the important information” – Inf\_3.

Based on the comments from informants, readable and understandable text in e-learning interface design is important for generation Z students, and they want important information is highlighted. An informant pointed out about formal and informal text in e-learning, and this issue would be asked in the survey later.

### D. Layout

The layout is an interface element that covers the entire e-learning page, which involves the arrangement of each interface element such as images, logo, symbols, typography and others. Informant’s feedbacks about layout are as follows:

- “The horizontal layout makes the content to appear clear and less fibrous. More importantly, the layout should be organized and balanced to facilitate students to find information and so on” – Inf\_1.
- “The vertical layout looks simpler, less compact and neat. Besides, students do not have to scroll to the side and only scroll down” – Inf\_7.
- “Prefer simple layout for e-learning” – Inf\_6.
- “E-learning interface layout should be consistent on every page. For example, the main menu position should be on the same place on each page. This is because changing places would make it difficult for students to find the main menu” – Inf\_7.

Most informants suggested that e-learning layout should be simple, organized, balanced, clear and neat to facilitate the students. In addition, the informants also emphasized the importance of consistent layout in e-learning.

### E. Navigation

Navigation is important to help students in completing tasks and achieving their goals. Below are comments from informants:

- “Prefer the position of the main menu on the left side of the screen and use the button as a link would look neater. Suggest the use of menu similar to Gmail where only main options are displayed on the side, and other options are available by request” – Inf\_7.
- “The main menu on the left side of e-learning pages look suitable and organized if buttons are used. Besides, prefer if e-learning has various options to navigate instead of limited options” – Inf\_2.
- “Knowing the position when navigating e-learning is necessary. It is a convenience if you know where you are while exploring e-learning” – Inf\_5.

In this interview, most informants focussed on the position of the main menu and users' location while using e-learning. Informants agreed with the position of the main menu on the left page. The position of the main menu on the left page is common in Malaysia, and informants found it suits with their preferences. Informants also suggested using the button as links compared using textual links. They also stated that they want to know their location while navigating e-learning. Breadcrumb is often used as a link to any pages without having to follow any sequence. Besides, the informants also suggested that e-learning should provide many navigation options. This allowed students to explore e-learning without limited option. Previous studies also focused on errors [24][49]. However, informants did not provide deeper feedback about the errors in e-learning on how it should be handled. Therefore, further characteristics of error would be asked in the questionnaire.

*F. Additional Elements*

The informants also asked whether the five elements of the interface were enough for e-learning. All informants suggested audio, and video elements should be placed in e-learning. The latest technology is the main reason audio and video elements must be used in e-learning. Learning-based audio and video are more attractive and have a lot of benefits compared to text-based learning [68]. Besides, an informant suggested adding simple animation in e-learning interface design.

Focus group interview was conducted to verify the elements of e-learning interface design that fulfil the needs of generation Z students. Five elements have been presented to the informants whereby audio and video element was proposed in order to engage generation Z students towards e-learning. During the interview, all informants can freely give their opinions about each element and would be interpreted as cultural-based design and characteristics. Table V shows the list of significant elements of e-learning interface design together with design characteristics from theoretical and empirical data.

Table V shows the interface elements and design characteristics that were gathered from theoretical and empirical (focus group interview) study. Most of the results from the theoretical study are not preferable by generation Z students. For example, previous research stated that images of a leader were supposed to be used in the interface [18][24]. However, generation Z students want images that are related to students or learning in e-learning interface design. The empirical study has confirmed that developing e-learning interface based on Hofstede results is irrelevant because it student preferences changes. Besides, a new interface element is discovered in this study which is an audio/video element. For now, only audio/video elements added to the list of significant elements. The next empirical study, further design characteristics of audio/video would be discovered.

TABLE V. E-LEARNING INTERFACE ELEMENTS AND DESIGN CHARACTERISTICS FROM THEORETICAL AND FOCUS GROUP

| Interface Elements | Theoretical study   | Empirical study (Focus group interview)  |
|--------------------|---|--|
| Graphic            | <ul style="list-style-type: none"> <li>Using images of important people or monument that represent history.</li> </ul>                            | <ul style="list-style-type: none"> <li>Using an image of students such as study group or graduating ceremony.</li> </ul> |
|                    | <ul style="list-style-type: none"> <li>Use group images.</li> </ul>   |  |
| Colour             | <ul style="list-style-type: none"> <li>Use white or black colour as the background to focus on the name and image of the institution.</li> </ul>  | -  |
|                    | <ul style="list-style-type: none"> <li>Use official colours of the institution.</li> </ul>  | <ul style="list-style-type: none"> <li>Use various colours instead of the official colour of the institution.</li> </ul> |
|                    | <ul style="list-style-type: none"> <li>Using various different colours (colourful).</li> </ul>  | <ul style="list-style-type: none"> <li>Using various different colours (colourful). *</li> </ul>                         |
|                    | <ul style="list-style-type: none"> <li>Use colours for classification.</li> </ul>   | -  |
| Layout             | <ul style="list-style-type: none"> <li>Complex menu with various paths and need to scroll.</li> </ul>   | <ul style="list-style-type: none"> <li>Simple menu, organized and consistent.</li> </ul>                                 |
|                    | <ul style="list-style-type: none"> <li>Vertical page scrolling.</li> </ul>  | <ul style="list-style-type: none"> <li>Vertical/horizontal page layout.</li> </ul>                                       |
|                    | <ul style="list-style-type: none"> <li>Main menu on the left panel of the page.</li> </ul>  | <ul style="list-style-type: none"> <li>Main menu on the left panel of the page.*</li> </ul>                              |
| Navigation         | <ul style="list-style-type: none"> <li>Using a text link to the main page from any pages.</li> </ul>  | <ul style="list-style-type: none"> <li>Using button link to the main page from any pages.</li> </ul>                     |
|                    | <ul style="list-style-type: none"> <li>Provide feedback on user location.</li> </ul>  | <ul style="list-style-type: none"> <li>Provide feedback on user location.*</li> </ul>                                    |
|                    | <ul style="list-style-type: none"> <li>Non-linear navigation.</li> </ul>  | <ul style="list-style-type: none"> <li>Non-linear navigation. *</li> </ul>   |
|                    | <ul style="list-style-type: none"> <li>Long path while navigating.</li> </ul>   | <ul style="list-style-type: none"> <li>Simple path while navigating.</li> </ul>  |
| Typography         | <ul style="list-style-type: none"> <li>The main text colour is black while the blue and red colour to highlight important information.</li> </ul> | <ul style="list-style-type: none"> <li>Bold, highlight or using colours to highlight important information.</li> </ul>   |
|                    | <ul style="list-style-type: none"> <li>Using animation through text movement.</li> </ul>  | -  |
| Audio/Video        | -   | -  |

V. CONCLUSION AND FUTURE WORK

This paper discussed the interface elements and together with their corresponding design characteristics that are suitable for generation Z students. The elements and design characteristics were gathered qualitatively through theoretical and empirical study. The proposed interface elements such as colour, graphic, navigation, layout and typography have been discussed and verified by generation Z students that are considered as valid. This is because they have been agreed and supported by the previous study, as mentioned in Table I. Besides, the findings from this study show the role of culture in generation Z's perception towards design characteristics of e-learning which differ from the previous study.

However, these findings still need further research, especially in the design characteristics of e-learning interface, which would be validated in large scale survey. In future research, detailing the features of each interface element is needed to ensure the quality of e-learning interface design prolonged. Besides that a comparative study should be conducted in the future in order to compare interface design preferences of generation Z with other generations. In the meantime, these findings could be an eye-opener to other researchers that every culture or generation has their own preferences concerning interface design, and this can motivate students to be more engage towards e-learning.

#### ACKNOWLEDGMENT

This research is funded by the Fundamental Research Grant Scheme (FRGS) by the Ministry of Higher Education, Malaysia (FRGS/1/2016/ICT01/UKM/02/2) and Universiti Kebangsaan Malaysia Internal Funding (PP-FTSM-2020). The authors also thank the informants who participated in this study.

#### REFERENCES

- [1] M. T. Mulyadi and N. A. Mat Zin, "MMORPG Game Framework based on Learning Style for Learning Computer Networking," *Asia-Pacific J. Inf. Technol. Multimed.*, vol. 8, no. 1, pp. 63–77, 2019.
- [2] M. Mirabolghasemi, N. A. Iahad, and S. H. Choshaly, "Microblogging in Higher Education: A Comparative Study," *Asia-Pacific J. Inf. Technol. Multimed.*, vol. 6, no. 2, pp. 65–75, 2017.
- [3] F. J. Fernández-Cruz and M. J. Fernández-Díaz, "Generation z's teachers and their digital skills," *Comunicar*, vol. 24, no. 46, pp. 97–105, 2016.
- [4] E. J. Cilliers, "the Challenge of Teaching Generation Z," *PEOPLE Int. J. Soc. Sci.*, vol. 3, no. 1, pp. 188–198, 2017.
- [5] C. Seemiller and M. Grace, "Generation Z: Educating and Engaging the Next Generation of Students," in *About Campus*, vol. 22, no. 3, 2017, pp. 21–26.
- [6] IBM, "Uniquely Generation Z," 2017.
- [7] C. Chun, "Teaching Generation Z at the University of Hawaii," 2016.
- [8] M. R. Hanifa and H. B. Santoso, "Evaluation and Recommendations for the Instructional Design and User Interface Design of Coursera MOOC Platform," in *2019 International Conference on Advanced Computer Science and Information Systems (ICACSIS)*, 2019, pp. 417–424.
- [9] M. Dokhani, B. Majidi, and A. Movaghar, "Visually Enhanced E-learning Environments Using Deep Cross-Medium Matching," in *The 7th International and 13th Iranian Conference on E-learning and E-teaching*, 2019, pp. 0–4.
- [10] S. Liu, T. Liang, S. Shao, and J. U. N. Kong, "Evaluating Localized MOOCs: The Role of Culture on Interface Design and User Experience," *IEEE Access*, vol. 8, pp. 107927–107940, 2020.
- [11] Z. A. M. Drus, D. Singh, M. R. Mokhtar, and R. A. Rashid, "Review of Computerized Cognitive Behavioural Therapy Based on Culture Centered Design for Substance Abuse in Malaysia," *Asia-Pacific J. Inf. Technol. Multimed.*, vol. 7, no. 1, pp. 119–132, 2018.
- [12] S. Ali, M. A. Uppal, and S. R. Gulliver, "A conceptual framework highlighting e-learning implementation barriers," *Inf. Technol. People*, vol. 31, no. 1, pp. 156–180, 2018.
- [13] M. Harju, T. Leppanen, and I. Virtanen, "Interaction and Student Dropout in Massive Open Online Courses," 2018.
- [14] R. Rusdi, S. Fadzilah, and M. A. T. Noor, "Usability Guidelines for Elderly Website Interface," *Asia-Pacific J. Inf. Technol. Multimed.*, vol. 6, no. 2, pp. 109–122, 2017.
- [15] F. P. Sari and N. S. @ Ashaari, "Usefulness Model for the Redesign of Graduate's Student Management Information System," *Asia-Pacific J. Inf. Technol. Multimed.*, vol. 6, no. 1, pp. 100–114, 2017.
- [16] Q. Conley, Y. Earnshaw, and G. Mewatters, "Examining Course Layouts in Blackboard : Using Eye-Tracking to Evaluate Usability in a Learning Management System," *Int. J. Human-Computer Interact.*, vol. 36, no. 4, pp. 1–13, 2019.
- [17] J. Díaz, C. Rusu, and C. A. Collazos, "Experimental Validation of a Set of Cultural-Oriented Usability Heuristics: e-Commerce Websites Evaluation," *Comput. Stand. Interfaces*, vol. 50, 2016.
- [18] F. Mahmood, W. A. W. Adnan, N. L. M. Noor, F. M. Saman, and Z. A. Nasruddin, "User Perception Towards Cultural-Based E-Government Portal Design," in *APIT 2019*, 2019, pp. 5–9.
- [19] R. Heimgärtner, "IUID Method-Mix : Towards a Systematic Approach for Intercultural User Interface Design ( IUID )," *J. Comput. Commun.*, vol. 7, pp. 162–194, 2019.
- [20] H. B. Santoso, "Cultural Consideration for Designing E-Commerce Site Interface," *2018 1st Int. Conf. Comput. Appl. Inf. Secur.*, pp. 1–5, 2018.
- [21] N. Saidin, D. Singh, Z. A. M. Drus, and Z. A. Mohd Drus, "Culture Centered Design : Reviews on Cultural Factors Influencing Interface Design Elements," *Pertanika J. Sch. Res. Rev.*, vol. 3, no. 1, pp. 42–54, 2017.
- [22] C. Sang-Hun, "South Koreans Connect Through Search Engine," *The New York Times*, pp. 7–9, 2007.
- [23] K. Reinecke and A. Bernstein, "Knowing what a user likes: a design science approach to interfaces that automatically adapt to culture," vol. 37, no. 2, pp. 427–453, 2013.
- [24] F. Lachner, M.-A. Nguyen, and A. Butz, "Culturally sensitive user interface design: A case study with German and Vietnamese users," in *ACM International Conference Proceeding Series*, 2018, pp. 1–12.
- [25] S. Nizamani, S. Nizamani, K. Khoubati, S. Nizamani, S. Memon, and N. Basir, "Cultural preferences of Pakistan for the university website design," in *2018 International Conference on Information and Communication Technology for the Muslim World (ICT4M)*, 2018, no. 70, pp. 323–328.
- [26] L. Aljasmí and H. Alobaidy, "The Cultural Impact on User Interface Design : The Case of e-Government services of Kingdom of Bahrain," in *2018 International Conference on Innovation and Intelligence for Informatics, Computing, and Technologies (3ICT)*, 2018, pp. 1–10.
- [27] M. A. Muhanna, R. N. Amro, and A. Qusef, "Using a new set of heuristics in evaluating Arabic interfaces," *J. King Saud Univ. - Comput. Inf. Sci.*, 2018.
- [28] Z. Jano, H. Hussin, A. N. Abdullah, and C. K. Mee, "Website interactivity in Malaysian and australian universities," *Asian Soc. Sci.*, vol. 11, no. 17, pp. 14–21, 2015.
- [29] G. Vitols and Y. Vitols-hirata, "Impact of Culture Dimensions Model on Cross-Cultural Website Development," in *Proceedings of the 20th International Conference on Enterprise Information Systems (ICEIS 2018)*, 2018, pp. 540–546.
- [30] E. Purwanto, "Moderation Effects of Power Distance on the Relationship between Job Characteristics, Leadership Empowerment, Employee Participation and Job Satisfaction: A Conceptual Framework," *Acad. Strateg. Manag. J.*, vol. 17, no. 1, pp. 1–9, 2018.
- [31] N. Saidin, D. Singh, Z. A. M. Drus, R. Hidayat, Z. Akramin Mohd Drus, and R. Hidayat, "Cultural Marker Identification for Web Application Design Targeted for Malaysian Multicultural Users," *Int. J. Adv. Sci. Eng. Inf. Technol.*, vol. 6, no. 2016, pp. 959–965, 2016.
- [32] R. Cermak and Z. Smutny, "A Framework for Cultural Localization of Websites and for Improving Their Commercial Utilization," in *Global Observations of the Influence of Culture on Consumer Buying Behavior*, 2018, pp. 206–232.
- [33] J. Amantha, B. Muniandy, and W. A. J. W. Yahaya, "Exploring the effects of visual aesthetics in e-learning for engineering," *Knowl. Manag. E-Learning*, vol. 10, no. 3, pp. 250–264, 2018.
- [34] A. S. Adnan, M. Ali, and R. Ahmad, "The Utilization of Visual Elements on Interface Design of e-learning .," in *International Conference on Information Technology & Society*, 2015, no. June, pp. 273–279.
- [35] J. Reyna, "The importance of visual design and aesthetics in e-learning," *Australian Institute of Training and Development*, no. 40, pp. 28–32, 2013.

- [36] N. Kamaruddin and J. Hamilton, "Malaysian interactive courseware : An analysis of interface design development," *Int. J. Excell. e-Learning*, vol. 3, no. 1, pp. 1–17, 2010.
- [37] R. Alexander, D. Murray, and N. Thompson, "Cross-Cultural Web Design Guidelines," in *W4A '17 Proceedings of the 14th Web for All Conference on The Future of Accessible Work*, 2017, pp. 1–4.
- [38] N. Alomar, V. Wanick, and G. Wills, "The design of a hybrid cultural model for Arabic gamified systems," *Comput. Human Behav.*, vol. 64, no. 2016, pp. 472–485, 2016.
- [39] H. Ullah and A. Alhusseni, "Optimized Web Design in the Saudi Culture," in *Science and Information Conference 2015*, 2015, pp. 906–915.
- [40] H. Almakky, R. Sahandi, and J. Taylor, "The effect of culture on user interface design of social media - A case study on preferences of Saudi Arabians on the Arabic user interface of Facebook," *Int. J. Soc. Educ. Econ. Bus. Ind. Eng.*, vol. 9, no. 1, pp. 107–111, 2015.
- [41] A. Calabrese, G. Capece, M. Corbò, N. L. Ghiron, and M. M. Marucchi, "Cross-Cultural Strategies for Web Design," *World Acad. Sci. Eng. Technol. Int. J. Humanit. Soc. Sci.*, vol. 6, no. 11, pp. 78–83, 2012.
- [42] M. A. Khanum, S. Fatima, and M. A. Chaurasia, "Arabic Interface Analysis Based on Cultural Markers," *Int. J. Comput. Sci. Issues*, vol. 9, no. 1, p. 255, 2012.
- [43] M. Ali and H. Lee, "The impact of culture and social interaction on weblog design: a Malaysian case," *J. Enterp. Inf. Manag.*, vol. 24, no. 5, pp. 406–423, 2011.
- [44] I. Kim and J. Kuljis, "Applying Content Analysis to Web-based Content," *J. Comput. Inf. Technol.*, vol. 18, no. 4, pp. 369–375, 2010.
- [45] H. C. L. Hsieh, C.-H. Chen, and S. D. Hong, "Incorporating Culture in Website Design: A Comparison," *Springer-Verlag Berlin Heidelberg*, 2013, pp. 393–403, 2013.
- [46] R. George, K. Nesbitt, M. Donovan, and J. Maynard, "Evaluating Indigenous Design Features Using Cultural Dimensions," in *13th Australasian User Interface Conference (AUIC2012)*, 2012, pp. 49–58.
- [47] C. Chang and Y. Su, "Cross-cultural interface design and the classroom-learning environment in Taiwan," *Turkish Online J. Educ. Technol.*, vol. 11, no. 3, 2012.
- [48] A. Baharum, P. Turumugon, N. H. Mat Zain, C. P. Yee, S. Dullah, and F. A. Lahin, "Evaluation of Localization for E-Learning Website : a Preliminary Study," *Proc. 6th Int. Conf. Comput. Informatics*, no. 059, pp. 541–546, 2017.
- [49] Z. Ishak and A. Jaafar, "Cultural Dimensions of Malaysian Teenagers and Their Relationship with Interface Design," *J. Theor. Appl. Inf. Technol.*, vol. 90, no. 2, pp. 220–227, 2016.
- [50] R. Alhendawi and K. Meyer, "The Importance of Cultural Adaptation of B2C E-Services Design in Germany," *Int. J. Soc. Behav. Educ. Econ. Bus. Ind. Eng.*, vol. 9, no. 9, pp. 2749–2753, 2015.
- [51] R. Heimgärtner, "Using Converging Strategies to Reduce Divergence in Intercultural User Interface Design," *J. Comput. Commun.*, vol. 05, no. 04, pp. 84–115, 2017.
- [52] N. Khashman and A. Large, "Arabic Website Design : User Evaluation," pp. 424–431, 2013.
- [53] J. Selthofer, "Visual presentation and communication of Croatian academic websites," *Inf. Res.*, vol. 23, no. 1, 2018.
- [54] E. Callahan, "Cultural Similarities and Differences in the Design of University Web sites," *J. Comput. Commun.*, vol. 11, no. 1, pp. 239–273, 2005.
- [55] N. Kamaruddin, S. Sulaiman, M. Of, and C. V. Between, "A conceptual framework for effective learning engagement towards interface design of teaching aids within tertiary education," *J. Adv. Res. Soc. Sci. Humanit.*, vol. 2, no. 1, pp. 35–42, 2017.
- [56] C. M. N. Faisal, M. Gonzalez-rodriguez, D. Fernandez-lanvin, and J. De Andres-suarez, "Web Design Attributes in Building User Trust , Satisfaction , and Loyalty for a High Uncertainty Avoidance Culture," *IEEE Trans. Human-Machine Syst.*, pp. 1–13, 2016.
- [57] A. H. Alsswey, H. Al-Samarraie, F. A. El-qirem, A. I. Alzahrani, and O. Alfarraj, "Culture in the design of mHealth UI An effort to increase acceptance among culturally specific groups," *Electron. Libr.*, vol. 38, no. 2, pp. 257–272, 2020.
- [58] M. A. Hamid, "Analysis of visual presentation of cultural dimensions: Culture demonstrated by pictures on homepages of universities in Pakistan," *J. Mark. Commun.*, vol. 23, no. 6, pp. 592–613, 2016.
- [59] N. Kamaruddin, "Challenges of Malaysian Developers in Creating Good Interfaces for Interactive Courseware," *Turkish Online J. Educ. Technol.*, vol. 9, no. 1, pp. 37–42, 2010.
- [60] J. Reyna, "Developing quality e-learning sites: A designer approach," in *ASciliate 2009*, 2009, pp. 837–838.
- [61] M. Zainuddin, I. Md.Yasin, I. Arif, and A. B. A. Hamid, "Alternative Cross-Cultural Theories : Why Still Hofstede?," in *Proceeding of ISERD - Science Globe International Conference*, 2018, no. December.
- [62] S. Wartna and C. Risse, "Designing for Culture," *Dynamic Design Magazine*, 2019.
- [63] M. A. Nasrul, K. M. Nor, M. Masrom, and A. Syarief, "Website user interface characteristics for multiracial settings in Malaysia," *ICIMTR 2012 - 2012 Int. Conf. Innov. Manag. Technol. Res.*, no. April, pp. 252–257, 2012.
- [64] E. W. Gould, N. Zakaria, and S. A. M. Yusof, "Applying Culture to Website Design : A Comparison of Malaysian and US Websites," *IEEE*, pp. 161–171, 2000.
- [65] J. Zanariah, M. Shamsuri, M. Saad, and H. Janor, "Analyzing a Reflection of Uncertainty Avoidance Index Between Malaysian and Australian University Websites," *Asian J. Inf. Technol.*, vol. 16, no. 1, pp. 88–94, 2017.
- [66] J. Zanariah and S. M. Noor, "The Portrayal of Masculinity/Feminity Between Malaysian and Australian University Websites," no. June, pp. 1–14, 2015.
- [67] S. Nizamani, K.-R. Khoubati, S. Nizamani, S. Memon, S. Nizamani, and N. Basir, "Exploring the impact of cultural preferences in website design," in *2018 International Conference on Information and Communication Technology for the Muslim World (ICT4M)*, 2018, pp. 329–334.
- [68] M. F. Odhaib, "Does E-Learning Give a Better Result than Traditional Learning?," *Int. J. Comput. Sci. Mob. Comput.*, vol. 7, no. 9, pp. 29–36, 2018.
- [69] Ratna Z.R. Model reka bentuk antara muka permainan komputer berdasarkan nilai budaya. Dr. of Philosophy. Universiti Kebangsaan Malaysia.2016.

# Cotton Leaf Image Segmentation using Modified Factorization-Based Active Contour Method

Bhagya M Patil<sup>1</sup>

School of Computer Science  
REVA University, Bengaluru, India

Dr. Basavaraj Amarapur<sup>2</sup>

Department of Electrical and Electronics Engineering  
PDA College of Engineering, Kalaburgi, India

**Abstract**—Cotton plant is one of the most widely cultivated crop across worldwide. The leaf is one of the important parts which help in the food production. There are different cotton leaf diseases like *Alternaria* spot, foliar, bacterial blight, etc. which affects the agricultural yield. In order to detect the diseases, leaf region extraction becomes a significant task and to achieve this we use image processing techniques. Henceforth in this paper, a novel method used to extract the leaf region from a complex background. The proposed method is used for leaf extraction from complex background. The algorithm used in this method is modified factorization based active contour (MFACM) which helps in getting better output images. The database images used for research are acquired from the field using a digital camera. The proposed work is compared with existing active contour algorithms like Gradient Vector Flow (GVF), Adaptive Diffusion Flow (ADF), and Vector Flow Convolution (VFC). From the experiment, it can be observed that the proposed method is better than the other active contour methods in terms of computation time and the number of iterations. In addition to that segmented result is analyzed using specificity, sensitivity, precision which showed that our proposed method is better than the other methods.

**Keywords**—Cotton leaf; active contour; Gradient Vector Flow (GVF); Adaptive Diffusion Flow (ADF); Vector Flow Convolution (VFC); Modified factorization based active contour (MFACM)

## I. INTRODUCTION

Plants play an important source of food for human beings. If plants get affected than yield will also get affected. Therefore, many researchers are into this field of plant phenotyping [7] and there has been extensive research in this area. Plant phenotyping is a research area where the quantitative measurements of the structural and functional properties are performed. The experiment is carried out using some of the computer vision methods like analyzing the image through software. Because of this there has been research in this field for plant disease classification, leaf counting, observing the development and growth of the plant. In literature, there are many researchers who did lot of work for segmenting the leaf images [15,25,26,27] for example leaf segmentation without background and with complex background, automatic identification of plant species. In some cases, segmentation is performed using edge based and local based method also is used [16].

Kiruba raji [2], proposed a method for the herbal leaves segmentation from the complex background. In this chan-ve-se method is used for the segmentation of leaf or leaves from

background when compared to other techniques like k means clustering, local adaptive mean color, without affecting the color, textures. J Praveen kumar [3], introduced the new edge enhancement technique and graph-based method to extract the leaf region. Later it involves counting the number of leaves using circular Hough transform. Manual Grand et al. [4], a review of comparative study of 13 different segmentation methods. The Guided active contour approach was one among them and it has good segmentation results. Jones De Kylder [6], proposed a new active contour framework for segmentation. In this, the author proposed that whenever at the edge of the object contour is placed at that moment it will maximize probability. The internal and external probability distribution functions are learned from a ground truth training set. So, using this segmentation is performed and it gave an outstanding result. Shivalika Sharma [7], presented the review of various methods available for leaf segmentation. Xiaodong Tang [10], introduced leaf extraction from a complicated background using marker-controlled watershed algorithm for image of the individual channel gradient image of Hue, Saturation and Intensity separately. The applying of the watershed algorithm leads to the segmentation of the leaf region and author used solidity measures, to know how the performance of the method is.

Though there are various techniques used for the leaf segmentation [1][12][17][28] but with complex background very less research is carried out. Taking this into consideration, we proposed a leaf segmentation from complex background using modified factorization based active contour for texture segmentation. The structure of the paper is as follows: Section II is Literature survey: discussed about previous work done; Section III is Methodology: explains regarding proposed method; Section IV is Results and discussion: the comparative study results; and Section V presents Conclusion.

## II. LITERATURE SURVEY

In this, we will be focusing on the previous work done by researchers in active contours, or snakes. Active contour are curves which are generated by the computer. The curve will move towards the edges of an object based on the image energies. There are various applications of active contours like in the medical field like to outline tumor in the brain or it may be for the number recognition etc. But it has some disadvantages too like if the object boundary is too far than its difficult of the curve to evolve towards the object boundaries.

Active contours which there are two kinds, namely parametric and geometric.

1) The parametric deformable model –curves which are moved under the influence of internal and external forces [19,22].

2) Geometric deformable model – level set function is used to represent curves and surfaces implicitly in higher-dimensional scalar function [18,20,21].

There are various applications of these models in wide range of medical field [24]. But apart from that its contribution is in leaf segmentation also. And in this paper, both methods are taken into consideration for experiment.

#### A. Gradient Vector Flow

Xu and Prince [5], introduced a gradient vector flow method to overcome the problems of traditional snake. The traditional snake had a problem with the convergence towards the concave region and also with the limited capture range. So, in this author introduced a new external force which helps the snake to move towards the concave regions.

The equation for the traditional snake is given by

$$E_{int} = \int_0^1 E_{snake}(v(s))ds$$

$$= \int_0^1 E_{int}(v(s)) + E_{ext}(v(s)) \quad (1)$$

Where  $E_{int}$  and  $E_{ext}$  are the internal and external energies.  $E_{ext}$  is calculated from the image. In this author proposed a new external force which is given by.

$$E_{ext}(v(s)) = -|\nabla[G_\sigma(x, y) * I(x, y)]|^2 \quad (2)$$

Where  $G_\sigma(x, y)$  is a two dimensional Gaussian function with standard deviation  $\sigma$  and  $\nabla$  is the gradient operator. This energy helps in moving the snake towards the edges.

It will give better performance with respect to deep concavities of an image.

#### B. Adaptive Diffusion Flow

Yuwei Wu [8], proposed a method of adaptive diffusion flow which is the modified version of gradient vector flow. Gradient vector flow has advantages with respect to reaching deep concavities but it had a problem with the leakage of weak edges and narrow concavity. So, to improve the author replaced the smoothness energy term of GVF to harmonic hyper surface minimal functional and to achieve the contour to the deep and narrow concavities the infinite laplacian functional is incorporated.

Following are the equations for the adaptive diffusion flow which helps in overcoming the problems of GVF.

Case i): To preserve the weak edges following is the harmonic hypersurface functional defined as.

$$E(v) = \iint \frac{1}{p|\nabla f|} \cdot (\sqrt{1 + |G_\sigma \otimes \nabla v|^2})^{p|\nabla f|} \quad (3)$$

Where  $f$  is the edge of the image  $I$ ,  $p(\cdot)$  is a decreasing monotonously function which ranges from 1 to 2. And

$p(\nabla f) = 1 + \frac{1}{1 + |\nabla G_\sigma \otimes f(x)|}$  such that when  $|\nabla G_\sigma \otimes f(x)| \rightarrow 0$ ,  $p(\cdot) \rightarrow 2$ , this function behaves like isotropic diffusion within homogenous regions. When  $|\nabla G_\sigma \otimes f(x)| \rightarrow \infty$ ,  $p(\cdot) \rightarrow 1$  behaves like total variation model. This basically helps in preserving the weak edges of an image.

Case ii): The laplacian energy functional helps in contour convergence of narrow and deep concavities of an image. Following is the laplacian energy functional.

$$E_\infty v = \int |\nabla G_\sigma \otimes \nabla v|_{L^\infty(\Omega)} d\Omega \quad (4)$$

So, using both the cases the adaptive diffusion method is given by.

$$E(u) = \iint [g \cdot (-m \cdot |G_\sigma \otimes \nabla u|_{L^\infty}) + (1 - m) \cdot \frac{1}{p(|\nabla f|)} \cdot (\sqrt{1 + |G_\sigma \otimes \nabla u|^2})^{p|\nabla f|}] + h \cdot (|u - f_x|^2) d\Omega \quad (5)$$

Where  $g, h, m$  are the weighting functions.

#### C. Vector Flow Convolution

Bing Li [9], introduced new external force called as vector field convolution (VFC). Here the external force vector field kernel  $k(x, y)$  in which all the vectors point to the kernel origin.

$$k(x, y) = m(x, y)n(x, y)$$

The convolution of kernel  $k(x, y)$  with edge map  $f(x, y)$  which is generated from the image  $I(x, y)$  which is considered as VFC external force.

$$fvfc = f(x, y) * k(x, y)$$

Here edges are more prominent compare to homogeneous regions. The VFC will result in edge map which doesn't depend on the origin of the vector field kernel but it depends on the magnitude of the vector field kernel  $m(x, y)$ . Author proposed a two magnitude functions given by.

$$m_1(x, y) = (r + \epsilon)^{-\gamma} \quad (6)$$

$$m_2(x, y) = \exp\left(\frac{-r^2}{\zeta^2}\right) \quad (7)$$

Where  $\gamma$  and  $\zeta$  are positive parameters to control the decrease of vector field,  $\epsilon$  is a positive constant to prevent division by zero at the origin.  $m_1(x, y)$  which influence the FOI as increases as  $\gamma$  decreases.  $m_2(x, y)$  is a Gaussian shape function and  $\zeta$  represents the standard deviation.

#### D. Level Set Method

Stanley Osher [29], introduced the level set method formulation for the evolution of the curve. The curve is represented implicitly as a level set of a scalar function referred as level set function which is defined as the set of points that have same function value. Fig. 1 shows the curve at zero level set of a function,  $\Phi(x)=0$ .

#### E. Drawbacks

One of the drawbacks, reinitialization of the curve won't be possible whenever the level set function is far away from a signed distance. Practically whenever the time step is not selected small enough at that time evolving level set function



can depart greatly from its value as signed distance in a small number of iteration steps.

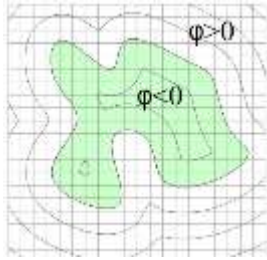


Fig. 1. Curve at Zero Level Set.

To maintain curve evolution stability and guaranteeing required results re-initialization has been broadly used as a numerical remedy. Later, various authors came with different methods which helped in overcoming the cost of reinitialization of curve during evolution.

Chumming Li [23], came with level set evolution without reinitialization cost which can be achieved by forcing the level set function to be close to signed distance function. The internal energy is used for deviation of the level set function from a signed distance function and external energy based on the desired feature like object boundary are used for the motion of zero level set curve.

### III. METHODOLOGY

#### A. Leaf Database

There are nearly 500 images for database. The images are captured under suitable conditions using Nikon digital camera of size 6000x4000. The images are resized nearly to 256 x 256 so that processing of images is faster. Fig. 2 shows the sample leaf images.

There are two different types of deformable models parametric and geometric. ADF, VFC, GVF are the methods for leaf segmentation and modified factorization based active contour.

#### B. Proposed Method

The proposed method is based on factorization-based active contour for texture segmentation [11][14] with modification. In geometric deformable the distance regularized active contour is used for the implementation [13]. The method is modified version of a factorization based active contour model for texture segmentation [11]. So, the process of extracting the leaf region shown in Fig. 3.

#### C. Modified Factorization based Active Contour Method (MFACM)

The images which are captured will undergo into preprocessing stage of flattening the field correction and using Gaussian pyramid.



Fig. 2. Sample Leaf Images.



Fig. 3. Process of Extracting the Leaf Region.

The Gaussian pyramid is a technique in image processing that breaks down an image into successively smaller groups of pixels, in repeated steps, for the purpose of blurring it. This is performed so that it is easier to detect the edges of an object.

The proposed method uses level set method introduced by chumming li [13], distance regularization is given by.

$$E_{regularization}(\phi) = \int_{\Omega} \frac{1}{2} (|\nabla\phi(x) - 1|)^2 dx \quad (7)$$

Where  $\nabla\phi$  denote the derivative of the level set method. The energy functional proposed by the author is defined as.

$$E(\phi, R) = \mu E_{data}(\phi, R) + \nu E_{regularization}(\phi) \quad (8)$$

where  $\mu$  and  $\nu$  are two positive constant which helps in balancing the total energy with respect to the corresponding terms in the equation. The equation is based on gradient based method.

The evolution process is defined as follows

$$\frac{\partial\phi}{\partial t} = -\frac{\partial E(\phi, R)}{\partial\phi} = -\delta_{\epsilon}(\phi)\mu(w_o - w_b) + \nu(\nabla^2\phi - \text{div}(\frac{\nabla\phi}{|\nabla\phi|})) \quad (9)$$

Where  $w_o$  and  $w_b$  are N dimension vectors.

Apart from level set method the proposed uses factorization-based texture segmentation introduced by Yaun et al. In this, spectral histograms are considered as texture features and a MxN matrix is calculated representing the feature matrix consisting of local window centered at each pixel. The matrix is denoted by Y which is given as.

$$Y = R\beta + \varepsilon$$

Where R is an MxL matrix gives you the columns which represent the features for the region to be segmented.  $\beta$  is the columns representing the weight vectors for every region and is of LxN matrix,  $\varepsilon$  is the additive noise.

The data energy can be represented as

$$E_{data}(\phi, R) = - \int_{\Omega} [H_{\varepsilon}(\phi)w_o(x, R) + (1 - H_{\varepsilon}(\phi))w_b(x, R)]dx \quad (9)$$

From the above equation 9 it can be observed that when the curve reaches the object boundary it data term will have minimum value at that point.

In this proposed method, we have included Gaussian pyramid along with factorization based active contour method.

#### Algorithm

- Step 1: Read the original image.
- Step 2: RGB to Gray scale image.
- Step 3: Histogram of RGB is matched with its gray scale image.
- Step 4: Flat field correction with sigma =1.
- Step 5: Gaussian low pass filter with hzise = [5 5] and sigma = 10.
- Step 6: Subtracting 'q' image from blurred image.
- Step 7: Reducing the flat filed corrected image using Gaussian pyramid.
- Step 8: Expanding the reduced image IR.
- Step 9: Fuse expanded image 'IR' and flat filed correction image 'O'.
- Step 10: Fuse Q and subtracted image q.
- Step 11: Flat field correction is applied to Q.
- Step 12: Image intensity values are adjusted with gamma = 0.01.
- Step 13: Flat field correction is applied with sigma = 5.
- Step 14: Required region should be selected from the original image.
- Step 15: Segmentation is performed using modified factorization active contour method.
- Step 16: Segmentation output.

#### IV. RESULTS AND QUANTITATIVE COMPARISONS

The images with complex background are captured from a various agriculture field of cotton plant leaf images. The images are captured from a digital camera of size 6000x4000.

The images were preprocessed using flat field correction and used Gaussian pyramid for the better visibility of the edges of the images. Later the processed image will be used for factorization based active contour. The results obtained using modified factorization based active contour will give better results when compare to other results.

The proposed method and parametric deformable model methods Gradient Vector Flow (GVF), Adaptive Diffusion Flow (ADF), Vector Field Convolution (VFC) are experimented with dataset created by Nikon digital camera of

image resolution 6000x4000. In this paper, Modified Factorization based Active Contour Method (MFACM) is compared with the existing parametric deformable models.

Fig. 4 and 5 shows the results obtained after applying the modified FAFM, GVF, ADF and VFC methods along with ground truth image. The performance of the methods are analyzed using time computation and segmentation measures like precision, accuracy, sensitivity, specificity. Table I gives the details regarding the time computation obtained from the various methods and we can observe that our proposed method takes less time when compare to other.

Fig. 6 shows the graph of time taken versus the number of iterations. Table II displays the output of different segmentation performance measures using ground truth images. From the table it can be seen that accuracy, precision, sensitivity, etc. are more compare to methods of parametric deformable models and it's been represented in the graph shown in Fig. 7.

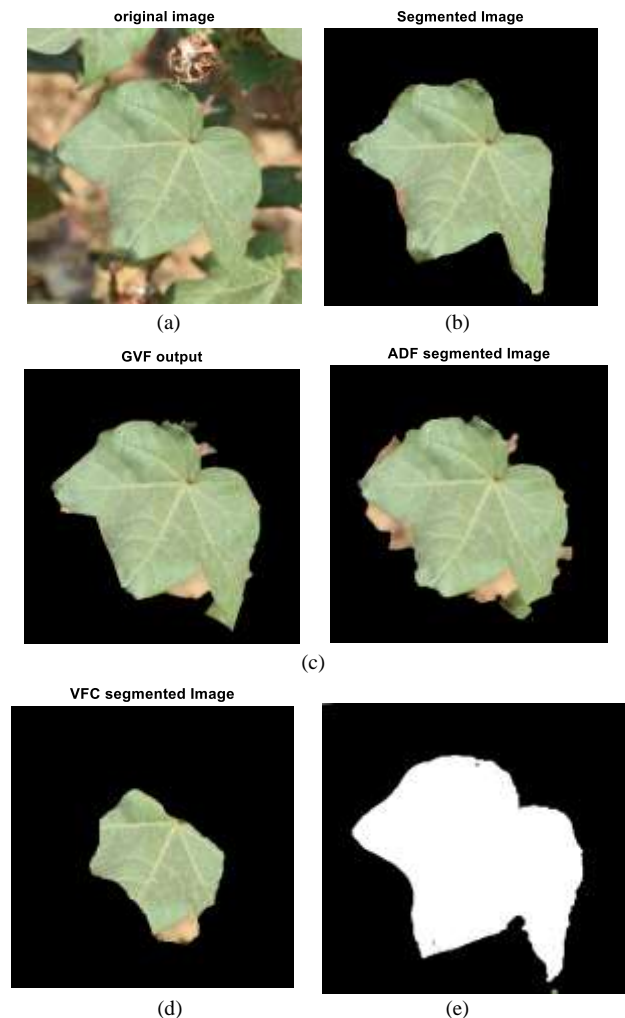


Fig. 4. (a) Input Image (b) MFACM Output (c) GVF Output (d) ADF Output (e) VFC Output (f) Ground Truth Image.

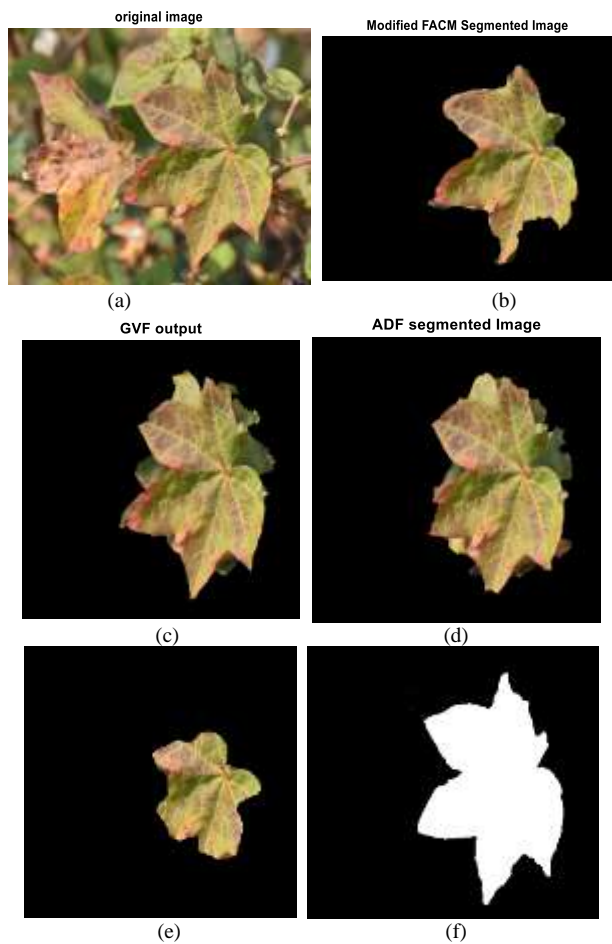


Fig. 5. (a) Input Image (b) Modified FACM Output (c) GVF Output (d) ADF Output (e) VFC Output (f) Ground Truth Image.

TABLE I. TIME SPENT FOR GVF, VFC, ADF AND MFACM

| Time spent | GVF   | VFC   | ADF   | MFACM |
|------------|-------|-------|-------|-------|
| Iterations | 60    | 60    | 60    | 60    |
| Leaf       | 24.68 | 29.14 | 29.41 | 12.8  |
| leaf1      | 23.73 | 40.67 | 32.54 | 11.06 |
| leaf2      | 22.43 | 41.71 | 27.1  | 9.26  |
| leaf3      | 23.63 | 37.01 | 30.91 | 11.59 |
| leaf4      | 17.96 | 36.23 | 17.35 | 12.74 |
| leaf5      | 21.68 | 29.32 | 22.92 | 11.54 |

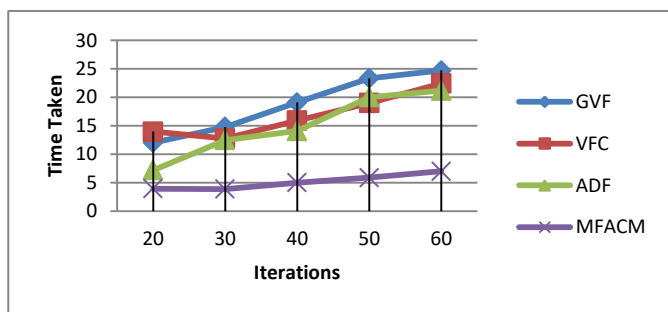


Fig. 6. Graph for Time Taken Versus Number of Iteration.

TABLE II. FALSE POSITIVE, FALSE NEGATIVE, TRUE POSITIVE, TRUE NEGATIVE, ACCURACY, SPECIFICITY, SENSITIVITY, PRECISION VALUES

|       | FP   | FN    | TP    | TN    | ACC    | SPE CI | SEN SI | Precisi on |
|-------|------|-------|-------|-------|--------|--------|--------|------------|
| GVF   | 2020 | 567   | 22115 | 40834 | 0.9605 | 0.9529 | 0.975  | 0.8864     |
| ADF   | 1838 | 1117  | 21565 | 41016 | 0.9549 | 0.9571 | 0.9508 | 0.9215     |
| VFC   | 704  | 11401 | 11363 | 42068 | 0.8153 | 0.9835 | 0.4992 | 0.9417     |
| MFACM | 682  | 1414  | 21268 | 42172 | 0.968  | 0.9841 | 0.9377 | 0.9689     |

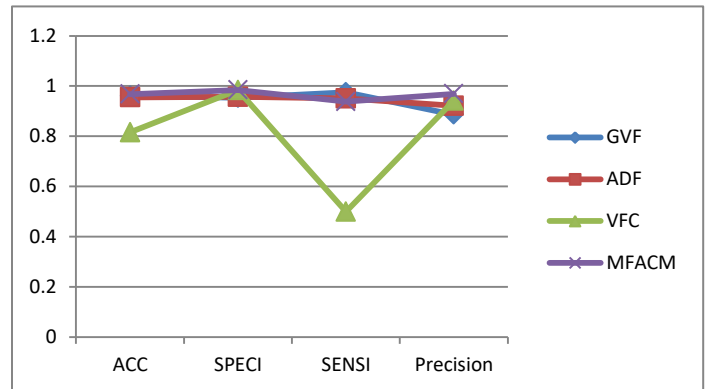


Fig. 7. Graph for Precision, Accuracy, Sensitivity, Specificity with Respect to MFACM, ADF, VFC, GVF.

## V. CONCLUSION

Leaf segmentation from complex background is performed using modified factorization based active contour for cotton leaf images. The segmented results are better than when compared to the parametric active contour methods. In this, active contour methods like GVF, VFC and ADF and level set method modified factorization based active contour method. From the results, we can observe that modified factorization based active contour method is good when compared to other methods in terms of time taken to perform segmentation. The precision, recall, sensitivity, etc. are calculated and the values are good when related to other methods. So, from this we can conclude that modified FACM is better than the other geometric and parametric active contour method.

## REFERENCES

- [1] Vijai Singh, A K Misra, "Detection of plant leaf diseases using image segmentation and soft computing techniques", Information processing in Agriculture, Volume 4, Issue 1, Pages 41-49, March 2017.
- [2] Kiruba raji, K K Thyagarajan, "An analysis of segmentation techniques to identify herbal leaves from complex background", International Conference on Intelligent Computing Communication and Convergence (ICCC 2014), Odisha, India.
- [3] J Praveen kumar, S. Domic, "Image based leaf segmentation and counting in rosette plants", Information processing in Agriculture, Volume 6, Issue 2, Pages 233-246, June 2019.
- [4] Manuel Grand-Brochier, Antoine Vacavant, Guillaume Cerutti, Camille Kurtz, Jonathan Weber, et al., "Tree leaves extraction in natural images: Comparative study of pre-processing tools and segmentation Methods". IEEE Transactions on Image Processing, Institute of Electrical and Electronics Engineers, 2015, 24 (5), pp.1549-1560.

- [5] Chenyang Xu and Jerry L. Prince, "Gradient vector flow: A external force for snakes", Proceedings of IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 1997.
- [6] Shivalika Sharma; Abhishek Gupta, "A review for the automatic methods of plant's leaf image segmentation", International journal of Intelligence and sustainable computing, Vol No 1, 2020.
- [7] Hanno Scharr, Massimo Minervini, Andrew P French, Christian Klukas, David M Kramer, Xiaoming Liu, Imanol Luengo, Jean Michel Pape, Gerrit Polder, Danijela Vukadinovic, Xi Yin, Sotirios A Tsafaris, "Leaf segmentation in plant phenotyping: a collation study", Machine vision applications, December 2015.
- [8] Yuwei Wu, Yuanquan Wang, Yunde Jia, "Adaptive diffusion flow active contours for image segmentation", computer vision and understanding, 2010.
- [9] Bing Li, Scott T Acton, "Vector field convolution for Image segmentation using snakes", 2006 International Conference on Image Processing, IEEE.
- [10] Xiaodong Tang, Manhua Liu, Hui Zhao, Wei Tao, "Leaf Extraction from Complicated Background", 2nd International Congress on Image and Signal Processing, IEEE, 2009.
- [11] Mingqi Gao, Hengxin Chen, Shenhai Zheng, Bin Fang, "A factorization based active contour model for texture segmentation", International conference on Image Processing, IEEE, 2016.
- [12] Zhang Jian-hua, KONG Fan tao, Wu Jian zhai, HAN Shu-qing, ZHAI Zhi-fen, "Automatic image segmentation method for cotton leaves with disease under natural environment", Journal of Integrative agriculture, 2018.
- [13] Chunming Li, Chenyang Xu, Changfeng Gui, and Martin D., "Level Set Evolution without Re-Initialization: A New Variational Formulation" Fox. IEEE CVPR, 2005.
- [14] Jiange Yaun, Deliang Wang, Anil M Chariyadat, "Factorization based Texture Segmentation", IEEE Transaction on Image Processing", vol 24, Issue 11, 2015.
- [15] Peng, Wu, Li Wenlin, and Song Wenlong. "Segmentation of Leaf images Based on Active Contour." International Journal of-u and-e Service, Science and Technology 8.6 (2015): 53-70.
- [16] Srikanth, Manassanan. "Active contours segmentation with edge based and local region based." Proceedings of the 21st International Conference on Pattern Recognition (ICPR2012). IEEE, 2012.
- [17] Zhang, Shanwen, et al. "Plant diseased leaf segmentation and recognition by fusion of superpixel, K-means and PHOG." *Optik* 157 (2018): 866-872.
- [18] V. Caselles, F. Catte, T. Coll, and F. Dibos, "A geometric model for active contours in image processing", *Numer.Math.*, vol. 66, pp. 1-31, 1993.
- [19] Cerutti G., Tougne L., Vacavant A., Coquin D., "A Parametric Active Polygon for Leaf Segmentation and Shape Estimation, Advances in Visual Computing. ISVC 2011.
- [20] J. A. Sethian, *Level Set Methods and Fast Marching Methods: Evolving Interfaces in Computational Geometry, Fluid Mechanics, Computer Vision, and Material Science*. Cambridge, UK: Cambridge University Press, 2nd ed., 1999.
- [21] I. Cohen, L. D. Cohen, and N. Ayache, "Using deformable surfaces to segment 3-D images and infer differential structures," *CVGIP: Imag. Under.*, vol. 56, no. 2, pp. 242-263, 1992.
- [22] Jonas De Vylder, Daniel Ochoa, Wilfried Philips, Laury Chaerle, and Dominique Van Der Straeten, "Leaf segmentation and tracking using probabilistic parametric active contour", vol 6930, 2011.
- [23] Xu, Chenyang, and Jerry L. Prince. "Gradient vector flow." *Computer Vision: A Reference Guide* (2020): 1-8.
- [24] Eltanboly, Ahmed, et al. "Level sets-based image segmentation approach using statistical shape priors." *Applied Mathematics and Computation* 340 (2019): 164-179.
- [25] Kumar, J. Praveen, and S. Domnic. "Image based leaf segmentation and counting in rosette plants." *Information Processing In Agriculture* 6.2 (2019): 233-246.
- [26] Wang, Ping, et al. "An maize leaf segmentation algorithm based on image repairing technology." *Computers and Electronics in Agriculture* 172 (2020): 105349.
- [27] Zhang, Shanwen, Zhuhong You, and Xiaowei Wu. "Plant disease leaf image segmentation based on superpixel clustering and EM algorithm." *Neural Computing and Applications* 31.2 (2019): 1225-1232.
- [28] Tian, Kai, et al. "Segmentation of tomato leaf images based on adaptive clustering number of K-means algorithm." *Computers and Electronics in Agriculture* 165 (2019): 104962.
- [29] S. Osher and J. A. Sethian, "Fronts propagating with curvature-dependent speed: algorithms based on Hamilton-Jacobi formulations", *J. Computational physics*, vol 79, pp. 12-49, 1988.

# Trading Saudi Stock Market Shares using Multivariate Recurrent Neural Network with a Long Short-term Memory Layer

Fahd A. Alturki<sup>1</sup>, Abdullah M. Aldughaiem<sup>2</sup>

Electrical Engineering Department  
King Saud University, Riyadh  
Saudi Arabia

**Abstract**—This study tests the Saudi stock market weak form using the weak form of an efficient market hypothesis and proposes a recurrent neural network (RNN) to produce a trading signal. To predict the next-day trading signal of several shares in the Saudi stock market, we designed the RNN with a long short-term memory architecture. The network input comprises several time series features that contribute to the classification process. The proposed RNN output is fed to a trading agent that buys or sells shares based on the share current value, current available balance, and the current number of shares owned. To evaluate the proposed neural network, we used the historical oil price data of Brent crude oil in combination with other stock features (e.g., previous day) opening and closing price of the evaluated share). The results indicate that oil price variations affect the Saudi stock market. Furthermore, with 55% accuracy, the proposed RNN model produces the next-day trading signal. For the same period, the proposed RNN trading method achieves an investment gain of 23%, whereas the buy-and-hold method obtained 1.2%.

**Keywords**—Time series; neural network; long short-term memory; stock price; Tadawul

## I. INTRODUCTION

Of all the presented works for forecasting stock markets, only very few have targeted the Saudi stock market. In this study, we presented a recurrent neural network (RNN) that utilizes the long short-term memory (LSTM) architecture for a multivariate time series prediction to generate a trading signal (buy, sell, or do nothing) for several Saudi stock indices that will be used in combination of a trading algorithm to buy and sell shares based on three factors: share current value, current available balance, and a current number of shares owned.

Neural networks have gained much attention in recent years, especially in stock market prediction. The nature of the randomness accosted with the stock market makes it hard to achieve high confidence in predicting the index price using normal statistical methods. By using neural networks with several futures, we can achieve a high prediction value. To study the effect of past historical prices on future prices and to develop a trading agent using neural networks, we tested the Saudi stock market for the weak form efficiency. In producing a trading signal, the developed neural network is an RNN with an LSTM architecture.

The remainder of the paper is organized as follows. Section II gives a literature review on the works undertaken to predict and forecast the stock market price. Section III tests the weak form of the Saudi stock market efficiency. (The test is useful for understanding the effect of the historical data of a share on future values.) Discussion on the proposed method and a brief neural network introduction is presented in Section IV. Section V evaluates the proposed method and compares it to a known trading method. Finally, we give the conclusions of this study in Section VI.

## II. LITERATURE REVIEW

Recently, the stock market prediction has been a hot topic in the research field. To predict stock prices, many researchers have developed methods, but only a few have developed a trading strategy. Some of the reviews of the developed methods are published [1, 2]. For example, Shah et al. classified stock prediction methods into four categories: statistical methods, pattern recognition, machine learning, and sentiment analysis.

The autoregressive integrated moving average (ARIMA) model, which is one of the well-known statistical methods, uses a class of models to model the time series based on historical values. The model is fitted to the historical values of a stock price in predicting (forecasting) the stock's future price. The model consists of three parts: (1) an autoregressive (AR) model, in which the forecasted value is a linear combination of past lagged values; (2) a moving average (MA) model that forecasts the future value using the past forecast errors; and (3) the difference operation of past and future values. The model is denoted by  $ARIMA(p, d, q)$ , where  $p$  is the order (number of time lags) of the AR model,  $d$  is the degree of differencing, and  $q$  is the order of the MA model.

Pattern recognition is closely related to machine learning but with a different implementation. Here, we focus on the methods of finding patterns in the stock's historical values. Then, by using computer algorithms, we predict future values using these patterns. Previous studies show an example of a pattern: the stock uptrend [3] and the open high–low close price candlestick charts [4].

Machine learning prediction uses historical data and the desired output as the training sets to build a mathematical model through an iterative process until an objective function is optimized. Previous studies have shown the usage of classification and regression as examples of machine learning in trading methods and the closing price of stock [5, 6, 7].

In sentiment analysis, it uses text information, such as news articles or social media feeds on stock markets. In predicting stock trends based on the feed provided, the analysis employs machine learning algorithms [8].

Idress et al. [9] built an ARIMA model to predict the Indian stock market, in which they found a deviation on a 5% mean percentage error.

Meanwhile, to predict the Saudi stock prices, Olatunji et al. [10] proposed an artificial neural network (ANN) model, applying on three major stock indices: Alrajhi bank, Saudi Telecom Company, and Saudi Basic Industries Corporation SABIC stocks. They only used the previous-day closing price as the model input. Moreover, the proposed model was used as an investment adviser, and it achieved a low root mean squared error (RMSE) of 1.8174 and a mean absolute percentage error of 1.6476.

Also, Jarrah and Salim [11] proposed an RNN and a discrete wavelet transform (DWT) to predict the Saudi stock price trends. The model consisted of two stages. The first stage uses DWT to break the stock price into both frequency and time domains to filter the noise associated with the signals, and the second stage is an RNN that performs the prediction. The model was tested to predict the next-seven days closing price of the Saudi stock. The prediction result was then compared with that obtained by a prediction process performed using the ARIMA model. Consequently, the proposed model (DWT + RNN) achieved an RMSE of 0.0522 when the RNN model used four batches and four neurons.

Alotaibi et al. [12] also used an ANN model to predict the Saudi stock market. Their ANN model consisted of three layers: input, hidden, and output layers. The input layer contained the historical close and open prices of the Saudi stock market and the historical close and open prices of oil. Bayesian regularization backpropagation was used for network training from 2003 to 2012. The test set training spanned from 2013 to the end of 2015.

Hua et al. [13] gave an introduction to deep learning with LSTM for time series prediction and proposed random connectivity for LSTM to overcome the computation cost.

Tilakaratne et al. [5] developed a neural network for predicting the trading signals of the Australian All Ordinary Index. Then, they compared an ANN to a probabilistic neural network (PNN), in which they found that the ANN outperformed the PNN.

On the basis of the previous studies mentioned above, many developed methods use historical information from the share itself without the combination of other factors (e.g., oil prices). These methods targeted different markets other than the Saudi stock market.

### III. WEAK FORM OF EFFICIENT MARKET TEST

The weak form of an efficient market hypothesis states that the future prices of a stock market with a weak efficiency cannot be predicted using historical information, such as trading volume, closing price, and earnings. It means that one cannot predict future values using the available information. Fama [14] divided the efficient market hypothesis into three: weak, semi-strong, and strong hypotheses.

Previous studies tested the Saudi stock market efficiency in its weak form and concluded the same; however, the presented studies are not up to date [15, 16].

To prove that the stock price under test can be predicted using historical values, we will be testing the Saudi stock indices used to evaluate the proposed RNN for the weak-form efficiency hypothesis. The weak form of the market efficiency for individual stocks is tested for randomness. If the stock does not follow a random walk, the hypothesis fails. The stock index can be predicted using historical data.

Several statistics tests are known for use in testing data randomness. Here, we used the Kolmogorov–Smirnov test (K–S test). The null hypotheses in the K–S test are that the data (stock returns) under the test follow a random walk, and the future value cannot be predicted. The alternative hypotheses are that the data under test are not random and that the data can be predicted using historical values.

Here, we used Alrajhi, Alinma, and SABIC stocks. The historical values are dated from January 2010 to the end of March 2020. The stocks' closing price was converted to the stock returns, as shown in Eq. (1), where  $R$  is the logarithmic stock return;  $l(i)$  is the day  $i$  closing price; and  $l(i - 1)$  is the previous closing price of the day  $i$ :

$$R(i) = \log\left(\frac{l(i)}{l(i-1)}\right) \quad (1)$$

#### A. Kolmogorov–Smirnov Test

The K–S test is a nonparametric test for data randomness. The null hypotheses of the test assume that the cumulative distribution function (CDF) of the data under test is equal to the hypothesized CDF. The CDF of the data was computed herein and compared with the hypothesized CDF using Eq. (2), where  $D_n$  is the maximum amount of the hypothesized CDF ( $F_n(x)$ ) exceeding the calculated CDF ( $G_n(x)$ ). When both CDFs are equal to some factors, the data are random, and the test fails to reject the null hypothesis that the test statistics converge to zero as  $n$  goes to infinity. Detailed mathematical background on the K–S test is provided in [17].

$$D_n = \max_x |F_n(x) - G_n(x)| \quad (2)$$

#### B. Market Weak form Test Results

We performed the test on the three stocks used to evaluate the proposed RNN. Table I shows the result of the K–S test performed with a significance level of 0.05. (The  $p$ -value is the probability value of the test.) Smaller values (typically  $<0.05$ ) indicate a strong rejection of the null hypothesis. The test statistic is a random variable calculated from the data under the test used in determining the null hypothesis rejection, whereas the  $z$ -value is the critical value. The K–S

test rejected the null hypotheses by comparing the  $p$ -value with the significance level. The null hypothesis is rejected if the  $p$ -value is less than the significance level (i.e., the data under test are not random).

Based on the test performed, Alinma, Alrajhi, and SABIC stock returns did not follow a random walk and were not independent of past values. This proved that the proposed stock prediction method and the trading agent could facilitate historical values to predict trading signals.

TABLE I. KOLMOGOROV-SMIRNOV TEST RESULTS

| Stock   | Hypothesis test result           | $p$ -value | Test statistic | $z$ -value |
|---------|----------------------------------|------------|----------------|------------|
| Alinma  | The null hypothesis is rejected. | 0.00       | 0.1191         | 0.0268     |
| Alrajhi | The null hypothesis is rejected. | 0.00       | 0.0904         | 0.0268     |
| SABIC   | The null hypothesis is rejected. | 0.00       | 0.1143         | 0.0268     |

#### IV. METHODOLOGY

Neural networks are a set of algorithms used to recognize underlying relationships in data sets. The process of a neural network is similar to the operation of a human brain. Here, we used an RNN with an LSTM architecture to produce a trading signal.

The input to the neural network is called a feature, which is a measurable characteristic of the observed data or a characteristic with an indirect effect on it. Accordingly, this section provides a brief introduction to neural networks. The introduction aims to familiarize the reader with the basics of neural networks and provide them the ability to understand some concepts. A detailed background regarding this matter is reported in [18].

##### A. RNN

RNNs are a class of neural networks best used in sequenced data sets, such as time series. An RNN has a one-to-one connection between its internal layers and the exact position in the time series [18]. An RNN can simulate any algorithm given sufficient data. These networks are based on the works by Rumelhart et al. [19], who described a new method for teaching a network through backpropagation. Unlike feedforward neural networks, RNNs have an advantage in using their internal memory to process a sequence of data, such as stock markets. Moreover, the network input (e.g., oil prices and index price) in RNNs are interrelated. On the contrary, an RNN suffers from exploding problems and gradient vanishing. Gradient vanishing is a term associated with neural network training, and a gradient is a vector of the calculated error during the network training process. The gradient is used to update the network weights to achieve a small error, such as an error in predicted stock value when compared with the actual value. The gradients in an RNN accumulate during the update process, which causes it to explode (i.e., it becomes large and goes to infinity).

Fig. 1 shows the basic building block of an RNN. The input to the block is a vectored time series  $x_t$ . In our case, we

used the stock price and associated features.  $h_t$  is the output from the block to be fed to the subsequent titration at time  $t + 1$ .  $h_{t-1}$  is the output from the previous block. Both  $h_t$  and  $h_{t-1}$  are called the hidden layer vectors.  $w_h$  and  $w_x$  are the weight vectors for the hidden connection and the input vector, respectively. The weight vectors are chosen by network training, which is achieved by comparing the output (predicted) with the actual value and adjusting the weight vectors to achieve the smallest error.  $F$  is an activation function within the block. Activation functions are mathematical equations that determine the block output based on preset conditions. The most important activation function is the  $\tanh$  function.  $b_t$  is a bias added to the block input. Equation (3) shows the math behind RNNs.

$$o_t = h_t = F(w_h h_{t-1} + w_x x_t + b_t) \tag{3}$$

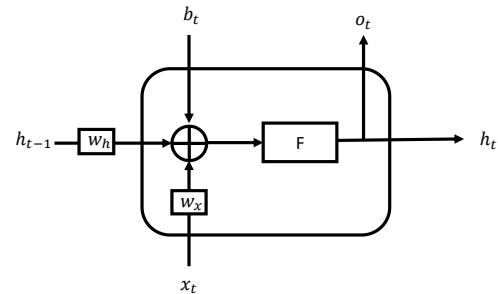


Fig. 1. The Basic Building Block of an RNN.

##### B. LSTM

Proposed by Hochreiter and Schmidhuber [20], LSTM is a type of RNN architecture used to solve the exploding and vanishing gradient problem that occurs in a normal RNN. The constant error carousel (CEC) LSTM was used to overcome the problems caused by the error back flow. The CEC controls the error flow by units, called gates, which are implemented in the memory block of the LMTS. The gates are categorized into the input gate, output gate, and forget gate, in which each gate has a function to achieve. The input gate controls the flow of the new sequence value. The output gate controls the usage of the value inside the cell using the activation function of the LSTM. The forget gate controls how long a value remains inside the memory cell.

Fig. 2 shows the building block of an LSTM unit, where  $C_t$  is the cell state,  $x_t$  is a vector input to the cell,  $f_t$  is the output from the sigmoid function that represents what cell state can be passed from adjacent cells, and  $i_t$  is the output from the sigmoid function that represents the output from the  $\tanh$  function of the input gate to the cell. This updates the cell state with new values.  $O_t$  is multiplied by  $\tanh$  of the cell state to choose what part to output to the adjacent cell.

Fig. 3 shows three hidden units for a vector input in an LSTM network. This number can be more than three, depending on the design. Equations (4)–(8) are the compact forms of the forward pass of an LSTM unit that contains a forget gate developed in [21]. In the equations,  $W$ ,  $U$ , and  $b$  denote the weights and biases determined by network training. Each layer produces a single output, called  $h_t$ , which is connected to a neuron at the final layer. The function of the neuron is to multiply each input by weight and sum them up to

produce an output  $\hat{y}_t$  with length n, where n is the number of classifications produced (Fig. 4).

$$f_t = \sigma_g(W_f x_t + U_f h_{t-1} + b_f) \quad (4)$$

$$i_t = \sigma_g(W_i x_t + U_i h_{t-1} + b_i) \quad (5)$$

$$g_t = \sigma_c(W_g x_t + U_g h_{t-1} + b_g) \quad (6)$$

$$o_t = \sigma_g(W_o x_t + U_o h_{t-1} + b_o) \quad (7)$$

$$\hat{f}_t = \sigma_g(W_f x_t + U_f h_{t-1} + b_f) \quad (8)$$

$$h_t = o_t \circ o_h \quad (9)$$

The output  $\hat{y}_t$  is connected to a softmax layer, which functions to convert the input vector  $\hat{y}_t$  of n elements to a normalized probability distribution with n probabilities. The element with the highest probability is the network output. The produced classifications are two training signals: buy and sell. An in-depth discussion on pattern recognition and classification is shown in [22].

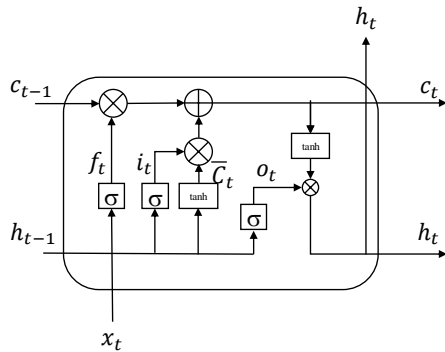


Fig. 2. LSTM basic Building Cell Called a Neuron or a Hidden Unit.

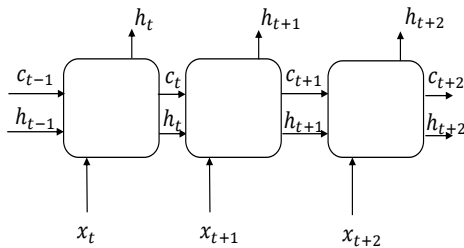


Fig. 3. The Network of the LSTM Units Known as Hidden Layers.

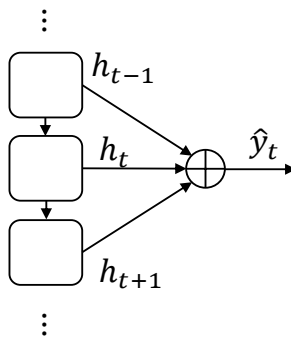


Fig. 4. Final Network Stage. The Output from each Cell is Added to Produce the Prediction.

### C. Trading

The proposed design was constructed using LSTM layers connected in series. The RNN input comprised a set of time series data representing the features associated with the stock and oil closing price. The network setup consisted of the training method, the number of hidden elements (LSTM units), and the number of training titrations. Fig. 5 shows a history of three stock prices in Saudi Riyals that was used in this study. The data will be divided into two sets. The first set will be used to train the classifier, and the other data will be used to evaluate the proposed classifier. Fig. 6 shows the history of the oil prices that will be used as an input to the proposed network. Table II lists the options used in constricting the network.

1) *Input features:* Table III lists the features used for the buy and sell classification network. Several methods can be used for feature selection. However, in this study, we used a trial-and-error method to find the best feature combination because some feature selection methods fail when chart technical indicators are used in the stock price.



Fig. 5. Historical Data of Three Stocks from March 21, 2012, to April 24, 2020.

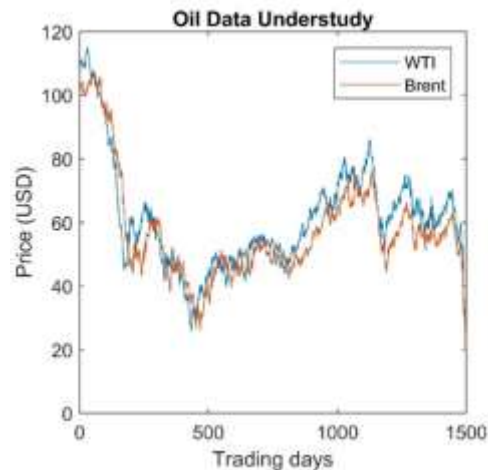


Fig. 6. Historical Data on Oil Prices in USD that are used in the Study. The Data are from March 21, 2012, to April 24, 2020.



TABLE II. LSTM NETWORK SETTINGS

| Option             | Description   | Value |
|--------------------|---|-------|
| Solver             | Training algorithm  | ADAM  |
| Epoch              | Number of full training data passes   | 750   |
| Hidden layers      | Number of LTMS cell per time series   | 200   |
| Gradient threshold | The gradient is clipped to the threshold if the gradient of the error passes the value                            | 1     |
| Initial learn rate | Specifying the rate of learning higher values will cause the learning to be faster, but could diverge the network | 125   |

TABLE III. FEATURE DESCRIPTION

| Features                                   | Description   |
|--|---|
| Stock closing price                        | The previous-day closing price of the stock   |
| WTI daily price                            | West Texas Intermediate oil price   |
| Brent daily price                          | Brent crude oil price   |
| No. of trades                              | Number of trades placed on a stock for the previous day   |
| Open price                                 | The opening price of the same day   |
| Highest price                              | The highest price of the previous day   |
| Lowest price                               | The lowest price of the previous day  |
| Month number                               | Current month in numerical form   |
| Number of days                             | Since the last trading session<br>Until the upcoming trading session  |
| Relative strength index                    | Relative strength index for 7 days<br>Relative strength index for 21 days   |
| Accumulation/distribution (A/D) oscillator | Momentum indicator for detecting the changes in the A/D line by measuring the momentum of the first signal of change of trend |
| Moving average convergence/divergence      | A trend-following momentum indicator that shows the relationship between two moving averages of a security's price            |
| Stochastic oscillator                      | An indicator comparing a closing price of a stock to a range of its prices over a certain period                              |
| Logarithmic return                         | The logarithm of the closing price divided by the previous closing price shown in Eq. (3)                                     |

2) *Network training*: To obtain the required gains and biases in the hidden network layers, we must train the neural network. A data set must be prepared to perform the training and evaluation processes of the RNN. The required data were divided into two sets: a training set and an evaluation set. The data set comprised the historical values of the proposed futures from March 21, 2012, to April 24, 2020, and the required response (trading signal) of that interval. The trading agent responses were obtained from the stock returns, in which a buy signal was generated from a positive return, and a sell signal was entreated from a zero or negative return. The data were normalized using Eq. (9).

$$x' = \frac{x - \min(x)}{\max(x) - \min(x)} \quad (9)$$

Each training run computes the generated responses with the required ones. An error is produced if the response is different, and the weights are updated in each training iteration. Adaptive moment estimation (ADAM), developed by Diederik Kingma and Jimmy Ba [23], was used as a solver to optimize the weights and biases of the neural network. The following lists the process undertaken to train the LSTM network.

- Initialize the LSTM network weights and biases randomly.
- Input the historical data to the network as a normalized time series.
- Compare the trading signal output with the required signal (buy and sell signal).
- Update the weights and biases using the ADAM solver and the computed error.
- Repeat the training process until the classification accuracy is higher than that in the previous run or stop when the required number of iterations has been satisfied.
- The evaluation data set was used to test the network after network training. This process is called the classification process.

3) *Trading agent*: The output of the neural network classification is connected to a trading agent. The presented trading agent strategy involves buying or selling a pre-defined number of shares in a trading session based on the number of shares and money currently owned. Fig. 7 depicts the trading process. The agent relies on the initial investment budget and the required shares to be bought and sold per trading session. These values are fixed in the current version of the trading agent.

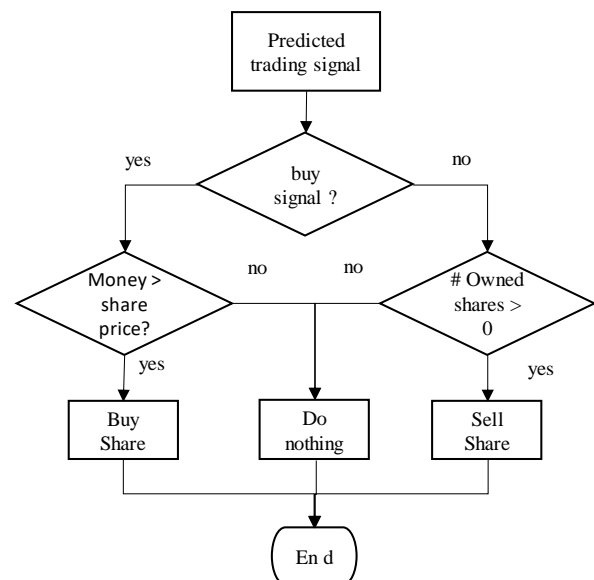


Fig. 7. Trading Agent Flow Chart.

V. RESULTS

The proposed neural network and trading agent were evaluated using three stock shares from the Saudi stock market (i.e., Alinma Bank, Alrajhi Bank, and SABIC). The evaluation data set comprised of historical values from June 2018 to August 2019. The performance of the proposed agent was compared with that of the buy-and-hold trading strategy. Table IV shows the accuracy of the trading signal, trading agent initial values, and investment gain. The trading gain was affected by the initial values used, which were optimized to achieve the highest gain.

Fig. 8 and 9 denote the output of the trading agent for the Alinma and Alrajhi stocks, respectively. The trading agent was effective for both the Alinma and Alrajhi shares, as shown by the output. The agent bought shares in an upward trend and sold them at the local maximum in several instances.

TABLE IV. TRADING OUTPUT RESULTS

| Stock   | Trading accuracy |       |         | Investment gain | Buy-and-hold gain |
|---------|------------------|-------|---------|-----------------|-------------------|
|         | Buy              | Sell  | Overall |                 |                   |
| Alinma  | 53.3%            | 50.3% | 57.3    | 28.24%          | 6.9%              |
| Alrajhi | 51.9%            | 60.4% | 57.3%   | 18.087 %        | 10.1 %            |
| SABIC   | 47.8%            | 61.4% | 57.3%   | 0.01%           | -23%              |

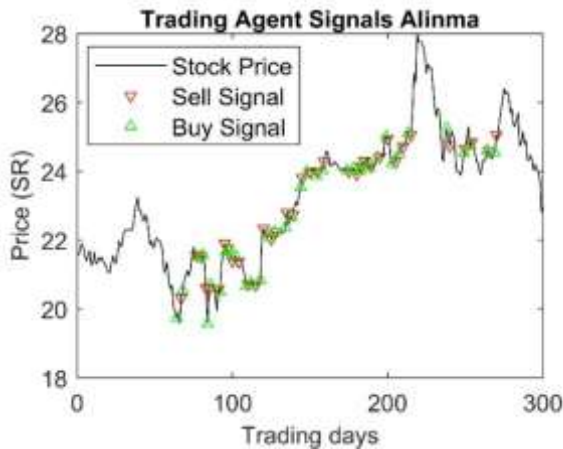


Fig. 8. Trading Agent Signal when used in the Alinma Stock Trading.



Fig. 9. Trading Agent Signal when used in the Alrajhi Stock Trading.



Fig. 10. Trading Agent Signal when used in the SABIC Stock Trading.

Fig. 10 shows the trading signal of the SABIC shares. The agent predicted the correct trading signals when trading the SABIC shares, but the gain was not high compared with that of the other two shares because of the fixed amounts of shares that can be bought per trading session. This low gain can be fixed if the number of shares is dynamic and linked to the classification layer output score.

VI. DISCUSSION AND CONCLUSION

To predict the Saudi stock trading signals, we proposed the usage of a multivariate RNN with an LSTM architecture. The model used historical stock information, such as closing prices, the volume of trades, number of trades, current-day opening prices, and oil price. The model result was satisfying compared with that obtained using the buy-and-hold trading method.

In future studies, we must consider more factors, such as the Fibonacci retracement, and develop a feature selection method to select the best feature among the presented features. Other financial trading methods may also be considered to train a neural network and develop a trading agent instead of relying on the prediction of future returns.

ACKNOWLEDGMENT

The authors would like to thank Deanship of scientific research in King Saud University for funding and supporting this research through the initiative of DSR Graduate Students Research Support (GSR).

REFERENCES

- [1] Shah, H. Isah and F. Zulkernine, "Stock Market Analysis: A Review and Taxonomy of Prediction Techniques," International Journal of Financial Studies, vol. 7, p. 26, 5 2019.
- [2] N. Singh, N. Khalfay, V. Soni and D. Vora, "Stock Prediction using Machine Learning a Review Paper," International Journal of Computer Applications, vol. 163, p. 36-43, 4 2017.
- [3] P. Parracho, R. Neves and N. Horta, "Trading in financial markets using pattern recognition optimized by genetic algorithms," in Proceedings of the 12th annual conference comp on Genetic and evolutionary computation - GECCO, Portland, 2010.
- [4] M. Velay and F. Daniel, "Stock Chart Pattern recognition with Deep Learning," 1 8 2018.
- [5] C. D. Tilakaratne, M. A. Mammadov and S. A. Morris, "Predicting Trading Signals of Stock Market Indices Using Neural Networks," in AI 2008: Advances in Artificial Intelligence, Springer Berlin Heidelberg, 2008, p. 522-531.

- [6] R. Dash and P. K. Dash, "A hybrid stock trading framework integrating technical analysis with machine learning techniques," *The Journal of Finance and Data Science*, vol. 2, p. 42–57, 3 2016.
- [7] A. H. Moghaddam, M. H. Moghaddam and M. Esfandyari, "Stock market index prediction using artificial neural network," *Journal of Economics, Finance and Administrative Science*, vol. 21, p. 89–93, 12 2016.
- [8] J.-L. Seng and H.-F. Yang, "The association between stock price volatility and financial news – a sentiment analysis approach," *Kybernetes*, vol. 46, p. 1341–1365, 9 2017.
- [9] S. M. Idrees, M. A. Alam and P. Agarwal, "A Prediction Approach for Stock Market Volatility Based on Time Series Data," *IEEE Access*, vol. 7, p. 17287–17298, 2019.
- [10] S. O. Olatunji, M. S. Al-Ahmadi, M. Elshafe and Y. A. Fallatah, "Forecasting the Saudi Arabia Stock Prices Based on Artificial Neural Networks Model," *International Journal of Intelligent Information Systems*, vol. 2, p. 77, 2013.
- [11] M. Jarrah and N. Salim, "A Recurrent Neural Network and a Discrete Wavelet Transform to Predict the Saudi Stock Price Trends," *International Journal of Advanced Computer Science and Applications*, vol. 10, 2019.
- [12] T. Alotaibi, A. Nazir, R. Alroobaea, M. Alotibi, F. Alsubeai, A. Alghamdi and T. Alsulimani, "Saudi Arabia Stock Market Prediction Using Neural Network," *International Journal on Computer Science and Engineering*, vol. 9, p. 62–70, 2 2018.
- [13] Y. Hua, Z. Zhao, R. Li, X. Chen, Z. Liu and H. Zhang, "Deep Learning with Long Short-Term Memory for Time Series Prediction," *IEEE Communications Magazine*, vol. 57, p. 114–119, 6 2019.
- [14] E. F. Fama, "Efficient Capital Markets: A Review of Theory and Empirical Work," *The Journal of Finance*, vol. 25, p. 383, 5 1970.
- [15] U. Awan and M. Subayyal, "Weak Form Efficient Market Hypothesis Study: Evidence from Gulf Stock Markets," *SSRN Electronic Journal*, 2016.
- [16] B. Asiri and H. Alzeera, "Is the Saudi stock market efficient? A case of weak-form efficiency," *Research Journal of Finance and Accounting*, vol. 4, no. 6, pp. 35–48, 2013.
- [17] F. J. Massey Jr, "The Kolmogorov-Smirnov test for goodness of fit," *Journal of the American statistical Association*, vol. 46, p. 68–78, 1951.
- [18] C. C. Aggarwal, *Neural Networks and Deep Learning*, Springer-Verlag GmbH, 2018.
- [19] D. E. Rumelhart, G. E. Hinton and R. J. Williams, "Learning representations by back-propagating errors," *Nature*, vol. 323, p. 533–536, 10 1986.
- [20] S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," *Neural Computation*, vol. 9, p. 1735–1780, 11 1997.
- [21] F. A. Gers, J. Schmidhuber and F. Cummins, "Learning to Forget: Continual Prediction with LSTM," *Neural Computation*, vol. 12, p. 2451–2471, 10 2000.
- [22] C. Bishop, *Pattern recognition and machine learning*, New York: Springer, 2006.
- [23] D. P. Kingma and J. Ba, "Adam: A Method for Stochastic Optimization," 22 12 2014.

# Implementing the Behavioral Semantics of Diagrammatic Languages by Co-simulation

Daniel-Cristian Crăciunean<sup>1</sup>

Computer Science and Electrical Engineering Department  
Lucian Blaga University of Sibiu, Sibiu, Romania

**Abstract**—Due to the multidisciplinary nature of cyber-physical systems, it is impossible for an existing modeling language to be used effectively in all cases. For this reason, the development of domain-specific modeling languages is beginning to become an integral part of the modeling process. This diversification of modeling languages often implies the need to co-simulate subsystems in order to obtain the effect of a complete system. This paper presents how behavioral semantics of a diagrammatic DSML can be implemented by co-simulation. For the formal specification of the language we used mechanisms from the category theory. To specify behavioral semantics, we introduced the notion of behavioral rule as an aggregation between a graph transformation and a behavioral action. The paper also contains a relevant example and demonstrates that the implementation of behavioral semantics of a diagrammatic model can be achieved by co-simulating standalone FMUs associated to behavioral rules.

**Keywords**—DSML; cyber-physical systems; behavioral semantics; standalone FMU; FMI; diagrammatic language

## I. INTRODUCTION

In the context of moving the effort from writing code to writing models, the development of modeling tools, appropriate to the domain of modeling, becomes an essential factor for increasing the efficiency of the modeling process. The diagrammatic syntax of domain-specific modeling languages (DSML) seems to be the most accessible for all parties involved in the model specification, because it is intuitive and can provide support in all phases of model development, starting with the informal model and ending with the executable model [1,2].

Models specified with these DSMLs must, in turn, interact with other models specified in other languages. Often the models specified with these DSMLs assemble heterogeneous components, which must be modeled with other languages. All these components can be specified in various modeling languages. But there is a need for a specific language to assemble the system components into a workflow [3] and coordinate the behavior of these components. In our opinion, these interaction problems can be solved elegantly by co-simulation [4].

One of the main objectives of building a model is to study the behavior of a system in order to analyze and optimize the modeled system. Due to the complexity of the systems, classical optimization methods cannot be used and therefore must be replaced by methods based on simulation or genetic algorithms. To achieve these objectives the model will have to

be executed by a simulator according to its behavioral rules to mimic the behavior of the system.

Complex systems such as Cyber-Physical Production Systems (CPPS) also have a high degree of heterogeneity and therefore involve components with different behaviors that cannot be efficiently specified in the same formalism. In these cases, we need a co-simulation environment that combines several simulators into one and that reproduces the behavior of the global system [5].

In order for these heterogeneous models to be coupled in the co-simulation process, they need to provide a common standardized interface. This interface is called Functional Mock-up Interface (FMI) [6] introduced in the European MODELISAR project, carried out in the period 2008-2011.

To achieve the goal of co-simulation, modeling tools must be able to generate co-simulation units with FMI interfaces, which are called Functional Mock-up Units (FMU). The orchestration of the components in order to obtain the behavior of the composite system is done by an orchestrator which is called master algorithm.

We believe that for the efficient implementation of a DSML, co-simulation mechanisms must be an integral part in the process of specifying and implementing a modeling tool. In this paper we present the methodology for specifying and implementing a DSML with FMU generation facility. To formalize the model, we use mechanisms from category theory. For co-simulation we used the INTO-CPS [7] tool chain. INTO-CPS is an EU-funded project that integrates a chain of tools for model-based CPS design and implementation by co-simulating components with an FMI-compatible interface.

In Section 2, we briefly specify the static metamodel of a diagrammatic model. In Section 3, we specify the behavioral syntax of the model and in Section 4, we deal with the semantic mapping of a model. In Section 5, we briefly present the mechanism for generating FMU components. Section 6 concludes the paper with original contributions and conclusions. All the mechanisms presented are exemplified with a simple model that was implemented on the ADOxx metamodeling platform.

## II. THE STATIC MODEL

In essence, a visual model of a system first defines the syntax of the static and behavioral model that represents the virtual and physical entities of the model and then the

semantics of the model represented by the significance of static constructions and a set of behavioral rules that represent the behavior of these entities.

Syntactically, a diagrammatic model is a graph with several types of nodes that represent different concepts in the domain of modeling and several types of arcs that represent links between these concepts [8,9]. When we want to associate models with a spatial representation, we can use a second graph, as a spatial dimension of them and thus we reach the notion of bigraph [10]. The models discussed in this paper have as syntactic representation a single graph.

Example 1. We consider a modeling language SML (Simple Modeling Language) that has the following concepts:

A buffer concept, which can store a single type of material, which we denote by  $B_1$ , and endow it with two attributes, namely: the stock attribute which represents the current quantity stored in the buffer and capacity which represents the maximum quantity that can be stored in the buffer. We associate to this concept the following graphic notation:



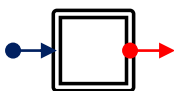
A buffer concept, which can store two types of materials, which we denote by  $B_2$ , and endow it with four attributes, namely: the attributes stock1, stock2 which represents the current quantity of each type stored in the buffer and capacity1, capacity2 which represents the maximum quantity of each type, which can be stored in the buffer. We associate to this concept the following graphic notation:



A processing or transfer activity concept, which we denote by  $W_1$ , and which can process or transfer materials from a type  $B_1$  buffer to a type  $B_2$  buffer, and endow it with three attributes, namely, the StockIn attribute which represents the quantity of material fed from buffer  $B_1$ , and the attributes Stock1Out, Stock2Out which represent the quantities of material of each type deposited in buffer  $B_2$ . We associate to this concept the following graphic notation:



A processing or transfer activity concept, which we denote by  $W_2$ , and which can process or transfer materials from a type  $B_2$  buffer to a type  $B_1$  buffer, and endow it with three attributes, namely: the attributes Stock1In, Stock2In which represents the quantities of material fed from each type and the StockOut attribute which represents the quantity of material deposited in buffer  $B_1$ . We associate to this concept the following graphic notation:



The SML model that we will specify is a graphical DSML for describing simple models in conformity with the requirements specified above.

SML models, therefore, are graphs with a set of syntactic restrictions on their components [11]. In the categorical model, the SML metamodel is a sketch that is composed of a graph and a set of constraints on the graph nodes [2,12,13].

Example 2. We will define a SML model as a graph  $G=(X,\Gamma,\sigma,\theta)$ , on the components of which we introduce four restrictions, namely:

- 1) The nodes of the graph are of two types and these types determine a partition on  $X$ , i.e.:  $X=B_1 \sqcup B_2$ ;
- 2) The arcs of the graph are of two types and these types determine a partition on  $\Gamma$ , i.e.:  $\Gamma=W_1 \sqcup W_2$ ;
- 3) Graph  $G$  has to be a connected graph;
- 4) There must be at most one arc between any two components.

A categorical sketch is a tuple  $\mathcal{S}=(\mathcal{G},\mathcal{C}(\mathcal{G}))$  where  $\mathcal{G}$  is a graph and  $\mathcal{C}(\mathcal{G})$  is a set of constraints on the set of nodes and arcs of the graph [2]. A model of the sketch  $\mathcal{S}=(\mathcal{G},\mathcal{C}(\mathcal{G}))$  is the image of this sketch through a functor in the Set category.

From the way of defining the SML model, from example 2 it results that the graph  $\mathcal{G}$  of the corresponding sketch  $\mathcal{S}=(\mathcal{G},\mathcal{C}(\mathcal{G}))$  is the one from Fig. 1.

The categorical sketches are based on the observation that a labeled diagram is an analogous construction of a logical formula that is mapped to the components of a graph, i.e. to the nodes and arcs of a graph [2,14].

We denote with Graph the category of graphs, i.e. the category that has graphs as objects and as arcs the homomorphisms between these graphs. We will also denote with  $\text{Graph}_0$  the set of objects of the Graph category and with  $\text{Graph}_1$  the set of arcs of the Graph category.

Constraints on the models specified by the categorical sketch are defined by a predicate signature diagram, which is composed of a set of predicates  $\Pi$ , and an application:  $\Pi \rightarrow \text{Graph}_0$ , which maps the indeterminate predicates to the nodes of a graph in  $\text{Graph}_0$  [2]. This predicate signature diagram, allows the definition of constraints on the models specified by a categorical sketch at the metamodel level.

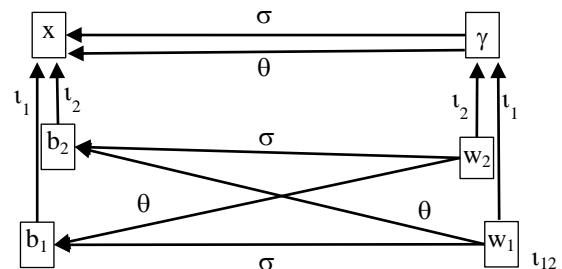


Fig. 1. The Graph of the SMM Sketch.

For example, for the graph  $G$  to be connected we will put the condition that the pushout of  $\sigma$  with  $\theta$  to be a terminal object in the Set category.

If we denote  $\text{Span}(x,y,z,r_{zx},r_{zy})=(x \xleftarrow{r_{zx}} z \xrightarrow{r_{zy}} y)$  then the pushout of  $\sigma$  with  $\theta$ , in the Set category, is the colimit of the diagram  $d:\text{Span}(1,2,3,r_{31},r_{32})\rightarrow\text{Set}$  where  $d(1)=x$ ,  $d(2)=x$ ,  $d(3)=y$ ,  $d(r_{31})=\sigma$ ,  $d(r_{32})=\theta$ . This constraints are imposed by the predicate  $P(n_1,n_2,n_3,r_{31},r_{32})$  with the shape graph arity of  $P_3$ ,  $\text{ar}(P(n_1,n_2,n_3,r_{31},r_{32}))=\text{Span}(1,2,3,r_{31},r_{32})$  defined as:  $\text{ar}(n_1)=1$ ,  $\text{ar}(n_2)=2$ ,  $\text{ar}(n_3)=3$ ,  $\text{ar}(a_{31})=r_{31}$ ,  $\text{ar}(a_{32})=r_{32}$ . In these conditions the predicate  $P(n_1,n_2,n_3,r_{31},r_{32})$  is defined as follows:  $P(n_1,n_2,n_3,r_{31},r_{32})=|\text{CoLim}(d)|=1$  where  $\text{CoLim}(d)$  is the colimit of diagram  $d$  in the Set category.

Therefore, the categorical sketch of the SML model has the following components: the graph of the sketch is the one from Fig. 1, and the set of constraints  $\mathcal{S}(\Pi)$  is obtained by mapping the shape graphs corresponding to the predicates from  $\Pi$  to the components of the sketch graph by means of diagrams, i.e.  $\mathcal{S}(\Pi)=\{\mathcal{S}(P_i) \mid P_i \in \Pi, i \geq 1\}$ . The categorical sketch  $\mathcal{S} = (\mathcal{G}, (\Pi))$  represents the abstract syntax of the SML models and at the same time the SML metamodel.

Each model specified by the categorical sketch  $\mathcal{S}$ , is the image of the graph  $\mathcal{G}$  of the sketch  $\mathcal{S}$  through a functor  $M$ , in the Set category, which respects the constraints imposed by the predicates  $(\Pi)$ . The predicates in the set  $(\Pi)$  will be mapped, at the level of each model  $M$ , from the Set category to a set of predicates as follows:  $\text{Set}(P_i)=\{(P_i; M \circ d \circ \text{ar}(P_i)) \mid d \text{ is a diagram}\}$ .

Thus, if we have the model  $M:\mathcal{S}\rightarrow\text{Sets}$ , where  $M(b_1)=B_1$ ,  $M(b_2)=B_2$ ,  $M(w_1)=W_1$  and  $M(w_2)=W_2$ , then the set of instances  $B_1, B_2, W_1, W_2$ , will respect the constraints defined by the set of predicates  $\text{Set}(P_i)$ . We notice that the graph of the categorical sketch contains besides the concepts from the modeling domain, also auxiliary nodes useful for imposing constraints.

If in the above model we have:  $B_1=\{B_{11},B_{12},B_{13}\}$ ;  $B_2=\{B_{21},B_{22}\}$ ;  $W_1=\{W_{11},W_{12}\}$ ;  $W_2=\{W_{21}, W_{22},W_{23}\}$  and  $\sigma(W_{11})=B_{11}$ ;  $\sigma(W_{12})=B_{12}$ ;  $\theta(W_{11})=B_{22}$ ;  $\theta(W_{12})=B_{22}$ ;  $\sigma(W_{21})=B_{21}$ ;  $\sigma(W_{22})=B_{21}$ ;  $\sigma(W_{23})=B_{22}$ ;  $\theta(W_{21})=B_{11}$ ;  $\theta(W_{22})=B_{12}$ ;  $\theta(W_{23})=B_{13}$ , then the SML model is like in Fig. 2.

We will consider that the nodes of the graph of the sketch  $\mathcal{S}$  are classes endowed with attributes. The graph nodes will be mapped by the functor  $M$  to sets of objects of the corresponding class type in the Set category, and the graph arcs will be mapped to functions between these sets. The semantics of such a static model is given by the significance of the attributes, the significance of the values of these attributes and the significance of the graph structure of the model.

A class defines a concrete modeling concept that can be used to specify a model in the modeling language. Therefore, each concrete concept of a model created with a tool implemented on the ADOxx platform is an instance of a class. Each concrete class has a distinct name.

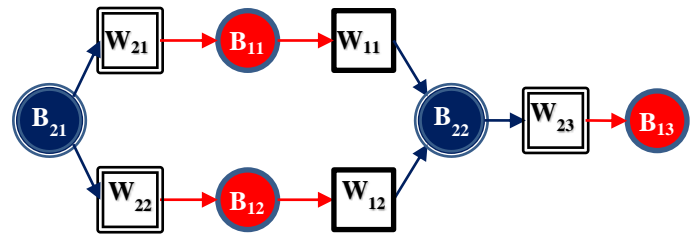


Fig. 2. Model Example.

### III. BEHAVIORAL SYNTAX OF THE MODELING LANGUAGE

Model transformation is one of the key techniques in MDE used especially for automating model management operations, such as code generation, model optimization, translation from one DSML to another, simulation, etc.

The transformation of diagrammatic models is based on, most often, the graph transformations defined by graph rules, also called productions. Such a production is a tuple  $p=(L, R)$ , consisting of two graphs; a left graph  $L$ , a right graph  $R$  and a mechanism that specifies the conditions and how to replace  $L$  with  $R$ .

In this paper we will use graph transformations to model the behavior of a diagrammatic model. Graph transformation rules are mechanisms that can express the local changes of a graph in successive transformation steps ordered by a relationship of causal dependence of actions and therefore can accurately define the behavior of a diagrammatic model.

In the approaches of implementing the transformations, of the left graph to the right graph, two distinct mechanisms are distinguished, namely, the double pushout (DPO) and single pushout (SPO) [15,16]. Graph transformations allow the simultaneous transformation of the structure of the diagrammatic model and of the attributes of the components. In this paper we will use the DPO variant to specify the behavioral dimension of a model.

The correct application of a production  $p$  is made under the conditions in which the squares in Fig. 3 are pushout squares. A set of productions related to each other form a graph transformation system and can be used in the process of transforming models without being integrated into a graph grammar [15,17,18].

The behavior of SML models is not based on structural transformations of the graph but only on changing attribute values and therefore we will use graph transformations only to specify the context necessary to locate the components involved in a behavioral rule and to locate critical regions in the simulation process.

We will define the graph transformations at the metamodel level and they will be applied for any static model specified by the sketch in the Set category.

In the case of our SML model we have two transformations at the level of the sketch  $\mathcal{S}$  from Fig. 1, namely,  $p_1$  (Fig. 4) and  $p_2$  (Fig. 5).

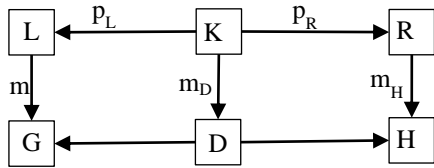


Fig. 3. A Double-Pushout Production.

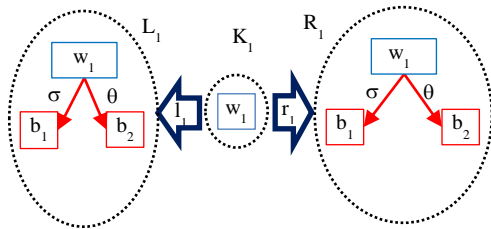


Fig. 4. Graph Transformation  $p_1$ .

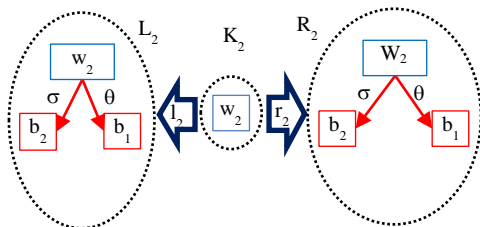


Fig. 5. Graph Transformation  $p_2$ .

#### IV. SEMANTICS OF THE MODELING LANGUAGE

A diagrammatic model is characterized by two dimensions: a static dimension, specified syntactically by the categorical sketch, and a behavioral dimension, specified syntactically by behavioral signatures. The two dimensions have dependent but still distinct semantics. Defining the semantics of the static dimension involves mapping the attributes with which the components of the graph of the categorical sketch are endowed to data domains and the graph structures of the model to well defined semantically structures, and defining behavioral semantics involves mapping the behavioral signature to mathematical functions.

The semantics of a model involves mapping attributes to their set of values, interpreting the graph structure of the model, while behavioral semantics highlights the structural and value transformations of the model.

The behavioral dimension of a SML model is defined by states and transitions. The states of the model are represented, in our approach, by the static models of the categorical sketch, and the transitions are represented by a set of mathematical functions associated with the behavioral rules. A behavioral rule can be applied only in the context in which a set of conditions are realized, represented by logical predicates that verify the state of the model. Therefore, a behavioral rule is defined as an association between a graph transformation and an action.

We will define the behavioral rules at the metamodel level by behavioral signatures represented by predefined actions

with formal parameters, which will be evaluated at the moment of execution of a model, when they will receive current parameters corresponding to the respective concrete model.

Let the sets  $Y_1, \dots, Y_m, U_1, \dots, U_n$ . Then an action is an application  $\text{Act}: U_1 \times \dots \times U_n \rightarrow Y_1 \times \dots \times Y_m$  defined as follows:  $(y_1, \dots, y_m) := \text{Act}(u_1, \dots, u_n) = (\omega_1(u_1, \dots, u_n), \dots, \omega_m(u_1, \dots, u_n))$  where  $\omega_i$  is an operation that calculates the value of  $y_i$ ,  $i=1, m$  depending on the values of the variables  $u_1, \dots, u_n$ .

We denote with  $\text{AGraph}$  the category of graphs with attributes, i.e. the category that has as objects graphs with attributes and as arcs the homomorphisms between these graphs. Also, if  $G \in \text{AGraph}_0$ , we will denote with  $\text{attr}(G)$  the set of attributes associated with the nodes and arcs of the graph  $G$ .

A diagram actions signature is a tuple  $\Delta = (\mathcal{A}, \text{ar})$  where  $\mathcal{A}$  is a set of actions and  $\text{ar}$  is a function  $\text{ar}: \mathcal{A} \rightarrow \text{Graph}_0$  which maps each action  $\text{Act} \in \mathcal{A}$  to two objects in the  $\text{AGraph}$  category as follows: if  $(y_1, \dots, y_m) := \text{Act}(u_1, \dots, u_n)$  then the outputs  $y_1, \dots, y_m$  will be mapped to the attributes of the graph  $R$  and the inputs  $u_1, \dots, u_n$  will be mapped to the attributes of the graph  $L$ . The pair  $(L, R)$  of graphs is called shape graph arity of  $\text{Act}$ ,  $\text{ar}(\text{Act}) = (L, R)$ . We will sometimes denote the image of  $\text{Act}$  through  $\text{ar}$  in the category  $\text{AGraph}$  with  $\text{Act}(L, R)$ .

The behavioral signature is a tuple  $\Sigma = (\mathcal{T}, C_L, \Delta, C_R)$  where  $\mathcal{T}$  is a set of graph transformation rules;  $C_L = (\Pi_L, \text{ar}_L)$  is a diagram predicate signature such that  $\text{ar}_L: \Pi_L \rightarrow \text{AGraph}_0$ , which we call the precondition signature;  $C_R = (\Pi_R, \text{ar}_R)$  is a diagram predicate signature such that  $\text{ar}_R: \Pi_R \rightarrow \text{AGraph}_0$ , which we call the postcondition signature and  $\Delta$  is a diagram actions signature, with the property that for any  $\text{Act} \in \Delta$  there is  $p \in \mathcal{T}$ ,  $p = L \xleftarrow{l} K \xrightarrow{r} R$  with  $\text{ar}(\text{Act}) = (L, R)$ , that specifies how to transform the attributes of graph  $L$  which is the domain of action into the components of graph  $R$  which composes the codomain of the action.

We now denote the graph in Fig. 7 with  $G_1(x_1, x_2, x_3)$  and the graph with a single node in Fig. 8 we denote it with  $G_2(x_1)$ . The shape graphs represent the local structures of a model and represent the areas of action of the behavioral rules in the context of a concrete model.

These shape graphs represent the local structure of the model, and the context in which a behavioral rule evolves. Behavioral signatures defined on the components of these shape graphs are mapped to behavioral transformations on the component elements of a model.

The behavior of the SML model can be specified by a behavioral signature that contains two behavioral rule signatures  $\sigma_1$  and  $\sigma_2$ . Since the set of behavioral rule signatures is equivalent to the behavioral signature, we will use the same notation  $\Sigma$  for the set of behavioral rule signatures.

So  $\Sigma = \{\sigma_1, \sigma_2\}$  where:

$$\sigma_1 = (L^1 \xleftarrow{l_1} K^1 \xrightarrow{r_1} R^1, C_L^1, \text{Act}^1, C_R^1); \sigma_2 = (L^2 \xleftarrow{l_2} K^2 \xrightarrow{r_2} R^2, C_L^2, \text{Act}^2, C_R^2);$$

$$L_1 = R_1 = L_2 = R_2 = G_1(1,2,3) ; K_1 = K_2 = G_2(1);$$

$$C_L^1 = (\Pi_L^1, ar_L^1); \Pi_L^1 = \{P_L^1(u_1, \dots, u_n)\}; ar_L^1(u_i) = a_i, i=1, n \text{ and } a_i \in \text{attr}(L^1);$$

$$C_R^1 = (\Pi_R^1, ar_R^1); \Pi_R^1 = \{P_R^1(y_1, \dots, y_m)\}; ar_R^1(u_i) = b_i, i=1, m \text{ and } b_i \in \text{attr}(R^1);$$

$$L_2 = R_2 = G_1(1,2,3) ; K_2 = G_2(1);$$

$$C_L^2 = (\Pi_L^2, ar_L^2); \Pi_L^2 = \{P_L^2(u_1, \dots, u_n)\}; ar_L^2(u_i) = a_i, i=1, n \text{ and } a_i \in \text{attr}(L^2);$$

$$C_R^2 = (\Pi_R^2, ar_R^2); \Pi_R^2 = \{P_R^2(y_1, \dots, y_m)\}; ar_R^2(u_i) = b_i, i=1, m \text{ and } b_i \in \text{attr}(R^2);$$

The behavioral signatures thus defined will be transformed into behavioral rules at the level of the metamodel, by mapping them to the components of the sketch  $\mathcal{S}$ , and will represent the behavioral model at the level of the metamodel, i.e. the abstract behavioral semantics of the models. The behavioral rules at the level of the sketch  $\mathcal{S}$  will then be mapped by matches at the level of the models.

The behavioral rule signatures must be mapped to the components of the graph  $\mathcal{G}$  of the sketch  $\mathcal{S}$  by sets of three diagrams, one for each of the graph forms L, R and K. These will be defined by three functors  $d_L$ ,  $d_K$  and  $d_R$ , where  $d_K$  is the restriction of the functors  $d_L$  and  $d_R$  at domain K;  $d_K = d_L/K = d_R/K$ ,  $l_s$  and  $r_s$  are monomorphisms that inject the graph K into L and R, respectively.

We will therefore define the diagrams corresponding to the signature of the behavioral rule  $\sigma_1$ :

$$d_L^1: G_1(1,2,3) \rightarrow G_1(\gamma_{12}, w_1, w_2) \text{ defined as } d_L^1(1) = \gamma_{12}; d_L^1(2) = w_1; d_L^1(3) = w_2;$$

$$d_R^1 = d_L^1; \text{ and } d_K^1: G_2(1) \rightarrow G_2(\gamma_{12}) \text{ defined as restriction } d_K^1 = d_L^1/K^1; d_K^1(1) = \gamma_{12}.$$

And for the signature of rule  $p_2$  we have the diagrams:

$$d_L^2: G_1(1,2,3) \rightarrow G_1(\gamma_{21}, w_2, w_1) \text{ defined as } d_L^2(1) = \gamma_{21}; d_L^2(2) = w_2; d_L^2(3) = w_1;$$

$$d_R^2 = d_L^2; \text{ and } d_K^2: G_2(1) \rightarrow G_2(\gamma_{21}) \text{ defined as restriction } d_K^2 = d_L^2/K^2; d_K^2(1) = \gamma_{21}.$$

Diagrams are functors that map the formal parameters defined by graph shapes to the concepts specified by the nodes of the sketch graph. The same graph shapes are, on the other hand, mapped to the components of a concrete model through matching applications.

A behavioral rule of the sketch  $\mathcal{S}$  is a tuple  $t = (L \xleftarrow{l} K \xrightarrow{r} R, d_L(C_L), \text{Act}(d_L(L); d_R(R)), d_R(C_R))$  where  $\sigma = (L \xleftarrow{l} K \xrightarrow{r} R, C_L, \text{Act}(L; R), C_R)$  is a signature of a behavioral rule. We used the following notations:  $d_L(C_L) = (\Pi_L, d_L(ar_L)); d_R(C_R) = (\Pi_R, d_R(ar_R)).$

Thus, starting from a behavioral signature, we generate a set of behavioral rules at the level of the metamodel, i.e. at the level of the components of the graph of the sketch.

If we denote with  $(\Sigma)$ , the set of behavioral rules induced by the behavioral signature  $\Sigma$ , then a behavioral metamodel is a tuple  $(\mathcal{G}, \mathcal{S}(\Sigma))$  where  $\mathcal{G}$  is the graph of the sketch  $\mathcal{S} = (\mathcal{G}, \mathcal{C}(\mathcal{G}))$ .

In our approach each of these behavioral rules will be implemented as an FMU component. The behavioral metamodel corresponding to the SML language is defined by two behavioral rules  $(\Sigma) = \{\mathcal{S}(\sigma_1), \mathcal{S}(\sigma_2)\}$  where:

$$(\sigma_1) = (L^1 \xleftarrow{l_1} K^1 \xrightarrow{r_1} R^1, \{P_L^1(d_L^1(L^1))\}, \text{Act}^1(d_L^1(L^1); d_R^1(R^1)), P_R^1(d_R^1(R^1)));$$

$$(\sigma_2) = (L^2 \xleftarrow{l_2} K^2 \xrightarrow{r_2} R^2, \{P_L^2(d_L^2(L^2))\}, \text{Act}^2(d_L^2(L^2); d_R^2(R^2)), P_R^2(d_R^2(R^2)));$$

Thus, for our SML metamodel we will implement two FMU components corresponding to the two behavioral rules  $(\sigma_1)$  and  $(\sigma_2)$ .

In order for the behavioral rules specified in the metamodel  $(\Sigma)$  to be applied at the level of a concrete model we will have to find the matches of each behavioral rule from  $(\Sigma)$  in a model from  $\text{Mod}(\mathcal{S}, \text{Set})$ .

A match of a graph  $\mathcal{G} = (N, A, s, t)$  in the image of a functor  $\phi: \mathcal{G} \rightarrow \text{Set}$  is a total monomorphism of graphs  $m: \mathcal{G} \rightarrow \phi(\mathcal{G})$  which maps the graph  $\mathcal{G}$  to the graph  $\mathcal{G}_m = (m(N), m(A), m(s), m(t))$  so that  $\forall y_i \in m(N) \Rightarrow \exists x_i \in N$  with  $y_i \in \phi(x_i)$  and  $\forall a_i \in m(A) \Rightarrow \exists r_i \in N$  with  $a_i \in \phi(r_i)$  respecting the conditions of homomorphism  $m(s(r_i)) = m(s)(m(r_i))$  and  $m(t(r_i)) = m(t)(m(r_i))$  for all  $r_i \in A$ . We will denote the set of graph matches  $\mathcal{G}$  in  $\phi(\mathcal{G})$  with  $m(\phi, \mathcal{G})$ .

In this way the graph transformations and the actions on the attributes will be executed on a concrete model.

Under these conditions, a behavioral model, in the Set category, contains all the behavioral rules induced by the behavioral signature  $\Sigma$ , in the Set category. We notice that the set of behavioral rules is specific to each concrete model, but they can be implemented generically at the metamodel level.

As we can see each behavioral rule  $\tau$  in Set, defines an application  $\tau: M_L \rightarrow M_R$ , where  $M_L, M_R: \mathcal{S} \rightarrow \text{Set}$  are functors which represents the domain and codomain of the rule  $\tau$  and all these behavioral applications together, maps the set  $\mathcal{S}(\Sigma)$  of behavioral rules of the sketch into a set of behavioral rules  $\text{Set}(\Sigma)$  in Set.

In the case of the SML language the atomic behavioral rules in Set,  $\tau$  are of the form  $\tau = (t, \mu, M_L, M_R) \in \mathcal{S}(\Sigma)$  where  $t \in \mathcal{S}(\Sigma)$  is a behavioral rule  $t = (L \xleftarrow{l} K \xrightarrow{r_s} R, d_L(C_L), \text{Act}(d_L(L); d_R(R)), d_R(C_R))$ ,  $M_L, M_R: \mathcal{S} \rightarrow \text{Mod}(\mathcal{S}, \text{Set})$  there are two functors, and  $\mu$  is a tuple of match  $\mu = (m_L, m_K, m_R)$ ,  $m_L \in m(M_L, L)$ ,  $m_R \in m(M_R, R)$ , and  $m_K$  is the restriction of  $m_L$  to K, so that the diagram in Fig. 6 is a double pushout.

For the model from Fig. 2 we have 5 behavioral rules in  $\text{Set}(\Sigma)$ , two for  $\mathcal{S}(\sigma_1)$  and three for  $\mathcal{S}(\sigma_2)$ :  $\text{Set}(\Sigma) = \{\tau_{11}, \tau_{12}, \tau_{21}, \tau_{22}, \tau_{23}\}$ .



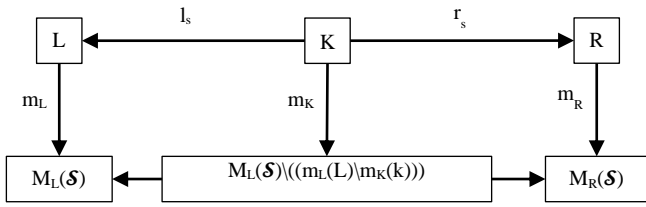


Fig. 6. A Double-Pushout Diagram.

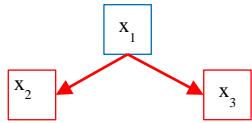


Fig. 7. Graph  $G_1(x_1, x_2, x_3)$ .



Fig. 8. Graph  $G_2(x_1)$ .

In our approach, behavioral rules are approached in two distinct phases. In the first phase, these rules are defined, at the metamodel level, by behavioral signatures, and in the second phase, these rules are applied at the level of each concrete model. If the behavioral rules of the model are faithful to the modeled system, then their successive application mimics the behavior of the modeled system.

### V. IMPLEMENTATION OF THE FMU GENERATOR

From the formalization presented in this paper results the fact that a diagrammatic metamodel has two dimensions, a static dimension represented by the categorical sketch and a behavioral dimension represented by the behavioral rules. A behavioral rule, as we have defined it, is an aggregation between a graph transformation on the structure of a model and a local action on the attributes of the model. If the functionalities of a metamodeling platform are designed to specify the graph structure of a metamodel and can be endowed with graph transformation facilities, behavioral actions are often performed by complex systems with a high degree of heterogeneity which implies the need for modeling on various modeling platforms. A solution to this problem is the assembly through co-simulation of  $n+1$  independent components where  $n$  is the number of behavioral rules. In other words, you can build an FMU component that manages the static dimension of the model together with the graph transformations on it and an FMU component for each behavioral rule that models the action corresponding to the behavioral rule.

Applying a transformation rule specified by a behavioral signature  $\sigma=(L \xleftarrow{l_s} K \rightarrow R, C_L, Act, C_R)$  is done as follows:

1) We first consider the diagrams  $d_L$  and  $d_R$  which maps the behavioral signature  $\sigma$  to the model sketch. In this way the components of the diagrams receive the types of components of the sketch.

When we are going to apply a behavioral rule  $\tau=(t, \mu, M^L, M^R) \in \mathcal{S}(\Sigma)$  we have the first component  $M_L$ , which represents the current state of the behavioral model, and we are going to determine the  $M_R$  component which represents the state in which the transition is made. Therefore we can find the matches  $m_L \in m(M_L, L)$  and  $m_K \in m(M_K, K) = m(M_L, L)$  which are the first two components of a match.  $\mu = (m_L, m_K, m_R)$  (Fig. 11).

2) The preconditions are verified, i.e. the fulfillment of the predicates defined by the  $C_L$  signatures, among which is the gluing condition. If the  $C_L$  conditions are met the graph transformation defined by the cospan  $L \leftarrow K \rightarrow R$  is executed in two steps, 1 and 2.

a) We calculate the complement  $M_L(\mathcal{S}) \setminus ((m_L(L), m_K(k)))$  of the pushout of  $l_s$  with  $m_k$ , from Fig. 9.

b) Now we can calculate the pushout of  $r$  with  $m_k$ , from Fig. 10 and therefore the functor  $M_R$  and the component  $m_R \in m(M_R, R)$  of the match  $\mu$ .

All these transformations are executed temporarily, i.e. with the possibility of being canceled.

3) In this phase, the Act action is temporarily executed.

4) If the postconditions are also verified then the transformations described at points 2 and 3 are permanently executed, otherwise they are canceled by a rollback operation.

Obviously, independent behavioral rules can be applied simultaneously, and also the same behavioral rule can be applied simultaneously to several areas of the model if these areas do not contain common elements.

The model was implemented on the ADOxx metamodeling platform (see Fig. 12). In the case of the SML metamodel we defined, as it results from the analysis of the graph  $\mathcal{G}$  of the sketch  $\mathcal{S}$ , two classes, corresponding to nodes  $b_1$  and  $b_2$  and two classes' relations corresponding to nodes  $w_1$  and  $w_2$ . Defining classes in ADOxx is done visually, but the metamodel can be exported in ADL language or XML format. We used the ADL language to generate from classes, C structure types for standalone FMU.

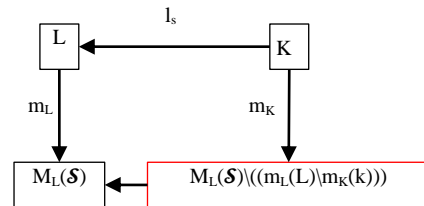


Fig. 9. The Complement of a Pushout.

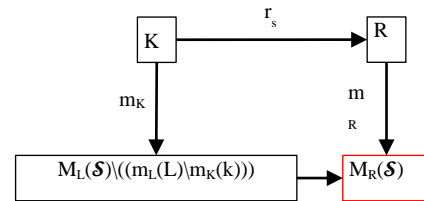


Fig. 10. A Pushout Diagram.

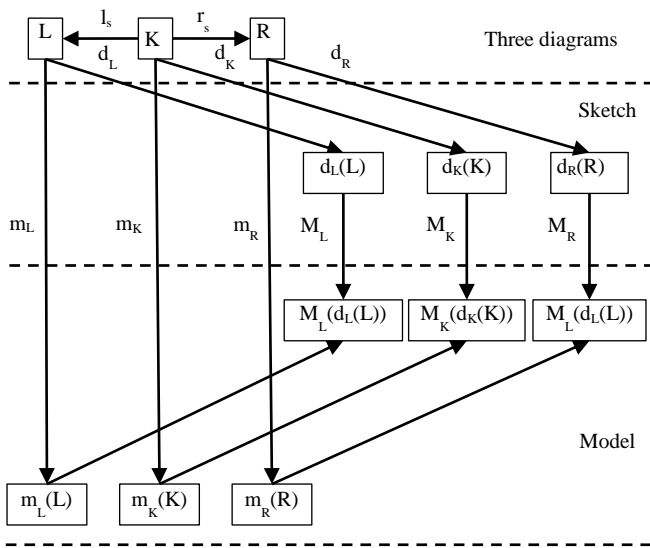


Fig. 11. Matching Three Diagrams.

The ADOxx metamodeling platform does not include behavior in classes and therefore we generated these classes as types of C structures. If the metamodeling platform would also include behavior this problem can be solved by function pointers. These structures were generated in a `<ModelName>_fmu_types.h` file, in our case `SML_fmu_types.h`. This type file is easy to write even manually, because it is written once and can then be used for all models specified with the implemented modeling tool.

We exported from ADOxx the model in ADL format from which we generated the FMU component corresponding to the static model, i.e. a graph structure corresponding to the specified model and with nodes that have corresponding types from the file `<ModelName>_fmu_types.h` and the name specified in the model. We put this graph structure in a file named `<ModelName>_fmu_structure.h`, in the case of the SML metamodel, `SML_FMU_structure.h`. Generating this graph structure is important because it is used in all specified models with the implemented modeling tool.

The behavioral part of the FMU component that manages the static dimension of the model was implemented manually in the C language. This is acceptable because it does not have a high complexity and is written only once for a modeling tool. The generation of this C code is possible but a translator from the ADOScript language to C should be implemented. To write this code we used FMU SDK [19] which can also be used in the case of generation from ADOScript.

FMU components corresponding to behavioral rules are usually written in another modeling tool. In the case of our SML metamodel, we specified the two components corresponding to the behavioral rules  $T_1$  and  $T_2$  in the VDM-RT language and exported them as standalone FMU.

Therefore, for the SML metamodel, we have 3 FMU components (Fig. 13) Which we briefly describe using the

notations from [20]. In this sense an FMU is defined as a tuple  $T_1 = \langle S_c, U_c, Y_c, set_c, get_c, doStep_c \rangle$ ; where:  $S_c$  is the space of states;  $U_c$  is the set of input variables;  $Y_c$  is the set of output variables;  $set_c: S_c \times U_c \times v \rightarrow S_c$  and  $get_c: S_c \times Y_c \rightarrow v$  are the input and output functions and  $doStep_c: S_c \times R_+ \rightarrow S_c$  is a function that calculates the state after a given step.

$$T_1 = \langle S_1, U_1, Y_1, set_1, get_1, doStep_1 \rangle;$$

$S_1 = N \times x | x \in N, 0 \leq x \leq capacity1 \} \times \{ (x, y) | x, y \in N, 0 \leq x \leq capacity1 \text{ and } 0 \leq y \leq capacity2 \}$ ;  $U_1 = T1u = \{ initial\_id, initial\_stock, initial\_stock1, initial\_stock2 \}$ ;  $Y_1 = T1y = \{ final\_id, final\_stock, final\_stock1, final\_stock2 \}$ ; The parameters  $PT_1$  of the component  $T_1$  are:  $PT_1 = \{ capacity, capacity1, capacity2, stockIn, stock1Out, stock2Out \}$ .

The  $doStep_1$  function implements only the action  $Act^1$  because we do not have structural transformations of the model. The precondition for the execution of the action  $Act^1$  is:  $initial\_stock \geq stockIn$  and  $capacity1 - initial\_stock1 \geq stock1Out$  and  $capacity2 - stock2 \geq stock2Out$ .

The action  $(final\_stock, final\_stock1, final\_stock2) = Act1(final\_stock, final\_stock1, final\_stock2)$  defines the operations;  $final\_stock = initial\_stock - stockIn$ ;  $final\_stock1 = initial\_stock1 + stock1Out$ ;  $final\_stock2 = initial\_stock2 + stock2Out$ . We will consider that we do not have postconditions in the case of the SML model.

$$T_2 = \langle S_2, U_2, Y_2, set_2, get_2, doStep_2 \rangle;$$

In the case of SML:  $S_2 = S_1$ ,  $U_2 = U_1$ ,  $Y_2 = Y_1$  and  $PT_2 = PT_1$ .

The  $doStep_2$  function implements the  $Act^2$  action as follows: The precondition for the execution of the  $Act^2$  action is  $initial\_stock1 \geq stock1In$  AND  $initial\_stock2 \geq stock2In$  AND  $capacity - initial\_stock \geq stockOut$ . The action  $(final\_stock, final\_stock1, final\_stock2) = Act^2(final\_stock, final\_stock1, final\_stock2)$  defines the operations:  $final\_stock1 = initial\_stock1 - initial\_stock1In$ ;  $final\_stock2 = initial\_stock2 - initial\_stock2In$ ;  $final\_stock = initial\_stock + initial\_stockOut$ . Even in the case of this behavioral rule we do not have a postcondition.

For component M we have the inputs and outputs identical to the inputs and outputs of the other two components in the case of the SML model. The  $doStep_M$  function finds the matches in the model and implements the distribution of activities to the  $T_1$  and  $T_2$  components. Of course, for the  $T_1$  and  $T_2$  components there will be several instances, one for each match. In the case of the example in Fig. 2 we have 2 instances of the  $T_1$  component and three instances of the  $T_2$  component. The distinction between the two types of instances is made by the value of the variables `initial_id` and `final_id`. For the co-simulation of the three components we used INTO-CPS. We performed the co-simulation on an example of data and we obtained the output from Fig. 14. The graphs show the stocks resulting from the two instances of the  $T_1$  component and three instances of the  $T_2$  component.

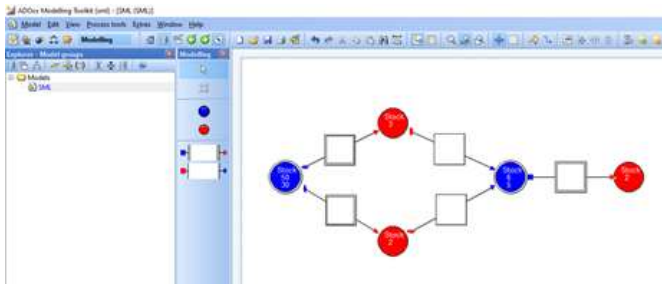


Fig. 12. SML Tool.

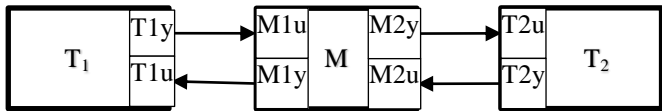


Fig. 13. FMU Components.

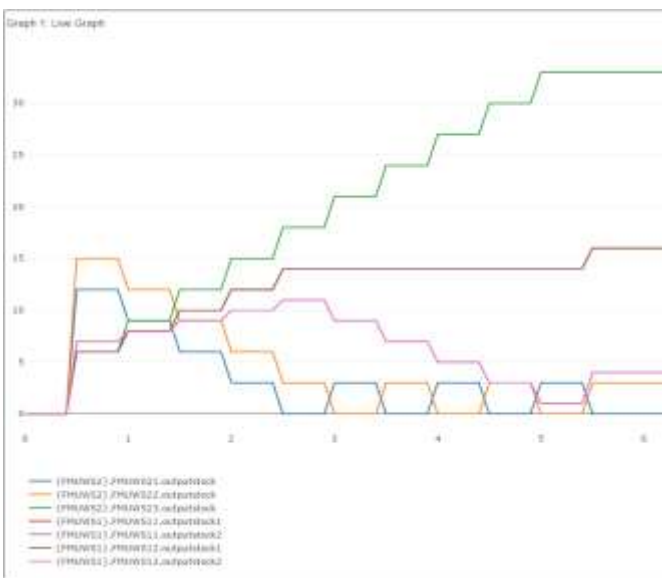


Fig. 14. Screenshot of the Output from INTO-CPS.

## VI. ORIGINAL CONTRIBUTIONS AND CONCLUSIONS

In this paper we used the mechanisms of category theory to specify diagrammatic models with co-simulation facilities. We introduced the concept of behavioral rule as an aggregation between a graph transformation and a behavioral action. We defined these behavioral rules by graph signatures at the metamodel level. We also implemented a simple example of diagrammatic language using the ADOxx [21] metamodeling platform. In all the phases of specification and implementation we highlighted the implementation of the constituent components of such an FMU. We performed the co-simulation of a concrete model specified with the SML language on the INTO-CPS platform.

Model transformations, if any, must be implemented in component M (Fig. 13), otherwise they could not be executed in parallel. As a result, the preconditions and postconditions should also be executed in component M. In principle, this is acceptable.

As it results from the previous observations, there are some important problems to be solved that we will deal with in future work such as: the implementation of a more complex model containing graphical transformations or the implementation of the export facility of a tool-wrapper for DSMLs implemented with ADOxx.

## REFERENCES

- [1] M. Fowler, R. Parsons, Domain Specific Languages, 1st ed. Addison-Wesley Longman, Amsterdam, 2010.
- [2] Uwe Wolter, Zinovy Diskin, The Next Hundred Diagrammatic Specification Techniques, A Gentle Introduction to Generalized Sketches, 02 September 2015 : <https://www.researchgate.net/publication/253963677>.
- [3] D.C. Crăciunean, D. Karagiannis, Categorical Modeling Method of Intelligent Workflow. In: Groza A., Prasath R. (eds) Mining Intelligence and Knowledge Exploration. MIKE Lecture Notes in Computer Science, vol 11308. Springer, Cham (2018).
- [4] D.C. Crăciunean, D. Karagiannis, A categorical model of process co-simulation, Journal of Advanced Computer Science and Applications(IJACSA), 10(2), (2019).
- [5] C. Gomes, C. Thule, D. Broman, P.G. Larsen, H. Vangheluwe - Co-simulation: State of the art, - ACM Computing Surveys, Vol. 1, No. 1, Article 1. Publication date: January (2016).
- [6] Functional Mock-up Interface for Model Exchange and Co-Simulation, Document version: 2.0.1 October 2nd 2019, <https://fmi-standard.org/>.
- [7] INTO-CPS Tool Chain User Manual, Deliverable Number: D4.3a Version: 1.0 Date: December, 2017 Public Document, <http://into-cps.au.dk>.
- [8] D. Karagiannis, H.C. Mayr, J. Mylopoulos, Domain-Specific Conceptual Modeling Concepts, Methods and Tools. Springer International Publishing Switzerland (2016).
- [9] Dominik Bork, Dimitris Karagiannis, Benedikt Pittl, A survey of modeling language specification techniques, Information Systems 87 (2020) 101425, journal homepage: [www.elsevier.com/locate/is](http://www.elsevier.com/locate/is).
- [10] R. Milner, The Space and Motion of Communicating Agents, Cambridge University Press, (2009).
- [11] D.C. Crăciunean, Categorical Grammars for Processes Modeling, International Journal of Advanced Computer Science and Applications(IJACSA), 10(1), (2019).
- [12] Michael Barr And Charles Wells, Category Theory For Computing Science- Reprints in Theory and Applications of Categories, No. 22, 2012.
- [13] Diskin Z., König H., Lawford M., 2018. Multiple Model Synchronization with Multiary Delta Lenses. In: Russo A., Schürr A. (eds) Fundamental Approaches to Software Engineering. FASE 2018. Lecture Notes in Computer Science, vol 10802. Springer, Cham.
- [14] Zinovy Diskin, Tom Maibaum- Category Theory and Model-Driven Engineering: From Formal Semantics to Design Patterns and Beyond, ACCAT 2012.
- [15] D. Plump, 'Checking graph-transformation systems for confluence', ECEASST, vol. 26, 2010. DOI: 10.14279/tuj.eceasst.26.367.
- [16] Hartmut Ehrig, Claudia Ermel, Ulrike Golas, Frank Hermann, Graph and Model Transformation General Framework and Applications, Springer-Verlag Berlin Heidelberg 2015.
- [17] D. Plump, 'Computing by graph transformation: 2018/19', Department of Computer Science, University of York, UK, Lecture Slides, 2019.
- [18] G. Campbell, B. Courtehoue and D. Plump, 'Linear-time graph algorithms in GP2', Department of Computer Science, University of York, UK, Submitted for publication, 2019. [Online]. Available: <https://cdn.gjcampbell.co.uk/2019/Linear-Time-GP2-Preprint.pdf>.
- [19] FMU SDK, <https://github.com/qtronic/fmusdk>.
- [20] Claudio Gomes, Casper Thule, Levi Lucio, Hans Vangheluwe, and Peter Gorm Larsen, Generation of Co-simulation Algorithms Subject to Simulator Contracts, <https://sites.google.com/view/cosimcps19>.
- [21] ADOxx, <https://www.adoxx.org>.

# A Cluster based Non-Linear Regression Framework for Periodic Multi-Stock Trend Prediction on Real Time Stock Market Data

Lakshmana Phaneendra Maguluri<sup>1</sup>, R. Ragupathy<sup>2</sup>

Department of Computer Science and Engineering, Annamalai University  
Annamalai Nagar, Chidambaram, Tamil Nadu 608002, India

**Abstract**—Trend prediction is and has been one of the very important tasks in the stock market since day one. For a sophisticated trend prediction using real time stock market data, stock sentiment news and technical analysis plays a vital role. While predicting the trend in the conventional way, technical indicators are delayed due to temporal data and less historic data. All the conventional stock trend predicting methods sustained without sentiment scores, technical scores and time periods for trend prediction. Considering the fact that all the previous conventional methods of stock trend predictions are bound to take single stock for trend prediction due to high computational memory and time, this prototype of highly functioning algorithms focus on trend prediction with multi stock data breaking all the conventional rules. This multi stock trend prediction model commissions and implements the effectively programmed algorithms on real time stock market data set. In this multi-stock trend prediction model, a new stock technical indicator and new stock sentiment score are proposed in order to improve the stock feature selection for trend prediction. In order to find the best real time feature selection model, a technical feature selection measure and stock news sentiment score are developed and incorporated. We used integrated stock market data to make a hybrid clustered model to find the relational multi stocks. Giving a final verdict, this is a cluster based nonlinear regression multi stock framework in order to predict the time-based trend prediction. The multi stock trend regression accuracy is bettered by 12% and recall by 11% while we cross check the experimental outcomes, henceforth making this model more accurate and precision furnished.

**Keywords**—Multi-stock trend prediction; stock market; clustering; nonlinear regression

## I. INTRODUCTION

Stock markets provide investors with the most profitable avenue to spend their money. The investors can't find another way to make a high rate of return growth from anywhere else. They will have to bear the loss if things go south. This implies that investment is quite a risky thing to be involved. There is every opportunity to make returns bigger and bigger, and to lose everything [1]. When looked into theories, any change in share prices is associated with change in fundamental variable relevant to share price assessment. For example, Stock earnings, size, dividend pay-out ratio, various economic variables etc. A moving average can be calculated using a predetermined period by using the mathematical analysis of the stock's average price value. Every time the price of the stock

changes, the average price either goes up or down [2]. Simple (arithmetic), triangular, exponential, and variable and weighted moving averages, calculated using open, close, low, high and stock price volumes are the different type of moving average indicators. The trend of the price maybe hiked or inflated at times at any point. These common outlines the flow of market trading. The investors gain profit when they buy a stock that shows uptrend. The uptrend stock here is value-appreciating stock. In uptrend, the value of the stock appreciates consistently over a period, even if there are consolidated. For example, a company's particular stock price began to say at ₹ 275, its price reached say ₹ 380, we say there is an uptrend for that stock, even though there were brief dips in stock prices between them. An uptrend can span hours, months, and years [3]. A downtrend is pretty much contrary of an uptrend. In downtrend, the stock price depreciates steadily over a period that may include some brief rises. If he follows these signals the investor can be benefited. Moving average trading is profitable if the price level shifts between selling and buying signals are adequate, otherwise it will result in losses. Bollinger Bands can be a great aid for dealers in dualistic decisions [4]. New openings to the trade can be opened by them. The market will likely get to seize around when the market approaches a Bollinger band. This information alone suffices for a dualistic decision to be won. Bollinger bands need a simple indication of how much they should make an attempt on the markets. Types of dualistic options with high outputs such as hierarchy options or one touch choices need this prediction, which is Bollinger Bands can turn a normal strategy into one that is highly profitable. Bollinger bands form essential levels of conflict and trend prediction. The relative strength index [5] is an indicator of momentum that measures the level of current price changes in order to assess the over-bought or over-sold conditions of a stock price or other quality. The RSI compares the momentum of stock price predictions [6].

When we select the appropriate technique of pre-processing data the sentiment analysis can be improved. This very fact makes pre-processing of data a crucial step in the process. Aside from the usual pre-processing techniques, some of the news articles require different pre-processing techniques because the content produced by the user community, for example, received messages from Twitter [7]. The views expressed in the news are either positive or negative or neutral opinion which plays an important role in the trend prediction. Analyzing sentiment is the task of selecting the sentiment label

for a given news article. It may also be considered a task to classify. As a result of this news reporting on a company, the opinions are formed among the investors, so they can make informed decisions about their share in that company's stock. All inputs are deemed independent from each other to predict the test class labels. The financial media reporters gather the information through reliable sources and the same would be disseminated in news article format [8]. The news articles which are published must be checked for trustworthiness. There are different methods of media where these news articles can be disseminated, and source from which the sentiment is to be derived from the news articles published in that source must be decided with utmost care. A Research Paper's result accuracy is based on the credible sources. Yahoo finance, the money control is few official websites we can say where the news articles are trustworthy. In the text classification approach, each word in the article is weighted with the frequency and this is classified within the specified group. Considering the importance of trading volume in understanding stock market microstructure, comprehensive empirical studies were conducted to research the relation between price, volatility and volume of trading. Market investors are often inclined to look for better investment options providing higher returns as the investment decision is made to gain better returns than other avenues available, or to expect a higher return than others. The probability of not achieving the anticipated or targeted return is commonly known as risk, but risk estimation is a difficult activity. Volatility is usually taken as the indicator of risk. Simple words volatility is a standard deviation in returns. Volatility can be actual volatility, historical volatility, the volatility implied and the volatility forward. Although considering the reasons for volatility, the economists argue that the market is moving according to the information provided to the market; others argue that volatility has little to do with the economic or external factors, and it is the reaction of the investors that exerts greater market impact. Investors are generally averse to risk. At the same time investment with volatile assets has to be made. The investment in security usually has varying purposes. Some buy stock and keep long to have the privilege of owning these capital assets. But some others are buying stock to sell and have the price differences. The return on equities varies with shifts in stock prices. Stock rates rarely remain the same. It is unpredictable. On the one hand, price volatility [26] is an opportunity, on the other, a threat to the investors. Price stability will reduce the risk stemming from price volatility. Yet stock prices cannot stay steady over time, because they are more prone to shifts in environmental factors. There are no bounds or barriers to the flow of funds in a globalized environment. The major players now in the Indian Stock Market are the FII (Foreign institutional investors). With the incurring losses, the risk of the stock increases, this is in fact measured by standard deviation statistic. This is the dispersion from what is required of the real. The larger the dispersion the greater the perceived security risk will be. The risk of a stock is viewed in relation to the market as well. Each safety is susceptible to market influence. Market influence may be greater or smaller. But the fact is, to a greater extent, the fortune of the individual stock is governed by the market. This part of the risk to the stock is called the systematic risk. All stocks on the market must share

that class of risk. Such risks are therefore also known as non-diversifiable risk, because they cannot be eliminated through diversification. The statistics used to measure this portion of overall risk are beta. A security beta tells how far the security is market related. Operation on the stock market became popular nowadays. Present investors don't find investing in the stock market as pointless. They find investment in stocks to be more remunerative than other opportunities. Formerly stock investment has not received due respect and it has been treated as somewhat speculative that even today some discounts for its social acceptability are considerable. Stocks give not just the institutional investors but also the small retail investors a better opportunity. But that doesn't mean everyone knows the surgery. Market operation transparency is still in jeopardy. The SEBI is trying hard to get things working.

Looking back at our previous contributions, we have developed a single stock trend prediction using the technical and news data in an intraday process. Now, we propose a single stock trend prediction model using the technical and news data in different periodic time intervals. In this contribution an advanced multi-stock trend prediction model is designed and implemented on real time stock market data in different periodic time intervals. In this paper, new multi-stock technical and sentimental scores are developed to improve the stock selection process. A multi-stock clustering algorithm and classification models are developed in order to predict the periodic multi-stock trend.

## II. RELATED WORK

Jeon et al. [9] demonstrated behavior next day by using a random subsample of collected tweets for stock market. They've gathered the NASDAQ, S&P 500 and DJIA tweet posts. For each day, they considered the factor of combined fear and hope, and analyzed the relationship between market indicators and these factors. They reported that the above mentioned stocks had been negatively associated with emotional tweets. Their findings have proved that stock market reaction on the very next day can be predicted by collecting emotional data [10].

Vu [11] proposed a new machine learning system by integrating features, consumer assurance and last 3 days data into the products. The cross-validation method has been adopted in a Decision Tree classifier for integrating all of the filtered features. Pre-processing steps include extracting noisy data, normalizing tweets and selecting data. The model was tested with NER (Named Entity Recognition Task) and without NER for Google, Apple, Amazon, and Microsoft companies stock and yielded 80.49 percent, 82.93 percent, 75.00 percent, and 75.61 percent for up and down NER (Named Entity Recognition Task) labels, respectively [11].

Vijh et al. [12] created a thorough study of stock prediction from data collection (how to collect it from twitter and tweet description), cloud storage, and then the process of opinion analysis (software and techniques) and finally the phase of prediction. Over time they examined the correlation between financial markets and social media data. They built a cloud-based system in JSON format to store various dimensions of public emotions contained in fetched tweets. The program was assessed for four companies listed under the UK Stock

Exchange, and the data checked were collected for 30 days. Their finding will help the firms assess the concerns of stakeholders and establish a new market strategy. Overall, the research enhanced the efficiency in the forecasting phase with emotional analysis and synthesizing.

Zhang et al. [13] used the twitter and survey index sentiments and attentive indicators, volatility and trading volume of S&P index 500 to forecast returns. Various supervised learning techniques and the Diebold-Mariano test were conducted and compared with autoregressive baseline model to confirm the significance of sentiments and attention-based predictions. They noted that tweet volume and sentiments were relevant to predicting lower-market capitalization portfolios. In addition, they show that Kalman Filter indicators and Twitter sentiment were helpful in forecasting some sentiment labels based on surveys.

Chen et al. [14] analyzed the data obtained from various networking networks called chat rooms, web forums, and micro blogs and found different characteristics to be present. They believed that chat room posts at the activity level are strongly correlated with the trend in stock and assumption that is true. For chat room post sentiments the same performance was achieved with short posts reported from previous studies. The result indicated that post sentiments improved stock price return forecasting as compared to using only historical prices. They also developed a trading strategy and reported a return of 21 per cent over seven months. Proposed Model The overall process of predicting stock market direction consists of different steps that include data collection, pre-processing of text and selection of features. The programming is required for the overall work to be carried out. R language, python, was used with Java. The packages of those tools have been used to implement the algorithms proposed. The data our problem requires are of two types [15]. The historic stock values and the news stories from which the emotions are to be derived. Unlike the other systems which used the static data, our system is based on both the streaming data and the static data. The crawler crawls on the specified website and extracts the specified company's news articles for which the future direction of the stock is to be predicted. Since the stock prices must be correlated with the news articles, the news articles must be extracted along with the time stamps. Such news articles then act as the input to the module for the study of sentiments. The researchers have been pursuing paths of sentiment analysis for many years, and have come up with many different algorithms to characterize the text's feeling. Every algorithm has some advantages and disadvantages. Choosing the algorithm for sentiment analysis [25] may depend on the available datasets, domain and prior experience. One approach is to be chosen among approaches, linguistic-based, lexicon-based, and machine learning. If the approach is selected, the correct algorithm must be determined in that approach [16]. It is very crucial to decide what data set is being used for the research. Our framework does not dispose of readily accessible data sets. Most of the data is data processing that is being processed and stored in the database. Our system has to have two types of data. One relates to historical stock values and the other set of data containing news articles which are published online. The data used is from a combination of

two different sources to study the correlation between news articles and stock prices: a dataset of historical data and a corpus of news articles. The initial source of data used to extract news articles is the website of money control, which has a large reservoir of critical news for the individual stocks. Money control is India's premier source of financial information. They derived historical values for 2012 from <http://ichart.finance.yahoo.com> for the Infosys stock in NIFTY. This data is then loaded into a table of databases that can then be queried and processed. The moneycontrol.com Website was used for the news articles. To predict future prices based on the sequence of events, historical data are extracted from the moneycontrol.com web site for all the companies listed in BSE (around 3000) [17]. The code (scraper) is written to extract the open price of each company, close prices for the years 2007 to 2014. For this system, events from disclosure records and pieces of content collected from [indiatimes.com](http://indiatimes.com), [moneycontrol.com](http://moneycontrol.com), [sebi.com](http://sebi.com), [watchoutinvestors.com](http://watchoutinvestors.com), [courts.gov.in](http://courts.gov.in), [cibil.com](http://cibil.com) are used as corpus.

Data pre-processing greatly decreases word space but there are still incentives for knowledge loss. Kang [18] measured the volatility of Indian stock market day-to-day returns. The study period was 1961-2005, and data were gathered together from the Economic Times Index and S&P CNX Nifty. The series observed volatility clustering quiet intervals of big returns were interspersed with cycles of volatility of great returns. The GARCH model was used to check the volatility effect asymmetry, and the result indicated a volatility asymmetry. It was known that high price movements started in response to strong economic fundamentals, and that the real reason for sudden movement was market imperfection [19].

Smruti et al. [20] proposed an extreme learning based PCA approach to predict the stock market data on limited training dataset. Shangkun et al. [21] proposed a gradient boosting approach to detect the trend in the china market. Shanoli et al. [22], proposed a novel time series model to predict the stock market data using the rule based approach on the training data.

### III. FILTER BASEDD STOCK TECHNICAL PREDICTION MODEL

In the paper [23], we have proposed a novel filtered based classification model on the technical and stock news datasets in order to predict the trend of the to find the bullish trend stocks on the real-time market data. This model is tested on the continuous type of technical data for trend prediction. In the proposed framework, a correlated multi-stock trend prediction model is designed and implemented on the real-time market data. In the initial phase, a real-time stock technical data and its related news are extracted from the money control and zerodha websites. These technical data and news data are pre-processed using the novel approaches developed in the papers [24]. In this work, an improved version of technical indicator and sentiment scores are defined based on the contextual information of the stock data. These scores and technical data are integrated to form the training data. A novel clustering measure is used to form the clusters based on the integrated data features. This clustering model is implemented to form the clusters based on the technical data and scores of an integrated dataset. Finally, this clustered data is given to classification

model to predict the trend of the multiple stocks based on the input test sample as shown in the Fig. 1.

**A. Stock Technical and Comments Data Collection and Pre-Processing**

In this phase, all the stock related technical and news data are extracted from the zerodha or trade view or money control websites for data collection. All these collected data are pre-processed using the models in the papers [20][21]. Text pre-processor is applied on the stock news data as text filtering. In this work, a modified version of stock technical score and sentiment score are developed on the technical and stock news datasets.

**B. Hybrid Stock News Score**

To each comment in the stock corpus S, we construct a dictionary of words that contains bullish and bearish words. In the stock news training data S, each input stock news is represented as sn[i] and term frequencies of the sn[i] is presented as tsn[i][j], where i, j represent the j<sup>th</sup> term of the ith stock news s. Here, the term frequency and normalized term frequencies are used to find the news score of the stock. T represents the total tokens in the i<sup>th</sup> stock news. This normalized data is scaled by using inverse document frequency (idf) and multi-stock scaling factor (mssf) is represented in Eq. (1)

$$tfidf(sn(i), S, mssf) = tf(tsn(i), S) \times idf(tsn(i), mssf)$$

$$idf(sn(i), S, mssf) = \log \frac{mssf}{1 + \max\{tsn[i][j]\}; j = 1..T.}$$

$$mssf = \max\{tf(S[i], tf(S[j]))\} / N * Prob(tf(S[i]/tf(S[j])); i \neq j \tag{1}$$

**C. Proposed New Stock Technical Indicator**

In the paper, a novel mutual information (MI) is proposed to find the contextual relationship of the bullish and bearish stocks using the technical indicators. Hybrid technical mutual information is represented in terms of bullish and bearish cases as shown in Eq. (2).

$$Technical\ Indicator = TI = stgv(l, bu) \log \frac{strv(l, be)}{rsi(bu) * rsi(be)} \tag{2}$$

Where bu represents the bullish and be represents the bearish stock type.

$$T1 = P(\beta) * CP_s - MA \left( \frac{n}{2} + 1 \right)$$

$$Bu: T1 \geq 0; n: time\ interval$$

$$Be: T1 < 0$$

$$T2 = P(\beta) * \frac{S_{ma(3)} + S_{ma(6)} + S_{ma(12)} + S_{ma(24)}}{4};$$

$$Bu: CP_s \geq T2$$

$$Be: CP_s < T2$$

$$Final\ Multi-Stock\ technical\ indicator = MT = (T1 + T2) / TI$$

where, S<sub>ma(i)</sub> is the ith stock moving average. Strv: Super trend red value, Stgv: Super trend green value. MA: moving average. CPs: Closed price of a stock s.

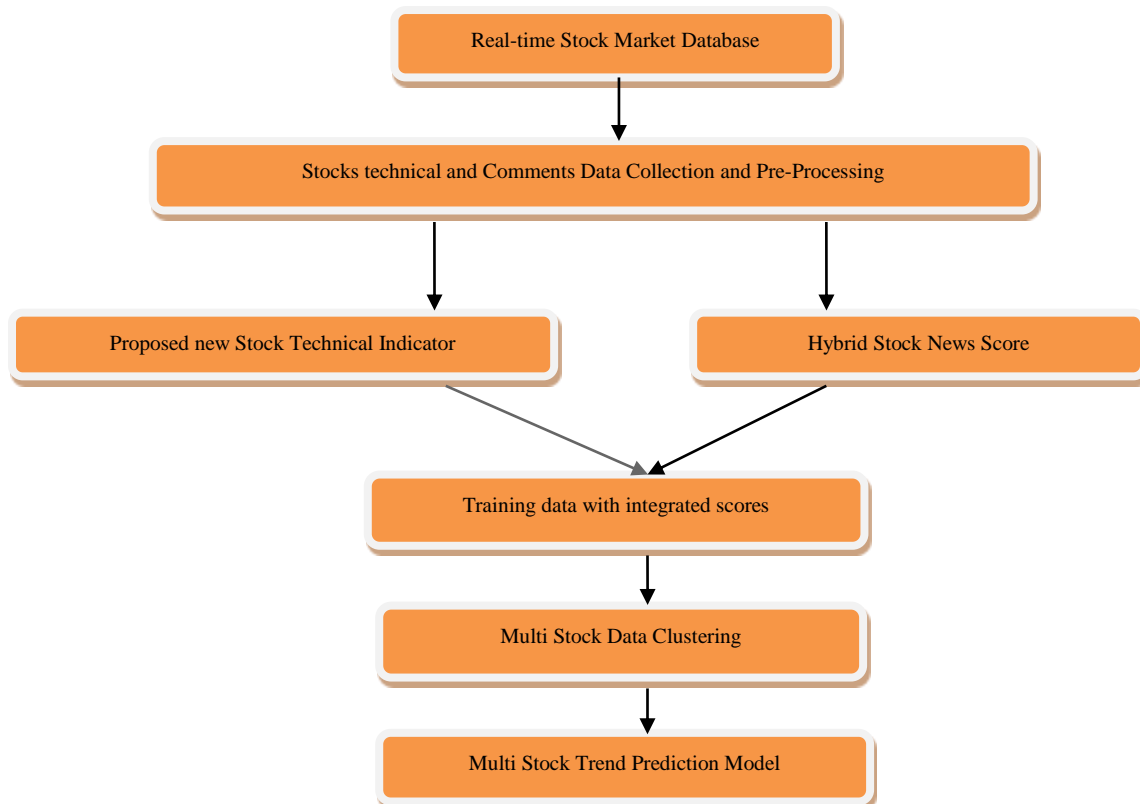


Fig. 1. Proposed Multi Stock Trend Prediction Framework.

#### D. Bullish Dictionary Words

Escalated, gain, enjoy, expansion, aggrandize, elevated, increment, rise, prefer, hallow, expand, supersize, idolize, positive, appreciate, plus, relish, accelerate, augment, raise, more, amplify, soar, adore, appreciative, approbatory, desire, esteem, approving, raised, swell, extend, addition, worship, climb, add, commendatory, venerate, augmentation, fancy, revere, friendly, proliferate, addendum, increased, escalate, proliferation, accumulate, love, stoke, complimentary, heightened, hype, uprising, accrual, boost, up, applauding, enlarge, admire, admiring, good, multiply, accretion.

#### E. Bearish Dictionary Words

Descend, recede, depreciative, abhor, drop, diminution, depreciatory, uncomplimentary, adverse, deplore, slide, detest, lower, plunge, lessen, unappreciative, depletion, dislike, dive, reduce, decrease, depressed, decreased, under, diminish, dip, low, derogatory, disapprove, unfavorable, negative, lowering, loathe, disfavor, unflattering, sink, receded, disdain, hate, decrement, unfriendly, subtract, loss, abate, decline, despise, fall, diminishment, lessening, downsize, abominate, minify, execrate, deprecate, inappreciative, dropped, shrinkage, reduction, wane, abatement, disapproving, dwindle, down.

#### F. Training Data with Integrated Scores

In this work, technical data is integrated with the newly computed technical score and the stocks news score as unlabeled data for data clustering. Here, multiple stock's technical data and scores are integrated to find the score-based data clustering on multiple stocks.

#### G. Multi-Stock Data Clustering

**Input parameters: SF: stock features, SC: stock clusters**

1: for each feature  $f_i$   $i = 1, \dots, N$  in the stock feature space SF

2: for each instance all points do

3: for all stock clusters  $SC_j$   $j = 1$  to  $k$

4: do

5: Compute

$$C_i = \frac{\frac{1}{n} \sum_{i=1}^n (f_k - \mu_f)^3}{\sigma_k^3} \parallel \max\{f_k\} - \mu_f \parallel^2$$

$$6. \text{ Find } C_i^*(f_n) = \arg \min_i \left( \frac{\frac{1}{n} \sum_{i=1}^n (f_k - \bar{X})^3}{\sigma_k^3} \parallel f_k - \mu_f \parallel^2 \right)$$

7: Update multi-stock cluster  $C_i = \{f_n \mid C_i^*(f_n) = i\}$

8: repeat until  $k$  clusters

9: end for //number of features

10: end for //number of clusters

#### H. Multi-Stock Trend Prediction Model

In the proposed multi-stock trend prediction model, a hybrid multi-linear regression model is designed and implemented to predict the trend of the multiple stocks. In this model, a new probability estimation based non-linear regression model is designed and implemented on the training clustered dataset. A Non-linear regression estimation using the time wise trend prediction is given as.

#### Parameters:

$f(s(\alpha))$ : input query stock.

$f(s(\beta))$ : Related list of multiple stocks.

$r_{s(\alpha)}, r_{s(\beta)}$ : RSI values of the input and multiple stocks.

$\varphi_{s(\alpha)}, \phi_{s(\beta)}$ : BB value of input and its related multiple stocks.

1. Read the clustered training dataset.
2. Compute the non-linear regression based estimation to the time based trend estimation as

Let  $m_{s(\alpha)}, m_{s(\beta)}$  are the test stock and its associated multi-stocks

$$f(m_s) \sim E(m(s), k(m_{s(\alpha)}, m_{s(\beta)}))$$

where

$$v = \text{Variance} = V(\text{RP}(\varphi_{s(\alpha)}, \gamma_{s(\beta)})) = \frac{\varphi_{s(\alpha)} \cdot \phi_{s(\beta)}}{(\varphi_{s(\alpha)} + \phi_{s(\beta)})^2 (\phi_{s(\beta)} + \gamma + 1)}$$

$$3. \sigma = \text{SD} = \text{SD}(\text{RP}(\varphi_{s(\alpha)}, \phi_{s(\beta)})) = \sqrt{\frac{\varphi_{s(\alpha)} \cdot \phi_{s(\beta)}}{(\varphi_{s(\alpha)} + \phi_{s(\beta)})^2 (\varphi_{s(\alpha)} + \phi_{s(\beta)} + 1)}}$$

$ED(\chi) := f(\chi) = k(r_{s(\alpha)}, r_{s(\beta)}) \cdot \chi e^{-\gamma \chi} / \sigma$  for  $x \geq 0$ ; // rsi based test stock and its correlated stocks data distribution estimation.

$$m(s) = ED[f(s(\alpha)), f(s(\beta))]$$

$$\text{Pr ob}(r_{s(\alpha)}, r_{s(\beta)}) = k(r_{s(\alpha)}, r_{s(\beta)}) / v$$



$$\begin{aligned}
 k(r_{s(\alpha)}, r_{s(\beta)}) &:= \frac{r_{s(\beta)}}{\pi[(x - r_{s(\alpha)})^2 + r_{s(\beta)}^2]} \\
 \text{ProbEsti} &= D\left(\frac{r_{s(\beta)}}{\pi[(x - r_{s(\alpha)})^2 + r_{s(\beta)}^2]} * \chi * e^{(-\chi x)} / \left(\sqrt{\frac{\varphi_{s(\alpha)} * \phi_{s(\beta)}}{(\varphi_{s(\alpha)} + \phi_{s(\beta)})^2 * (\varphi_{s(\alpha)} + \phi_{s(\beta)} + 1)}}\right)\right) \\
 &= \frac{\partial}{\partial x} \left(\frac{r_{s(\beta)}}{\pi[(x - r_{s(\alpha)})^2 + r_{s(\beta)}^2]} * \chi * e^{(-\chi x)} / \left(\sqrt{\frac{\varphi_{s(\alpha)} * \phi_{s(\beta)}}{(\varphi_{s(\alpha)} + \phi_{s(\beta)})^2 * (\varphi_{s(\alpha)} + \phi_{s(\beta)} + 1)}}\right)\right) \\
 &= \left(\frac{r_{s(\beta)}}{\pi[(x - r_{s(\alpha)})^2 + r_{s(\beta)}^2]}\right) e^{-\chi x} \frac{\partial}{\partial x} (-\chi x) \cdot \sqrt{\frac{\varphi_{s(\alpha)} * \phi_{s(\beta)}}{(\varphi_{s(\alpha)} + \phi_{s(\beta)})^2 * (\varphi_{s(\alpha)} + \phi_{s(\beta)} + 1)}} \\
 &= \left(\frac{r_{s(\beta)}}{\pi[(x - r_{s(\alpha)})^2 + r_{s(\beta)}^2]}\right) e^{-\chi x} - \chi \cdot \sqrt{\frac{\varphi_{s(\alpha)} * \phi_{s(\beta)}}{(\varphi_{s(\alpha)} + \phi_{s(\beta)})^2 * (\varphi_{s(\alpha)} + \phi_{s(\beta)} + 1)}} \\
 &= -\chi \cdot e^{-\chi x} \left(\frac{r_{s(\beta)}}{\pi[(x - r_{s(\alpha)})^2 + r_{s(\beta)}^2]}\right) \cdot \left(\sqrt{\frac{\varphi_{s(\alpha)} * \phi_{s(\beta)}}{(\varphi_{s(\alpha)} + \phi_{s(\beta)})^2 * (\varphi_{s(\alpha)} + \phi_{s(\beta)} + 1)}}\right)
 \end{aligned}$$

Estimate non-linear least square equation for

$$\begin{aligned}
 \text{Trend } T^* &= \text{Max} \left\{ -\chi \cdot e^{-\chi x} \left(\frac{r_{s(\beta)}}{\pi[(x - r_{s(\alpha)})^2 + r_{s(\beta)}^2]}\right) \cdot \left(\sqrt{\frac{\varphi * \phi}{(\varphi + \phi)^2 * (\varphi + \phi + 1)}}\right) * \max\{b(s(\alpha)), b(s(\beta))\} \right\} \\
 b(s(\alpha)) &= \frac{\sum_{i=1}^n (p(s(\alpha))_i - \frac{\sum_{i=1}^n p(s(\alpha))_i}{n}) (q(s(\beta))_i - \frac{\sum_{i=1}^n q(s(\beta))_i}{n})}{\sum_{i=1}^n (p(s(\alpha))_i - \frac{\sum_{i=1}^n p(s(\alpha))_i}{n})^2} \\
 b(s(\beta)) &= \frac{1}{n} \left( \sum_{i=1}^n q(s(\beta))_i - b(s(\alpha)) \cdot \sum_{i=1}^n p(s(\alpha))_i \right)
 \end{aligned}$$

4. If ( $T^* \geq 1M(MBB) \& T^* > 1M(MACDS)$ )
5. then
6. stock s has Trend Positive.
7. Get stocks S\* with similar T\* value which satisfies the given condition on 1 minute time frame.
8. Else
9. stock s has Trend Negative.
10. Get stocks S\* with similar T\* value which satisfies the given condition on 1 minute time frame.
11. If ( $T^* \geq 3M(MBB) \& T^* > 3M(MACDS)$ )
12. then
13. stock s has Trend Positive.
14. Get stocks S\* with similar T\* value which satisfies the given condition on 3 minute time frame.
15. Else
16. stock s has Trend Negative.
17. Get stocks S\* with similar T\* value which satisfies the given condition on 3 minute time frame.
18. If ( $T^* \geq 5M(MBB) \& T^* > 5M(MACDS)$ )
19. then
20. stock s has Trend Positive.
21. Get stocks S\* with similar T\* value which satisfies the given condition on 5 minute time frame.
22. Else
23. stock s has Trend Negative.
24. Get stocks S\* with similar T\* value which satisfies the given condition on 5 minute time frame.

25. If  $(T^* \geq 10M(MBB) \& \& T^* > 10M(MACDS))$
  26. then
  27. stock  $s$  has Trend Positive.
  28. Get stocks  $S^*$  with similar  $T^*$  value which satisfies the given condition on 10 minute time frame.
  29. Else
  30. stock  $s$  has Trend Negative.
  31. Get stocks  $S^*$  with similar  $T^*$  value which satisfies the given condition on 10 minute time frame.
  32. If  $(T^* \geq 15M(MBB) \& \& T^* > 15M(MACDS))$
  33. then
  34. stock  $s$  has Trend Positive.
  35. Get stocks  $S^*$  with similar  $T^*$  value which satisfies the given condition on 15 minute time frame.
  36. Else
  37. stock  $s$  has Trend Negative.
  38. Get stocks  $S^*$  with similar  $T^*$  value which satisfies the given condition on 15 minute time frame.
  39. If  $(T^* \geq 30M(MBB) \& \& T^* > 30M(MACDS))$
  40. then
  41. stock  $s$  has Trend Positive.
  42. Get stocks  $S^*$  with similar  $T^*$  value which satisfies the given condition on 30 minute time frame.
  43. Else
  44. stock  $s$  has Trend Negative.
  45. Get stocks  $S^*$  with similar  $T^*$  value which satisfies the given condition on 30 minute time frame.
  46. If  $(T^* \geq 1H(MBB) \& \& T^* > 1H(MACDS))$
  47. then
  48. stock  $s$  has Trend Positive.
  49. Get stocks  $S^*$  with similar  $T^*$  value which satisfies the given condition on 1 hour time frame.
  50. Else
  51. stock  $s$  has Trend Negative.
  52. Get stocks  $S^*$  with similar  $T^*$  value which satisfies the given condition on 1 hour time frame.
  53. If  $(T^* \geq 1D(MBB) \& \& T^* > 1D(MACDS))$
  54. then
  55. stock  $s$  has Trend Positive.
  56. Get stocks  $S^*$  with similar  $T^*$  value which satisfies the given condition on 1 day time frame.
  57. Else
  58. stock  $s$  has Trend Negative.
  59. Get stocks  $S^*$  with similar  $T^*$  value which satisfies the given condition on 1 day time frame.
  60. If  $(T^* \geq 1W(MBB) \& \& T^* > 1W(MACDS))$
  61. then
  62. stock  $s$  has Trend Positive.
  63. Get stocks  $S^*$  with similar  $T^*$  value which satisfies the given condition on 1 week time frame.
  64. Else
  65. stock  $s$  has Trend Negative.
  66. Get stocks  $S^*$  with similar  $T^*$  value which satisfies the given condition on 1 week time frame.
  67. If  $(T^* \geq 1Mt(MBB) \& \& T^* > 1Mt(MACDS))$
  68. then
  69. stock  $s$  has Trend Positive.
  70. Get stocks  $S^*$  with similar  $T^*$  value which satisfies the given condition on 1 month time frame.
  71. Else
  72. stock  $s$  has Trend Negative.
  73. Get stocks  $S^*$  with similar  $T^*$  value which satisfies the given condition on 1 month time frame.
-

In the proposed multi-stock trend prediction algorithm, a non-linear regression model is used to predict the trend of the input stock with different time frames. In this work, we have used 1m,3m,5m,10m,15m,30m,1H,1D,1W,1Mt time frames in order to predict the trend of the given stock based on the clustered stock market dataset. Here, the computed non-linear regression estimator value is tested against the MACD signal value and middle line Bollinger line values to predict the similar type of trends in the real-time market.

IV. EXPERIMENTAL RESULTS

Experimental results are simulated using java environment and real-time market data. Proposed model is compared to the traditional stock market classification models to verify the performance of the hybrid feature selection-based clustering and classification model to the traditional models. Also, proposed model is compared to the traditional techniques by using various statistical performance measures such as accuracy, true positive rate, recall, precision, false positive rate, runtime etc. These performance metrics are analyzed and compared by using third party java libraries. Different types of statistical metrics such as recall, precision, accuracy, F-measure are evaluated on the stock market sentiment data along with the technical data. These statistical measures are evaluated based on the confusion matrix as described in Table I.

Accuracy: It is the ratio of correctly labelled stock predictions class labels to the entire stock class labels as shown in Eq. (3):

$$Stok\ Accurecy(SA) = \frac{(STP+STN)}{(STP+SFP+SFN+STN)} \quad (3)$$

Precision: It is the ratio of correctly classified positive stock classes to the all actual positive and negative labelled stock classes as shown in Eq. (4):

$$Stock\ Precision(SP) = \frac{STP}{(STP+SFP)} \quad (4)$$

Recall: It is the ratio of correctly classified positive stock classes' labels to the all predicted positive and negative labelled stock classes as shown in Eq. (5):

$$Stock\ Recall(SR) = \frac{STP}{(STP+SFN)} \quad (5)$$

F-Measure: It is the harmonic average of recall and precision as shown in Eq. (6):

$$Stock\ F - Measure(SF) = \frac{2*SR*SP}{SR+SP} \quad (6)$$

TABLE I. STOCK STATISTICAL MEASURES

| Model Predicted values |                |                                      |                                      |
|------------------------|----------------|--------------------------------------|--------------------------------------|
| Stock positive         |                | Stock negative                       |                                      |
| Actual stock values    | Stock positive | Stock true positive prediction(STP)  | Stock false positive prediction(SFP) |
|                        | Stock negative | Stock false negative prediction(SFN) | Stock true negative prediction(STN)  |

Fig. 2 illustrates the 5 min candlestick pattern graph in the ZERODHA brokerage website. As shown in the Fig. 2, it is noted that the reliance industries stock is uptrend in the afternoon session.

Fig. 3 illustrates the 10 min candlestick pattern graph in the ZERODHA brokerage website. As shown in the Fig. 3, it is noted that the reliance industries stock is downtrend in the morning session and slightly uptrend in the afternoon session.

Fig. 4 illustrates the 15 min candlestick pattern graph in the ZERODHA brokerage website. As shown in the Fig. 4, it is noted that the reliance industries stock is downtrend in the morning session and slightly uptrend in the afternoon session.

Fig. 5 illustrates the 5 min candlestick pattern graph in the ZERODHA brokerage website. As shown in the Fig. 5, it is noted that the HDFC bank stock has uptrend in the afternoon session.



Fig. 2. Zerodha: 5 Min Reliance Industries Candlestick Chart.



Fig. 3. Zerodha: 10 Min Reliance Industries Candlestick Chart.



Fig. 4. Zerodha: 15 Min Reliance Industries Candlestick Chart.



Fig. 5. Zerodha: 5 Min HDFC Bank Candlestick Chart.

Fig. 6 illustrates the 10 min candlestick pattern graph in the ZERODHA brokerage website. As shown in the Fig. 6, it is noted that the HDFC bank stock is downtrend in the morning session and slightly uptrend in the afternoon session.

Fig. 7 illustrates the 15 min candlestick pattern graph in the ZERODHA brokerage website. As shown in the Fig. 7, it is noted that the HDFC bank stock is downtrend in the morning session and slightly uptrend in the afternoon session.

Fig. 8 illustrates the 5 min candlestick pattern graph in the ZERODHA brokerage website. As shown in the Fig. 8, it is noted that the nifty index is uptrend in the entire session.

Fig. 9 illustrates the 10 min candlestick pattern graph in the ZERODHA brokerage website. As shown in the Fig. 9, it is noted that the Nifty index is downtrend in the morning session and slightly uptrend in the afternoon session.

Fig. 10 illustrates the 15 min candlestick pattern graph in the ZERODHA brokerage website. As shown in the Fig. 10, it is noted that the Nifty index is downtrend in the morning session and slightly uptrend in the afternoon session.



Fig. 10. Zerodha: 15 Min NIFTY Index Candlestick Chart.

Table II describes the performance of computational runtime (ms) of stock trend feature extraction using the proposed approach on large datasets. From the Table II, it is clearly shown that the present feature extraction procedure has low computation runtime as compared to the conventional approaches.

Table III illustrates the proposed multi-stock feature selection measures on the input data. From the Table III, it is observed that the proposed multi-stock feature selection has better filtering than the conventional feature selection measures.

Table IV describes the historical data of the multiple stocks from the real-time market. From the Table IV, periodic levels of stock such as stock volume, stock MACD, stock volume change, low and high details are taken from the real-time market.



Fig. 6. Zerodha: 10 Min HDFC Bank Candlestick Chart.



Fig. 7. Zerodha: 15 Min HDFC Bank Candlestick Chart.



Fig. 8. Zerodha: 5 Min NIFTY Index Candlestick Chart.



Fig. 9. Zerodha: 10 Min NIFTY Index Candlestick Chart.

TABLE II. RUNTIME OF THE STOCK FEATURE COMPUTATION MEASURES ON THE REAL-TIME DATASET

| Stocks | Roughset | PCA  | PSO  | MultiStockFS |
|--------|----------|------|------|--------------|
| #1     | 6238     | 6637 | 6030 | 3777         |
| #2     | 6868     | 6354 | 6086 | 3682         |
| #3     | 7780     | 7637 | 6301 | 3858         |
| #4     | 6321     | 5105 | 7771 | 3679         |
| #5     | 6641     | 6633 | 6813 | 3995         |
| #6     | 6290     | 6609 | 5357 | 3756         |
| #7     | 6362     | 7051 | 7573 | 4742         |
| #8     | 6939     | 5876 | 7640 | 4292         |
| #9     | 5349     | 6132 | 5701 | 4231         |
| #10    | 6942     | 6153 | 7816 | 4780         |
| #11    | 6860     | 5317 | 5191 | 3575         |
| #12    | 6552     | 6384 | 5442 | 3722         |
| #13    | 7709     | 7463 | 6909 | 4518         |
| #14    | 7457     | 5957 | 6097 | 3666         |
| #15    | 7755     | 6264 | 5348 | 3688         |
| #16    | 6157     | 6746 | 5625 | 3659         |
| #17    | 5796     | 5386 | 5113 | 4546         |
| #18    | 6833     | 7465 | 6554 | 3369         |
| #19    | 5379     | 5538 | 5078 | 4593         |
| #20    | 6142     | 6056 | 6448 | 3358         |

TABLE III. STOCK TREND FEATURES EXTRACTION USING THE PROPOSED MODEL

| Stocks | Roughset | PCA | PSO | MultiStockFS |
|--------|----------|-----|-----|--------------|
| #1     | 62       | 66  | 58  | 47           |
| #2     | 58       | 67  | 64  | 55           |
| #3     | 65       | 65  | 62  | 52           |
| #4     | 63       | 61  | 60  | 54           |
| #5     | 67       | 66  | 62  | 48           |
| #6     | 67       | 59  | 63  | 53           |
| #7     | 66       | 63  | 63  | 53           |
| #8     | 64       | 67  | 67  | 48           |
| #9     | 59       | 63  | 60  | 50           |
| #10    | 64       | 64  | 64  | 50           |
| #11    | 67       | 61  | 62  | 53           |
| #12    | 62       | 61  | 64  | 55           |
| #13    | 65       | 64  | 67  | 54           |
| #14    | 62       | 61  | 58  | 48           |
| #15    | 68       | 67  | 65  | 48           |
| #16    | 65       | 63  | 66  | 49           |
| #17    | 66       | 61  | 62  | 54           |
| #18    | 68       | 60  | 64  | 48           |
| #19    | 67       | 68  | 62  | 53           |
| #20    | 68       | 61  | 64  | 54           |

Table V describes the historical data of the multi-stock technical data with MACD less than zero for clustering and classification. From the Table V, periodic levels of stock such as stock high, stock low, stock change, and stock volume, stock MACD are taken from the real-time market.

Table VI describes the performance of F1-measure of multi-stock trend classification using the proposed framework on large datasets. From the Table VI, it is clearly shown that the present framework has better efficiency F1-measure as compared to the conventional approaches.

Fig. 11 describes the performance of recall of multi-stock trend classification using the proposed learning framework on large datasets. As shown in the Fig. 11, it is clearly shown that the present framework has better efficiency recall as compared to the conventional frameworks.

Table VII describes the performance of precision of multi-stock trend classification using the proposed framework on large datasets. From the Table VII, it is clearly shown that the present framework has better efficiency precision measure as compared to the conventional approaches.

Fig. 12 describes the performance of accuracy of stock trend classification using the proposed multi-stock trend prediction framework on large datasets. As shown in the Fig. 12, it is clearly shown that the present framework has better efficiency accuracy as compared to the conventional frameworks.

TABLE IV. HISTORICAL SAMPLE MULTIPLE STOCKS AND ITS TECHNICAL DATA IN THE TRAINING DATA

| Symbol      | % change | price   | volume  | High    | low    | MACD(>0) |
|-------------|----------|---------|---------|---------|--------|----------|
| UNICHEMLAB  | 19.99    | 282.1   | 1284742 | 282.1   | 239.1  | 17.86    |
| AMBER       | 16.69    | 1730    | 465661  | 1749.95 | 1490   | 25.69    |
| VIMTALABS   | 14.63    | 119.5   | 1574820 | 122.85  | 101.55 | 5.36     |
| KREBSBIO    | 12.74    | 101.75  | 178157  | 105     | 85.65  | 2.55     |
| SOLARA      | 11.2     | 846.2   | 989205  | 897     | 778.3  | 31.99    |
| MASTEK      | 10.95    | 635     | 942715  | 649.5   | 573.9  | 52.34    |
| PGEL        | 9.95     | 45.85   | 106109  | 45.85   | 43.15  | -0.28    |
| KOPRAN      | 9.92     | 52.65   | 908150  | 52.65   | 47.4   | 3.55     |
| ONMOBILE    | 9.91     | 32.15   | 255665  | 32.15   | 28.85  | 0.437    |
| JUBILANT    | 9.81     | 873     | 1649341 | 895     | 791    | 45.33    |
| SMSPHARMA   | 9.71     | 82.45   | 2135488 | 83.9    | 74     | 5        |
| RATNAMANI   | 9.56     | 1144.95 | 40209   | 1192.85 | 1043.2 | 16.2     |
| NEOGEN      | 9.39     | 590.15  | 109135  | 593.45  | 538.35 | 12.35    |
| BLISSGVS    | 8.9      | 116.85  | 1111067 | 118     | 106.85 | 2.4      |
| CAPLIPPOINT | 8.31     | 467.85  | 922628  | 472     | 435    | 21.79    |
| NEULANLAB   | 8.2      | 845.1   | 449593  | 850     | 765    | 63.7     |
| ERIS        | 8.09     | 531.05  | 778165  | 550     | 481.05 | 4.84     |
| INFOBEAN    | 8.06     | 117.35  | 75037   | 119.3   | 107.4  | 3.01     |
| WOCKPHARMA  | 8.01     | 298.1   | 2887437 | 301.9   | 276    | 3.49     |
| SMSLIFE     | 7.97     | 362.9   | 63598   | 368     | 331    | 10.01    |
| GRANULES    | 7.85     | 294.75  | 8874157 | 304.3   | 274.4  | 20.76    |

|            |      |         |          |         |         |        |
|------------|------|---------|----------|---------|---------|--------|
| RKFORGE    | 7.8  | 156.25  | 243706   | 160     | 143     | -5.32  |
| SASKEN     | 7.79 | 584.95  | 114039   | 595.6   | 536     | 26.77  |
| ALPA       | 7.77 | 27.05   | 666647   | 27.6    | 24.35   | 1      |
| LAURUSLABS | 7.63 | 1004.5  | 10940287 | 1075.15 | 955     | 91.21  |
| HIMATSEIDE | 7.48 | 66.1    | 385132   | 66.65   | 61.5    | 0.93   |
| BASF       | 7.45 | 1432.15 | 156248   | 1447.3  | 1332    | 46.21  |
| INDOCO     | 7.42 | 243.35  | 558803   | 254.4   | 229.7   | 3.57   |
| DIXON      | 7.4  | 8198.9  | 138471   | 8358    | 7720.15 | 519.66 |
| LINCOLN    | 7.3  | 206.6   | 565273   | 211.8   | 191.7   | 8.03   |
| BALPHARMA  | 7.23 | 46      | 85579    | 46.45   | 42.5    | 0.276  |
| VISAKAIND  | 6.98 | 292.1   | 142707   | 295     | 268.85  | 12.82  |
| TATAMOTORS | 6.88 | 111.85  | 1.17E+08 | 113.5   | 102.9   | 1.34   |
| TEXINFRA   | 6.64 | 37.75   | 625822   | 38.8    | 36.55   | -0.327 |
| ALKEM      | 6.49 | 2835    | 443912   | 2836    | 2680.35 | 65.93  |
| GESHIP     | 6.41 | 240.85  | 727951   | 243     | 228.4   | 3.11   |
| POLYMED    | 6.36 | 415     | 258325   | 421     | 383.1   | 24.5   |
| SHILPAMED  | 6.16 | 584     | 462587   | 597     | 560     | 15.65  |
| JUBLINDS   | 6.1  | 113.1   | 77728    | 116.8   | 106.5   | 0.826  |
| SASTASUNDR | 5.93 | 85.8    | 9005     | 85.85   | 80.1    | 1.22   |
| GRPLTD     | 5.71 | 740     | 919      | 751.05  | 682.35  | 15.78  |
| SANOFI     | 5.67 | 8208.45 | 59322    | 8287    | 7849    | 14.37  |
| ZYDUSWELL  | 5.67 | 1702.1  | 194025   | 1739.95 | 1615    | 87.75  |
| FDC        | 5.65 | 316.05  | 3007566  | 324.4   | 300.5   | 9.79   |
| PANACEABIO | 5.5  | 227.25  | 436753   | 231     | 217.4   | 5.27   |
| ADVENZYMES | 5.1  | 199.8   | 904277   | 202.4   | 190.2   | 5.3    |
| FELDVDR    | 5    | 17.85   | 262859   | 17.85   | 17.15   | 0.578  |
| PSL        | 5    | 1.05    | 18633    | 1.05    | 1.05    | 0.132  |
| TATACOMM   | 5    | 797.7   | 169136   | 797.7   | 770.25  | 42.36  |
| CREATIVE   | 4.99 | 94.65   | 9760     | 94.65   | 93.7    | 3.93   |

TABLE V. HISTORICAL MULTI-STOCK TECHNICAL DATA WITH MACD LESS THAN ZERO

| symbol     | % change | price    | volume   | High    | Low    | MACD(<0) |
|------------|----------|----------|----------|---------|--------|----------|
| MAANALU    | 18.3     | 60.45    | 178577   | 61.3    | 48     | 0.23     |
| KGL        | 14.29    | 0.4      | 197734   | 0.4     | 0.35   | -0.022   |
| UVSL       | 11.11    | 0.5      | 18958673 | 0.5     | 0.45   | -0.013   |
| NTL        | 11.11    | 0.5      | 4300     | 0.5     | 0.5    | -0.057   |
| FCSOFT     | 11.11    | 0.5      | 2335985  | 0.5     | 0.45   | -0.004   |
| METKORE    | 10       | 0.55     | 7913     | 0.55    | 0.55   | -0.066   |
| CUPID      | 9.64     | 227.55   | 447482   | 231.8   | 209.15 | 3.9      |
| ALICON     | 8.85     | 270.1    | 11552    | 283.5   | 250    | 0.403    |
| GAMNINFRA  | 7.69     | 0.7      | 559919   | 0.7     | 0.65   | -0.017   |
| EROSMEDIA  | 7.65     | 19.7     | 753262   | 20      | 17.6   | 0.324    |
| INDNIPPON  | 7.1      | 286.5    | 158200   | 300.1   | 268.95 | -1.39    |
| TATAMTRDVR | 6.67     | 40       | 13128562 | 40.3    | 37.25  | -0.752   |
| TASTYBITE  | 6.63     | 12568.95 | 3575     | 12589.8 | 11601  | 14.29    |

|            |      |         |         |         |        |        |
|------------|------|---------|---------|---------|--------|--------|
| NECLIFE    | 6.01 | 22.05   | 1142308 | 22.45   | 20.85  | 0.132  |
| AVTNPL     | 5.96 | 40      | 512506  | 40.85   | 37.75  | 0.663  |
| NBVENTURES | 5.95 | 48.95   | 436139  | 49.5    | 46.35  | 0.012  |
| MERCATOR   | 5.88 | 0.9     | 723310  | 0.9     | 0.8    | -0.088 |
| EUROMULTI  | 5.88 | 0.9     | 13420   | 0.9     | 0.8    | 0.028  |
| BALAMINES  | 5.79 | 603.5   | 480666  | 614.75  | 566.45 | 20.17  |
| MAYURUNIQ  | 5.74 | 231.05  | 172540  | 246.9   | 217.2  | 6.96   |
| LOTUSEYE   | 5.71 | 27.75   | 2274    | 29      | 27     | -0.761 |
| INSPIRISYS | 5.59 | 26.45   | 9846    | 26.75   | 25.1   | 0.346  |
| BILENERGY  | 5.56 | 0.95    | 619895  | 0.95    | 0.9    | -0.057 |
| BFINVEST   | 5.53 | 313.9   | 40009   | 315     | 298    | 2.82   |
| LUMAXIND   | 5.49 | 1309    | 85835   | 1319.6  | 1211   | 32.92  |
| GSCLCEMENT | 5.46 | 30.9    | 221319  | 31.35   | 28.85  | 0.268  |
| IGARASHI   | 5.29 | 271.7   | 132410  | 275.65  | 258.05 | -0.776 |
| JMCPROJECT | 5.24 | 47.2    | 72962   | 48.5    | 44.85  | -0.804 |
| GEOJITFSL  | 5.17 | 36.6    | 261160  | 36.8    | 34.3   | 0.929  |
| ASTEC      | 5.11 | 987     | 190376  | 1032.9  | 936.1  | 47.17  |
| BOMDYEING  | 5.08 | 62      | 1591922 | 62.4    | 58.3   | -2.2   |
| GUJAPOLLO  | 5.05 | 171.65  | 4630    | 175.95  | 160    | -1.31  |
| OPTIEMUS   | 5    | 21      | 56611   | 21      | 21     | -0.677 |
| NATHBIOGEN | 5    | 350.7   | 10431   | 350.7   | 340    | 3.9    |
| UNITECH    | 5    | 2.1     | 2111540 | 2.1     | 2.05   | 0.024  |
| SHIRPUR-G  | 5    | 7.35    | 39968   | 7.35    | 6.8    | -0.09  |
| DQE        | 5    | 1.05    | 1150    | 1.05    | 1.05   | -0.151 |
| AFFLE      | 5    | 1772.95 | 58185   | 1772.95 | 1681.8 | 39.91  |
| BFUTILITIE | 4.99 | 229.35  | 133305  | 229.35  | 217.8  | 7.73   |
| ESTER      | 4.98 | 59      | 233721  | 59      | 56.5   | 1.72   |
| TRIVENI    | 4.98 | 56.9    | 148572  | 56.9    | 56.9   | 0.735  |
| WABAG      | 4.98 | 120.15  | 40710   | 120.15  | 120    | 2.93   |
| PREMEXPLN  | 4.97 | 115     | 2873    | 115     | 108.5  | 2.46   |
| DWARKESH   | 4.97 | 25.35   | 1119087 | 25.35   | 24.3   | 0.229  |
| SKIPPER    | 4.96 | 37      | 74287   | 37.3    | 35.25  | -0.528 |
| AURIONPRO  | 4.96 | 48.65   | 17269   | 48.65   | 45.1   | -0.829 |
| CIMMCO     | 4.96 | 20.1    | 60846   | 20.1    | 19.2   | 0.417  |
| BODALCHEM  | 4.96 | 77.25   | 1070456 | 78.9    | 72.4   | 3.5    |
| BCG        | 4.96 | 6.35    | 1785857 | 6.35    | 6.05   | -0.588 |
| VAKRANGEE  | 4.96 | 29.65   | 342576  | 29.65   | 29.65  | -0.822 |

TABLE VI. PERFORMANCE ANALYSIS OF F1-MEASURE USING DIFFERENT TRADITIONAL CLASSIFICATION LEARNING FRAMEWORKS

| StockName   | SVM  | RF   | NB   | NN   | CNN  | MultiStockFS |
|-------------|------|------|------|------|------|--------------|
| LT          | 0.91 | 0.95 | 0.9  | 0.91 | 0.89 | 0.97         |
| ASIAN_PAINT | 0.94 | 0.89 | 0.94 | 0.9  | 0.9  | 0.97         |
| AXIS_BANK   | 0.92 | 0.92 | 0.92 | 0.92 | 0.93 | 0.97         |
| BAJAJ_AUTO  | 0.93 | 0.87 | 0.94 | 0.91 | 0.87 | 0.98         |
| BAJFINANCE  | 0.89 | 0.91 | 0.91 | 0.94 | 0.89 | 0.98         |

|            |      |      |      |      |      |      |
|------------|------|------|------|------|------|------|
| BAJAJFINSV | 0.93 | 0.91 | 0.95 | 0.95 | 0.89 | 0.97 |
| BPCL       | 0.91 | 0.9  | 0.9  | 0.89 | 0.9  | 0.96 |
| BHARTIARTL | 0.94 | 0.92 | 0.88 | 0.91 | 0.92 | 0.99 |
| INFRATEL   | 0.88 | 0.88 | 0.92 | 0.92 | 0.92 | 0.96 |
| BRITANNIA  | 0.94 | 0.92 | 0.94 | 0.92 | 0.91 | 0.97 |
| CIPLA      | 0.93 | 0.89 | 0.9  | 0.93 | 0.91 | 0.97 |
| COALINDIA  | 0.94 | 0.91 | 0.91 | 0.88 | 0.94 | 0.98 |
| DRREDDY    | 0.9  | 0.88 | 0.87 | 0.88 | 0.93 | 0.97 |
| EICHERMOT  | 0.89 | 0.95 | 0.91 | 0.95 | 0.95 | 0.96 |
| GAIL       | 0.88 | 0.87 | 0.9  | 0.91 | 0.92 | 0.98 |
| GRASIM     | 0.88 | 0.9  | 0.94 | 0.92 | 0.9  | 0.98 |
| HCLTECH    | 0.88 | 0.94 | 0.94 | 0.93 | 0.92 | 0.97 |
| HDFCBANK   | 0.89 | 0.89 | 0.87 | 0.89 | 0.89 | 0.96 |
| HEROMOTOCO | 0.94 | 0.88 | 0.89 | 0.89 | 0.89 | 0.97 |
| HINDALCO   | 0.93 | 0.87 | 0.89 | 0.91 | 0.95 | 0.98 |

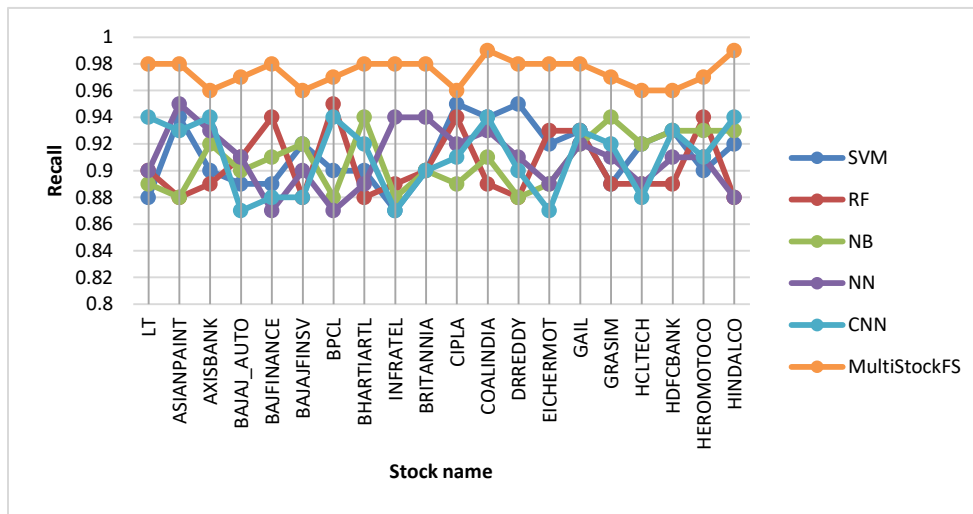


Fig. 11. Performance Analysis of Recall using different Traditional Feature Selection based Classification Frameworks.

TABLE VII. PERFORMANCE ANALYSIS OF PRECISION MEASURE USING DIFFERENT TRADITIONAL CLASSIFICATION LEARNING FRAMEWORKS

| StockName  | SVM  | RF   | NB   | NN   | CNN  | MultiStockFS |
|------------|------|------|------|------|------|--------------|
| LT         | 0.91 | 0.88 | 0.89 | 0.9  | 0.88 | 0.99         |
| ASIANPAINT | 0.87 | 0.92 | 0.88 | 0.93 | 0.88 | 0.99         |
| AXISBANK   | 0.93 | 0.94 | 0.89 | 0.88 | 0.93 | 0.97         |
| BAJAJ_AUTO | 0.88 | 0.93 | 0.87 | 0.93 | 0.89 | 0.99         |
| BAJFINANCE | 0.9  | 0.92 | 0.89 | 0.88 | 0.95 | 0.98         |
| BAJAJFINSV | 0.88 | 0.94 | 0.89 | 0.92 | 0.89 | 0.97         |
| BPCL       | 0.88 | 0.94 | 0.91 | 0.93 | 0.9  | 0.97         |
| BHARTIARTL | 0.88 | 0.92 | 0.91 | 0.88 | 0.95 | 0.98         |
| INFRATEL   | 0.89 | 0.93 | 0.95 | 0.92 | 0.9  | 0.99         |
| BRITANNIA  | 0.9  | 0.94 | 0.93 | 0.93 | 0.95 | 0.98         |
| CIPLA      | 0.92 | 0.91 | 0.88 | 0.88 | 0.88 | 0.98         |
| COALINDIA  | 0.93 | 0.95 | 0.87 | 0.89 | 0.89 | 0.96         |
| DRREDDY    | 0.93 | 0.9  | 0.94 | 0.95 | 0.94 | 0.98         |



|            |      |      |      |      |      |      |
|------------|------|------|------|------|------|------|
| EICHERMOT  | 0.88 | 0.89 | 0.91 | 0.93 | 0.9  | 0.97 |
| GAIL       | 0.95 | 0.94 | 0.91 | 0.92 | 0.94 | 0.99 |
| GRASIM     | 0.91 | 0.92 | 0.89 | 0.91 | 0.92 | 0.96 |
| HCLTECH    | 0.9  | 0.9  | 0.89 | 0.94 | 0.92 | 0.99 |
| HDFCBANK   | 0.89 | 0.89 | 0.92 | 0.91 | 0.9  | 0.98 |
| HEROMOTOCO | 0.95 | 0.93 | 0.91 | 0.9  | 0.93 | 0.98 |
| HINDALCO   | 0.93 | 0.89 | 0.92 | 0.91 | 0.91 | 0.97 |

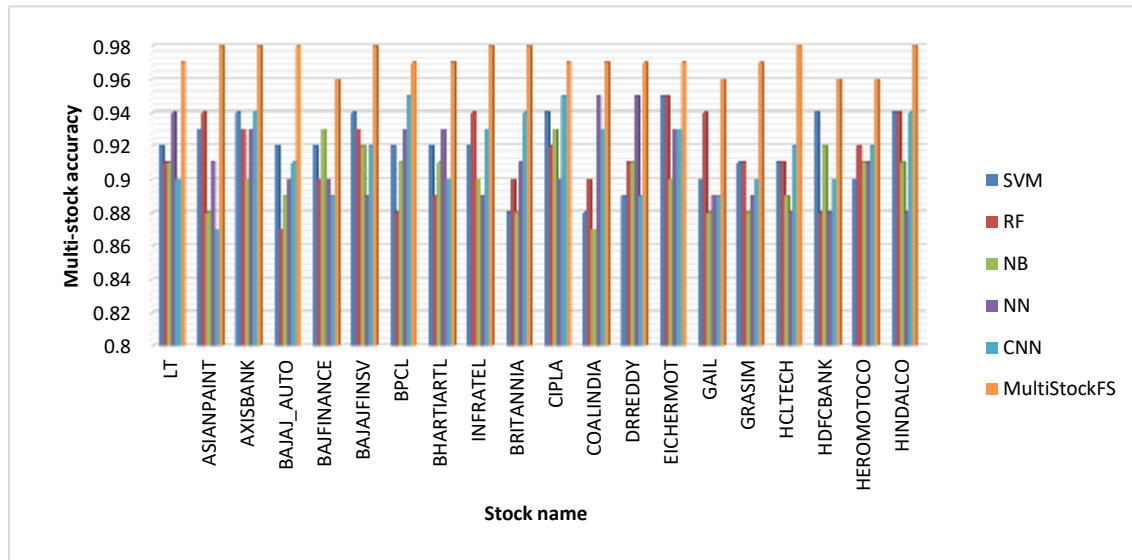


Fig. 12. Performance Analysis of Accuracy using different Traditional Classification Learning Frameworks.

## V. INFERENCE

Performance of various segments such as precision, recall, F-Measure, accuracy and runtime are improved due to data filtering and feature selection in the above model as we can see in the results from the above tables. From the above tables, it is clearly identified that the proposed feature extraction and scoring approach optimizes the stock sentiment of the social media comments and its technical data. The proposed stock feature has less runtime and more efficiency in the real time stock market databases when compared with the traditional feature extraction measures. When the traditional classifiers and the proposed non-linear classifiers are compared from the above tables, it is observed that the performance of the non-linear classifiers are better than the traditional classifiers in terms of recall precision, accuracy and runtime(ms). 12% of accuracy is obtained through the proposed model when compared to the traditional stock market prediction classifiers.

## VI. CONCLUSIONS

In this paper, a hybrid real-time multi-stock trend prediction model is designed and implemented on the stock market data. Since, most of the conventional single stock trend prediction models are depend on data size and limited feature space, it is difficult to find a novel feature selection measure on the stock technical data and stock news data. Also, these models are independent of temporal features for stock trend prediction. In this work, an advanced time based multi-stock trend prediction model is developed on the real-time data. In

this model, a new technical stock feature selection indicator and sentiment scores are computed for the clustering method. Finally, a cluster based non-linear regression framework for periodic multi-stock trend prediction is applied on the real time stock market data. Experimental results proved that the present model has better efficiency than the traditional technical indicators in terms of accuracy, f-measure, precision and recall. From the experimental results, it is observed that the proposed stock market trend prediction model has 9% of runtime (ms) and 12% of average classification accuracy as compared to the traditional trend prediction models on training and test dataset.

## REFERENCES

- [1] O. Aladesanmi, F. Casalin, and H. Metcalf, "Stock market integration between the UK and the US: Evidence over eight decades," *Global Finance Journal*, vol. 41, pp. 32–43, Aug. 2019, doi: 10.1016/j.gfj.2018.11.005.
- [2] J. Bley and M. Saad, "An analysis of technical trading rules: The case of MENA markets," *Finance Research Letters*, vol. 33, p. 101182, Mar. 2020, doi: 10.1016/j.frl.2019.04.038.
- [3] A. C. Briza and P. C. Naval, "Stock trading system based on the multi-objective particle swarm optimization of technical indicators on end-of-day market data," *Applied Soft Computing*, vol. 11, no. 1, pp. 1191–1201, Jan. 2011, doi: 10.1016/j.asoc.2010.02.017.
- [4] O. Bustos and A. Pomares-Quimbaya, "Stock market movement forecast: A Systematic review," *Expert Systems with Applications*, vol. 156, p. 113464, Oct. 2020, doi: 10.1016/j.eswa.2020.113464.
- [5] Z. Dai, X. Dong, J. Kang, and L. Hong, "Forecasting stock market returns: New technical indicators and two-step economic constraint method," *The North American Journal of Economics and Finance*, vol. 53, p. 101216, Jul. 2020, doi: 10.1016/j.najef.2020.101216.

- [6] S. R. Das, D. Mishra, and M. Rout, "Stock market prediction using Firefly algorithm with evolutionary framework optimized feature reduction for OSELM method," *Expert Systems with Applications: X*, vol. 4, p. 100016, Nov. 2019, doi: 10.1016/j.eswx.2019.100016.
- [7] S. Das, R. K. Behera, M. kumar, and S. K. Rath, "Real-Time Sentiment Analysis of Twitter Streaming data for Stock Prediction," *Procedia Computer Science*, vol. 132, pp. 956–964, Jan. 2018, doi: 10.1016/j.procs.2018.05.111.
- [8] D. P. Gandhmal and K. Kumar, "Systematic analysis and review of stock market prediction techniques," *Computer Science Review*, vol. 34, p. 100190, Nov. 2019, doi: 10.1016/j.cosrev.2019.08.001.
- [9] B. N. Jeon and B.-S. Jang, "The linkage between the US and Korean stock markets: the case of NASDAQ, KOSDAQ, and the semiconductor stocks," *Research in International Business and Finance*, vol. 18, no. 3, pp. 319–340, Sep. 2004, doi: 10.1016/j.ribaf.2004.04.006.
- [10] H. M. G. E.a., V. K. Menon, and S. K.p., "NSE Stock Market Prediction Using Deep-Learning Models," *Procedia Computer Science*, vol. 132, pp. 1351–1362, Jan. 2018, doi: 10.1016/j.procs.2018.05.050.
- [11] M. S. Pagano, L. Peng, and R. A. Schwartz, "A call auction's impact on price formation and order routing: Evidence from the NASDAQ stock market," *Journal of Financial Markets*, vol. 16, no. 2, pp. 331–361, May 2013, doi: 10.1016/j.finmar.2012.11.001.
- [12] M. Vijh, D. Chandola, V. A. Tikkiwal, and A. Kumar, "Stock Closing Price Prediction using Machine Learning Techniques," *Procedia Computer Science*, vol. 167, pp. 599–606, Jan. 2020, doi: 10.1016/j.procs.2020.03.326.
- [13] Y. Zhang and L. Wu, "Stock market prediction of S&P 500 via combination of improved BCO approach and BP neural network," *Expert Systems with Applications*, vol. 36, no. 5, pp. 8849–8854, Jul. 2009, doi: 10.1016/j.eswa.2008.11.028.
- [14] M.-Y. Chen, C.-H. Liao, and R.-P. Hsieh, "Modeling public mood and emotion: Stock market trend prediction with anticipatory computing approach," *Computers in Human Behavior*, vol. 101, pp. 402–408, Dec. 2019, doi: 10.1016/j.chb.2019.03.021.
- [15] R. Jacinto, E. Reis, and J. Ferrão, "Indicators for the assessment of social resilience in flood-affected communities – A text mining-based methodology," *Science of The Total Environment*, p. 140973, Jul. 2020, doi: 10.1016/j.scitotenv.2020.140973.
- [16] M. M. Rounaghi and F. Nassir Zadeh, "Investigation of market efficiency and Financial Stability between S&P 500 and London Stock Exchange: Monthly and yearly Forecasting of Time Series Stock Returns using ARMA model," *Physica A: Statistical Mechanics and its Applications*, vol. 456, pp. 10–21, Aug. 2016, doi: 10.1016/j.physa.2016.03.006.
- [17] X. Li, P. Wu, and W. Wang, "Incorporating stock prices and news sentiments for stock market prediction: A case of Hong Kong," *Information Processing & Management*, vol. 57, no. 5, p. 102212, Sep. 2020, doi: 10.1016/j.ipm.2020.102212.
- [18] J. H. Yu, J. Kang, and S. Park, "Information availability and return volatility in the bitcoin Market: Analyzing differences of user opinion and interest," *Information Processing & Management*, vol. 56, no. 3, pp. 721–732, May 2019, doi: 10.1016/j.ipm.2018.12.002.
- [19] Z. Zhou, M. Gao, Q. Liu, and H. Xiao, "Forecasting stock price movements with multiple data sources: Evidence from stock market in China," *Physica A: Statistical Mechanics and its Applications*, vol. 542, p. 123389, Mar. 2020, doi: 10.1016/j.physa.2019.123389.
- [20] S. R. Das, D. Mishra, and M. Rout, "Stock market prediction using Firefly algorithm with evolutionary framework optimized feature reduction for OSELM method," *Expert Systems with Applications: X*, vol. 4, p. 100016, Nov. 2019, doi: 10.1016/j.eswx.2019.100016.
- [21] S. Deng, C. Wang, M. Wang, and Z. Sun, "A gradient boosting decision tree approach for insider trading identification: An empirical model evaluation of China stock market," *Applied Soft Computing*, vol. 83, p. 105652, Oct. 2019, doi: 10.1016/j.asoc.2019.105652.
- [22] S. S. Pal and S. Kar, "Time series forecasting for stock market prediction through data discretization by fuzzistics and rule generation by rough set theory," *Mathematics and Computers in Simulation*, vol. 162, pp. 18–30, Aug. 2019, doi: 10.1016/j.matcom.2019.01.001.
- [23] L.P.Maguluri, and R. Ragupathy "A New Sentiment Score Based Improved Bayesian Networks For Real-Time Intraday Stock Trend Classification". *International Journal of Advanced Trends in Computer Science and Engineering*, vol. 8, no. 4, 2019, pp. 1045-1055. *The World Academy Of Research In Science And Engineering*, doi:10.30534/ijatcse/2019/10842019. Accessed 20 July 2020.
- [24] L.P.Maguluri, and R. Ragupathy, "An Efficient Stock Market Trend Prediction Using the Real-Time Stock Technical Data and Stock Social Media Data", *International Journal of Intelligent Engineering and Systems*, Vol.13, No.4, 2020, DOI: 10.22266/ijies2020.0831.28.
- [25] M. Syamala, and N.J. Nalini, "A deep analysis on aspect based sentiment text classification approaches", *International Journal of Advanced Trends in Computer Science and Engineering*, Volume 8, No.5, September - October 2019, <https://doi.org/10.30534/ijatcse/2019/01852019>.
- [26] M. Syamala, and N.J. Nalini, "A filter based improved decision tree sentiment classification model for real-time amazon product review data", *International Journal of Intelligent Engineering and Systems*, Vol.13, No.1, 2020 DOI: 10.22266/ijies2020.0229.18.

# Development of a Graphic Information System Applied to Quality Statistic Control in Production Processes

Laura Vázquez<sup>1</sup>, Alicia Valdez<sup>2</sup>, Griselda Cortes<sup>3</sup>, Mariana Rosales<sup>4</sup>  
Research Center Autonomous University of Coahuila  
Coahuila, Mexico

**Abstract**—One of the advantages that organizations have when using an Information System is the control of their activities. This article develops an Information System that will allow an organization to graphically obtain the real results of a production process by applying Nelson's eight rules to determine if any measured variable is out of control. The software architecture pattern used is the Model View Controller (MVC) to keep the functionality of the application separate. The front-end, that is, the part that interacts with the users, was developed in ASP.NET as a web platform to provide the required services, JavaScript, HTML 5, Razor and Bootstrap. The back-end, which is the part that processes the entry of the front-end and performs the calculations, operations, communication with the database and reading of files, was developed with the C sharp programming language, the SQL Server database management system and the entity framework. As a result, the system sends an e-mail as an alarm with an explanation of what has happened when it detects that some measured variable is out of control by applying Nelson's rules. This allows the organization to make effective decisions in the processes involved.

**Keywords**—Information system; Nelson's rules; Model View Controller pattern; C#; ASP.NET

## I. INTRODUCTION

Technology can help all kinds of businesses improve the efficiency and effectiveness of their business processes, managerial decision making, and workgroup collaboration, which strengthens their competitive positions in rapidly changing marketplaces [1].

Also, technology is used to solve problems, through its systems that must be adapted to the needs of the organizations.

Information Systems and Technologies are vital components of successful businesses and organizations some would say they are business imperative [1]. Information System is defined as group of elements organized with the purpose of supporting management and operational decision making [2].

It is important to mention that Information Systems can be developed using fourth-generation software tools; their functionality is that users in organizations can access data, create, and interpret reports quickly to facilitate processes.

The benefits that an organization has when using an Information System in the long term are the following: automation of operational processes, provision of an

information platform for decision making, and achievement of competitive advantage.

In the production processes, different factors and elements are combined on which measurements must be made, that is to say, in each process it is necessary to control variables such as pressure, weight, flow, etc. to guarantee the quality of a product.

According to the above, it is necessary to develop a system as a support tool that allows us to obtain graphically the real results of a productive process, applying Nelson's eight rules. As it will be explained later, these rules will allow controlling the production process to determine if any measured variable is out of control.

The theoretical fundamental for the development of the graphic Information System is shown below, considering the MVC software architecture model and the software used for the back-end and front-end. The results expected at the end of the processing as output are directly related to the characteristics and processing development of the proposed Information System, which guarantees the efficiency of its effectiveness.

## II. RELATED WORKS

For the area of education, the work in [3] presented the Development of Information System for a University.

The student Information System of a university stores and tracks all student data which are needed by the faculty and staff to manage the operations of the university. Information such as grades, attendance records, admission information, and financial aid are tracked through these platforms [3].

In this paper [3], students Information System has been developed to maintain the information and other content of digitized Information using ADO .NET technology and Microsoft SQL. This system is mainly intended to be used by the staff from the student affairs department and faculties of the university.

In the paper titled "Implementation of Scrum work framework in the development of quality assurance Information System" uses the Scrum agile development methodology.

The purpose of this research is to develop a quality assurance Information System by implementing the Scrum

Framework. Scrum is one of the popular frameworks in Agile Development Methodology. In this way, the development of productivity increases significantly. In this Applied Research, the Action Research approach is used [4].

The work in [5] presented a system that was designed using PHP and MySQL as the programming language for the database. The system can classify the brown sugar by calculating the weight of the criterion. Besides, the classification process is performed to determine the optimal value.

### III. THEORICAL FUNDAMENTAL

#### A. Information System

An Information System can be any organized combination of people, hardware, software, communications networks, data resources, and policies and procedures that stores, retrieves, transforms, and disseminates information in an organization [1].

The requirements of an Information System are determined by the objectives of the organization for which the system is being designed and built [6].

The functions of Information Systems according the reference [2] are the following:

- As a source of information to help in effective decision making by managers.
- A contributor to productivity efficiency and customer satisfaction.
- The Information System is useful to achieve success in various functions such as Finance, operations, marketing, human resource, and store management.

In addition, Bagad in his book "Managenent Information Systems" [2] comments on Information System Activities. It is indicated that in a business process, the different information processing activities take place. For example:

- 1) The input of data resources
- 2) Processing of data into information
- 3) The output of information products
- 4) Storage of data resources
- 5) Control of system performance

#### B. Nelson's Rules

The Nelson rules were published first in the October 1984 issue of the Journal of Quality Technology in an article by Lloyd S. Nelson. Nelson rules are methods in process control of determining if some measured variable is out of control [7].

These zones are used in conjunction with a set of pattern analysis rules to determine when a process has gone out of control [9]. In general, when identifying these rules, the region between the usual  $\pm 3$  sigma limits are divided into six region and the pattern is explained with respect to  $\pm 1, 2$  and  $3$  sigma limits as shown in the Fig. 1 [8].

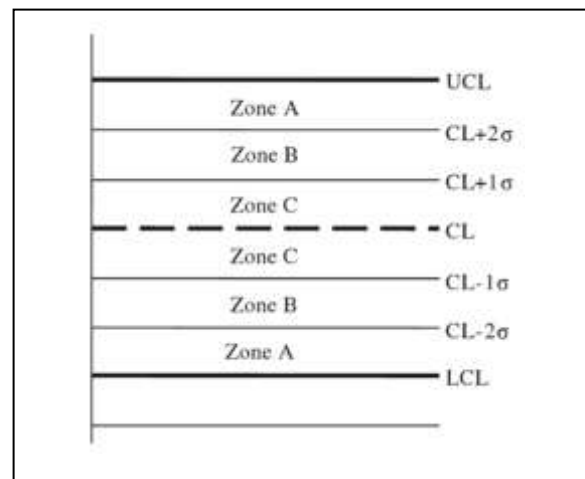


Fig. 1. Zones in a Control Chart. Source [8].

Nelson's rules shown in his 1984 article [7] are as follows:

- Rule 1. One point is more than three standard deviations from the mean.
- Rule 2. Nine (or more) points in a row are on the same side of the mean.
- Rule 3. Six (or more) points in a row are continually increasing (or decreasing).
- Rule 4. Fourteen (or more) points in a row alternate in direction, increasing then decreasing.
- Rule 5. Two (or three) out of three points in a row are more than two standard deviations from the mean in the same direction.
- Rule 6. Four (or five) out of five points in a row are more than one standard deviation from the mean in the same direction.
- Rule 7. Fifteen points in a row are all within one standard deviation of the mean on either side.
- Rule 8. Eight points in a row exist, but none within one standard deviation of the mean, and the points are in both directions from it.

The rules apply to an XS control chart in which the magnitude of some variable is plotted against time. The rules are based on the mean value and standard deviation of the samples considered through time.

#### C. Software Development Tools

##### 1) Integrated Development Environment

IDE support means the tools can generate code, help you write code, and provide features and artifacts that accelerate your coding. Here is where many third-party languages often fall short. It takes a lot to provide IDE support to build the many application types Visual Studio enables [10]. Visual Studio is part of the family of integrated development environments (IDE) [11].

## 2) Programming Language

Visual C# is a programming language designed for those who are familiar and comfortable programming in C-Style languages (such as C, C++, and Java). C# is type -safe, object-oriented, and targeted for rapid application development. C# developers tend to spend more of their time inside the Visual Studio code editor and less time with the designers [10].

ASP.NET is a web platform that provides all the services that you require to build enterprise-class server-based Web applications. One of the major benefits of ASP.NET is the change from interpreted code, previously used for Classic ASP (the programming model before ASP.NET), to compiled code, allowing web application to have better performance [12].

In addition, ASP.NET includes the following features [12]:

- A page and controls frameworks
- The ASP.NET compiler
- Security infrastructure
- Application configuration
- Health monitoring and performance features
- Debugging support

Razor is a template syntax that allows you to combine code and content in a fluid and expressive manner. Razor lets you write code using languages such as C# or Visual Basic.NET [13].

ASP.NET Razor uses a simple programming syntax that lets you embed server-based code into a web page [12].

ASP.NET MVC is a free fully supported framework for building web applications that use the model-view-controller pattern. The MVC pattern itself makes it easier to manage complexity by clearly separating the functionality of the application into three core parts, the model, the view, and the controller [12].

Fig. 2 shows a simple implementation of the MVC pattern.

The straight arrows indicate direct associations, whereas curved arrows identify indirect associations.

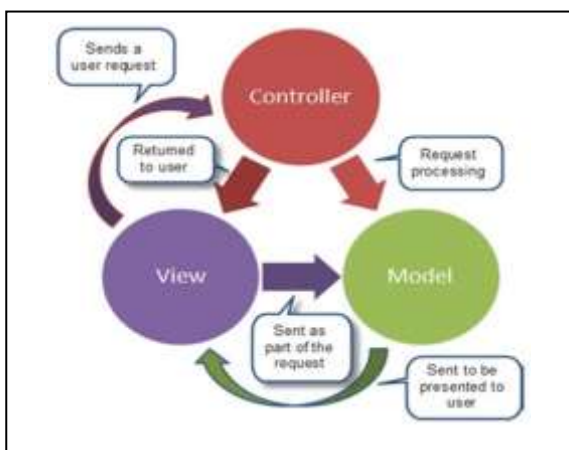


Fig. 2. Model-view-Controller Implementation. Source [12].

Explanation of Fig. 2 is shown below:

Models: models are basically a collection of classes where you will be working with data logic as well as business logic. Views: views are nothing but a pure HTML which decides how the UI (User Interface) is going to look like. Controller: Controllers are basically a bunch of classes which handles communication from the user [14].

ASP.NET MVC works well for web applications that are supported by large teams of developers and for web designers who need a high degree of control over the HTML [12]. ASP.NET is a complete and effective low-cost answer. It allows processes to be separated effectively continued with its effectiveness. This is done through the MVC solution. The ease of maintenance and scalability of ASP.NET depends on the skills of the developer. If the developer has enough programming knowledge, in the development of C#, the framework will be simple to use and take advantage of its main characteristics.

JavaScript is a programming language, light, interpreted, object-oriented, prototype-based and first-class functions, better known as the Web's scripting language. It is the programming language that Netscape created to bring your browser to life (on the front end client side) [15]. JavaScript is a scripting language that allows you to create dynamically updated content, control multimedia, animate images, and just about everything else.

The main goal of Bootstrap is to provide a web frontend framework for responsive developing with cross-browser compatibility [16]. Bootstrap is an open source product. It has evolved from being an entirely CSS-driven project to include a host of JavaScript plugins and icons that go hand in hand with forms and buttons [17].

One of the highlights is the build tool on Bootstrap's website, where you can customize the build to suit your needs, choosing which CSS and JavaScript features you want to include on your site [17].

CSS is the recommended format for pages written in HTML based on "Cascading Style Sheets" standards. CSS allows the use of methods to create structures by treating the styles separately, reducing the rendering time, offering the client a faster and more efficient connection.

Also, HTML documents using CSS are small, since the style design is used. Therefore, the page increases the speed of transmission of content information, benefiting customers, page owners, and Web server administrators.

SQL Server is an enterprise-class database management system (DBMS) that is capable of running anything from a personal database only a few megabytes in size on a handheld Windows Mobile device up to a multiserver database system managing terabytes of information [18].

The database is effectively the highest-level object that you can refer to within a given SQL Server [19]. Besides, SQL Server offers a variety of administrative tools to ease the burdens of database development, maintenance, and administration.

#### IV. METHODOLOGY

Fig. 3 shows the diagram of the Information System developed.

The MVC software design pattern was used in the development of the system. The part of the front-end, that is, the one that interacts with the users was developed in ASP.NET, JavaScript, HTML 5, Razor and Bootstrap. On the other hand, for the back-end, which is the part that processes the input from the front-end and performs the calculations, operations, communication with the database and file reading, it was developed with C#, SQL Server and entity framework, as described.

ASP.NET is the development environment used for the creation of the web page, using C# as the back-end.

Fig. 4 shows the design of the controller view model of the Information System developed.

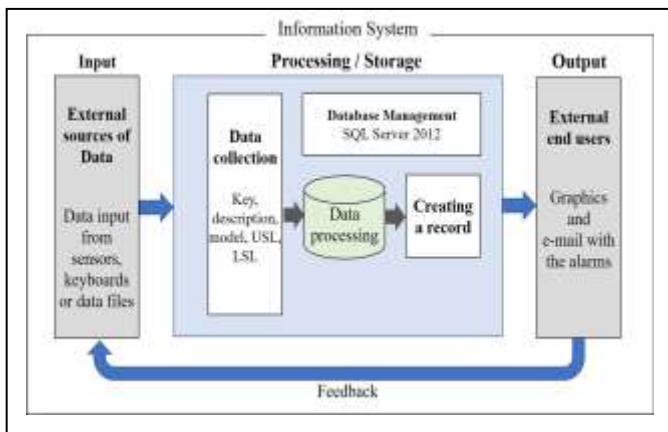


Fig. 3. Diagram of the Graphic Information System.

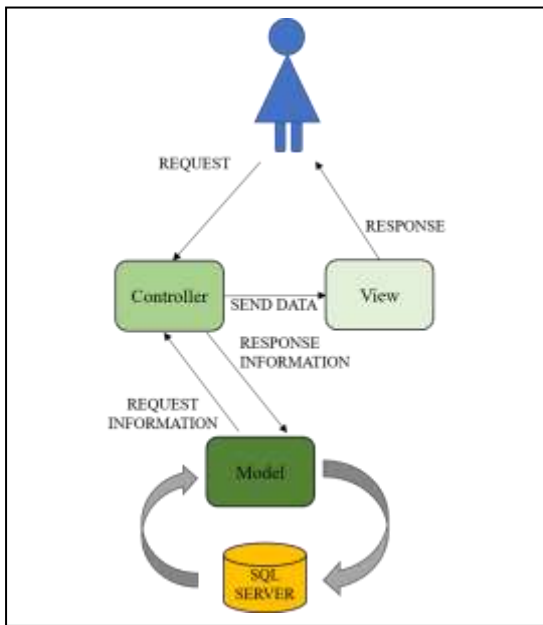


Fig. 4. Model Vista Controller Pattern Designed for the System.

In ASP.NET a web server was created with the style of MVC software architecture. When the client tries to navigate to a specific URL, the URL is taken by the controller, it uses the models (classes) and returns a view that is the page shown to the client.

To create this page, environments and languages were used such as HTML 5 to create the graphics, JavaScript to perform calculations on the client-side, Bootstrap and CSS for the design part (fonts and colors).

For the back-end, C# was used to communicate with databases in SQL Server 2012 and to perform calculations. It is significant to mention that this database management system was used because ASP.Net has stable communication with SQL Server and offers access to .net functionalities.

Fig. 5 shows the tables used in the developed system.

As mentioned above, one element that Information Systems must have is data entry. In the system designed the ways to make the data entries are the sensors, the keyboard or from data files.

To start working in the system, the corresponding user must capture the information requested in the window shown in Fig. 6.

This information corresponds to the configuration of the variables of the production processes that will be subject to measurement. The data are headed as follows: key, description, model, upper specification limit (USL) and lower specification limit (LSL), and will form the records in the database used.

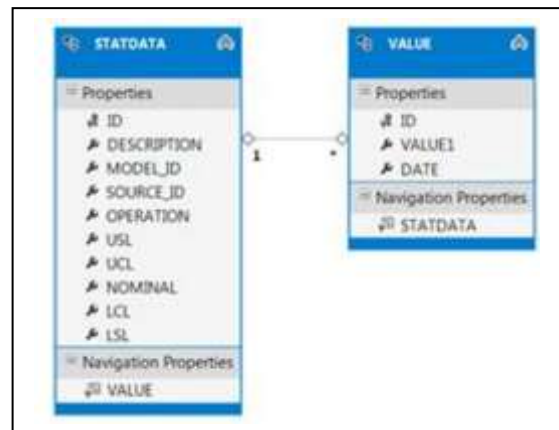


Fig. 5. System Tables.

The screenshot shows a web form titled 'PRODUCTION PROCESS DESCRIPTION'. It contains the following input fields:
   
- KEY
   
- DESCRIPTION
   
- MODEL
   
- UPPER SPECIFICATION LIMIT (USL)
   
- LOWER SPECIFICATION LIMIT (LSL)
   
At the bottom, there are two status indicators: a green checkmark and a red X.

Fig. 6. Configuration Data Recording.

V. RESULTS

Once the system has the necessary data to work, such as: the keys, the description, the model, the USL and the LSL, form the data collection that will be processed and stored. This information corresponds to the variables that will be measured in the production processes of the organization. Nelson's eight rules will apply to these variables for their corresponding evaluation.

The variables considered for evaluation are determined by the productive processes of each organization.

Fig. 7 shows some of the records captured in the system. The data considered are the following: total weight, total volume, humidity test, noise level, resistance, temperature, etc. In this same window, it is possible to make some changes in the data of the records.

The above information will be stored in the database of the SQL Server 2012 management system. Once the information is processed, the system will show the data output graphically, as it will be explained later.

Fig. 8 shows the data captured in the system, which corresponds to the measured variable "weight" of a production process.

The system considers the observations of the sample made, calculates the arithmetic mean and the standard deviation of the set of results. It also uses the established factors to build control diagrams, applying the corresponding coefficients to calculate the graph control limits in the formulas allowing the creation of XS control charts.

| KEY  | DESCRIPTION              | MODEL       | USL | LSL | CHARTEDIT  |
|------|--------------------------|-------------|-----|-----|--|
| 4421 | TOTAL WEIGHT             | 87DJW-14597 | 300 | 298 | <input type="checkbox"/> <input checked="" type="checkbox"/> |
| 2514 | TOTAL VOLUME             | 14ERF-35984 | 470 | 460 | <input type="checkbox"/> <input checked="" type="checkbox"/> |
| 3682 | METAL DETECTION TEST     | 78ERF-88690 | 25  | 20  | <input type="checkbox"/> <input checked="" type="checkbox"/> |
| 1694 | HUMIDITY TEST            | 45NBQ-1478  | 120 | 117 | <input type="checkbox"/> <input checked="" type="checkbox"/> |
| 8547 | NOISE LEVEL              | 23BQ5-3625  | 58  | 55  | <input type="checkbox"/> <input checked="" type="checkbox"/> |
| 3425 | RESISTANCE               | 21LKS-7594  | 27  | 25  | <input type="checkbox"/> <input checked="" type="checkbox"/> |
| 7864 | TEMPERATURE              | 23LON-3222  | 131 | 129 | <input type="checkbox"/> <input checked="" type="checkbox"/> |
| 2500 | VERIFICATION OF SYMMETRY | 25V5W-3651  | 544 | 541 | <input type="checkbox"/> <input checked="" type="checkbox"/> |

Fig. 7. Data Storage.

| Sample (n) | Time  | Sample observations |                |                |                |                | x̄      | s       |
|------------|-------|---------------------|----------------|----------------|----------------|----------------|---------|---------|
|            |       | x <sub>1</sub>      | x <sub>2</sub> | x <sub>3</sub> | x <sub>4</sub> | x <sub>5</sub> |         |         |
| 1          | 08:00 | 299.76              | 299.20         | 299.84         | 299.64         | 299.77         | 1676.32 | 0.32279 |
| 2          | 08:00 | 299.72              | 299.60         | 299.76         | 299.72         | 299.70         | 1676.22 | 0.31847 |
| 3          | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 4          | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 5          | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 6          | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 7          | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 8          | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 9          | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 10         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 11         | 08:00 | 299.72              | 299.60         | 299.76         | 299.72         | 299.70         | 1676.22 | 0.31847 |
| 12         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 13         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 14         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 15         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 16         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 17         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 18         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 19         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 20         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 21         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 22         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 23         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 24         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 25         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 26         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 27         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 28         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 29         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 30         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 31         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 32         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 33         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 34         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 35         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 36         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 37         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 38         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 39         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 40         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 41         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 42         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 43         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 44         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 45         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 46         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 47         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 48         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 49         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 50         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 51         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 52         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 53         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 54         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 55         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 56         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 57         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 58         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 59         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 60         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 61         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 62         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 63         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 64         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 65         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 66         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 67         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 68         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 69         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 70         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 71         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 72         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 73         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 74         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 75         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 76         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 77         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 78         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 79         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 80         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 81         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 82         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 83         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 84         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 85         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 86         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 87         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 88         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 89         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 90         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 91         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 92         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 93         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 94         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 95         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 96         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 97         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 98         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 99         | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| 100        | 08:00 | 299.84              | 299.20         | 299.60         | 299.60         | 299.59         | 1676.24 | 0.31860 |
| AVG        |       | 299.32              |                | 0.60966        |                |                |         |         |

Fig. 8. Calculations made by the Information System.

Once the configuration of the corresponding parameters is done, and the data entered in the Information System completes, the system makes the control charts to analyze and determine if any variable is out of control.

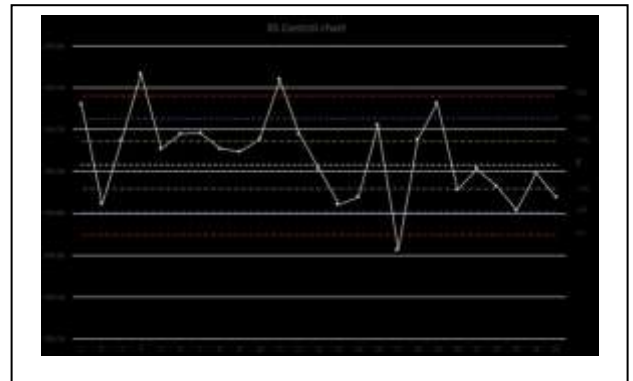
Nelson's eight rules apply to each of the variables involved in the production process. The magnitude of the variable gets plotted against time. As mentioned above, the rules bases on the arithmetic mean and standard deviation of the samples.

As a result, the developed Information System performs the processing and analysis of the variables and displays the XS control chart. Fig. 9 shows the XS control graph corresponding to the measured variable "weight", where the real results of the arithmetic mean applied to the samples are observed. In addition, the upper and lower control limits and the additional 1 and 2 sigma limits are plotted on the control chart, as shown.

Fig. 10 shows the XS control chart of the previous variable, but with the standard deviation data calculated. The control limits and the 1 and 2 sigma limits are the same.

Remembering Nelson's Rule 6 which says: Four (or five) out of five points in a row are more than one standard deviation from the mean in the same direction [7], as shown in Fig. 11.

As can be seen, in the graph in Fig. 9, Nelson's rule number 6 is presented, since samples 5, 6, 7, 8, and 9, which are points in a row, are more than one standard deviation from the mean in the same direction, as shown in Fig. 12.



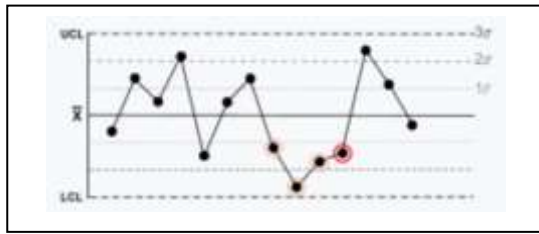


Fig. 11. Nelson's Rule 6.

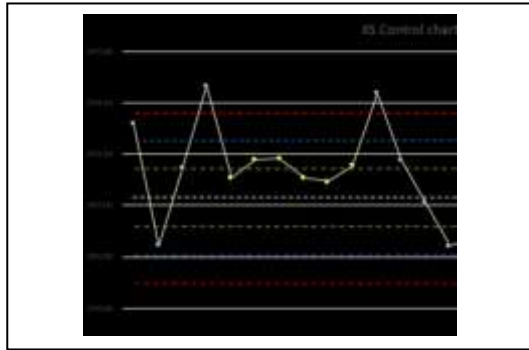


Fig. 12. First Segment where Nelson's Rule 6 was Detected.

At another point in the same graph, the Information System detects Nelson's Rule 6, as shown in Fig. 13 and displays it as a result in the alarm, as can be seen in Fig. 14.

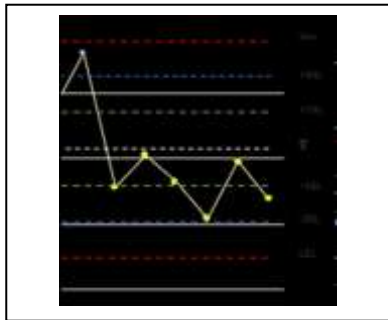


Fig. 13. Second Segment where Nelson's Rule 6 is Detected.

The Information System shows as a result once the data analysis is done and graphed, an e-mail with the alarms that have occurred when applying Nelson's eight rules and detecting if any measured variable is out of control, as shown in the Fig. 14. In this case, the variable measured was the total weight that has a KEY 4421.

As explained above, Nelson's rules are a set of methods that detect trends in the analyzed data. The system also shows the following values: Upper Specification Limit (USL) and Lower Specification Limit (LSL), as well as the automatically calculated values: Upper and Lower Control Limits, average and standard deviation.

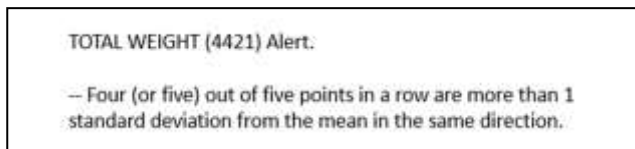


Fig. 14. Alarm Presented by the Information System.

The developed Information System allows to analyze in a graphic way the real measures in the productive processes, facilitating the decision making in real time by issuing the corresponding alarms.

## VI. CONCLUSION

As a conclusion, the Information Systems allows the organization to have effective control over activities. In addition, such systems help to increase the effectiveness in the productive processes, since they anticipate and foresee situations so that the decision making is effective.

Nelson's rules allow us to know and analyze the variables that are measured in the production processes and that are out of control, improving the quality of the products.

The use of these rules facilitates the interpretation of the control charts, but the experience of the human expert should indicate to what extent the results should be taken.

A trend in Information Systems is the use of technology as an essential part of corporate strategy, achieving an administrative advantage by facilitating decision making. In addition, technology is important for design and production control.

The use of technology will transform the organization, changing its structure. The information technology will significantly support the redesign of the processes through the Information Systems.

It is important to mention that the future work for the application and development of this paper could be to develop usable systems to facilitate the monitoring of process quality from the design of a product until it leaves the factory and is delivered to the final customer. On the other hand, it can be used in complementary areas for quality improvement in production processes, such as sales, inventory and purchases.

## VII. DISCUSSION

With the development of Science and technology, the Information System is essential in our daily life, according to the reference [3].

The use of Information Systems will continue to be useful and effective in decision-making in various areas of each organization.

Organizations should use a system that ensures that their manufactured products meet specified quality standards, keeping in mind that the system must be tailored to the needs of that organization.

The author in [4] indicates an important factor that must be taken into account for the successful development of the Information System proposed in its research: human resources management, that is, the importance of the composition of the staff teams involved.

## REFERENCES

- [1] R. Behl, J. A. O'Brien and G. M. Marakas, Management Information Systems, 11th ed. India: Mc GrawHill Education.
- [2] V. S. Bagad, Management Information Systems, 4th ed. Pune: Technical Publications Pune, 2009, pp. 1,12.



- [3] T. Htwe and M. Phyo Aung, "Development of Information System for a University", *International Journal of Scientific and Research Publications*, vol. 9, no. 6, pp. 494-499, June 2019, doi: 10.29322/IJSRP.9.06.2019.p9071.
- [4] M. Broto, B. Indiarso and D. Prayitno, "Implementation of Scrum work framework in the development of quality assurance Information System", *Jurnal Penelitian Pos dan Informatika*, vol. 9, no. 2, pp. 125-139, November 2019, doi: 10.17933/jppi.2019.090204.
- [5] I. Taufik, A. Saleh, C. Slamet, D. S. Maylawati, M. A. Ramdhani and B. A. Muhammad, "Decision support system design for determining brown sugar quality with weighted product method", *Journal of Physics: Conference Series*, vol. 1280, no. 2, pp. 1-8, doi:10.1088/1742-6596/1280/2/022019.
- [6] A. Olivé, *Conceptual Modeling of Information Systems*. Berlin: Springer, 2007, pp. 2.
- [7] L. S. Nelson, "The Shewhart Control Chart — Tests for Special Causes", *Journal of Quality Technology*, 16 no. 4, pp. 237-239, October 1984.
- [8] C. Kahraman and S. Yanik, *Intelligent Decision Making in Quality Management: Theory and Applications*. Switzerland: Springer, 2016, pp. 41-42.
- [9] S. Boslaugh, *Statistics in a Nutshell*, 2nd ed. United States of America: O'Reilly, 2012, p.p. 345.
- [10] M. Snell and L. Powers, *Microsoft Visual Studio 2012 Unleashed*, 1st ed. Indiana: Pearson Education, Inc., 2012.
- [11] B. S. Guérin, *ASP.NET en C# con Visual Studio 2015: Diseño y desarrollo de aplicaciones Web*. Barcelona: Ediciones ENI, 2016, pp. 18.
- [12] J. R. Guay, *Beginning ASP.NET MVC 4, The expert's voice in .NET*. Apress, pp. 1,2, 6-9.
- [13] J. Chadwick, *Programming Razor: Tools for Templates in ASP.NET MVC or WebMatrix*, 1st ed. United States of America: O'Reilly, 2011, pp. 3.
- [14] R. Sahay, *Hands on with ASP.NET MVC: Covering MVC 6*, Quills Ink Publishing, 2014.
- [15] A. A. Castillo, *Curso de Programación Web: JavaScript, Ajax y jQuery*, 2nd ed. ITCampus Academy, 2017, pp. 13-14.
- [16] S. Moreto, M. Lambert, B. Jakobus and J. Marah, *Bootstrap 4 – Responsive Web Design*, Birmingham: Packt Publishing, 2016, pp. 1.
- [17] J. Spurlock, *Bootstrap: Responsive Web Development*, 1st ed. United States of America: O'Reilly, 2013.
- [18] R. Rankins, P. Bertucci, C. Gallelli and A. T. Silverstein, *Microsoft SQL Server 2012 Unleashed*, 2nd ed. United States of America: Pearson Education, Inc., 2014, pp. 9.
- [19] P. Atkinson and R. Vieira, *Beginning Microsoft SQL Server 2012 Programming*, Indianapolis: John Wiley & Sons, Inc., 2012, pp. 2.

# Netnography and Text Mining to Understand Perceptions of Indian Travellers using Online Travel Services

Dashrath Mane<sup>1</sup>, Dr. Prateek Srivastava<sup>2</sup>

Department of CSE, Sir Padampat Singhania University [SPSU]  
Udaipur, India

**Abstract**—Advancements in the electronic commerce industry have helped online travel services (OTA) in many ways. The paper examines the overall impact of traveller's using online services and their sentiments derived from a collection of reviews for online travel service providers known as online travel agents (OTA) in India. Customer reviews from different identified sources are collected and the satisfaction of travellers using various online travel services is analyzed using netnographic analysis and text mining. This paper also covers a detailed process of data collection, analysis using netnography and text mining methods which helps us for the analysis and deriving sentiments from collected reviews. Various results obtained are presented as part of token lists, keyword analysis, and service-specific analysis. The statistical analysis of different results is tested to understand the relationship between various services and OTA.

**Keywords**—Consumer; travellers; netnography; text mining; OTA; sentiment; perception

## I. INTRODUCTION

### A. Tourism Market in India

India, the country is known for its rich history, cultures. India placed 40th out of 136 nations as per the report of Travel & Tourism Competitiveness for the year 2017 [1]. The tourism industry in India is playing a significant role in the economy [2].

The current scenario shows that India has a bright future in the tourism sector as lots of development taking place in the Indian tourism industry. In the last few years, the Ministry of Tourism, Government of India, has initiated various attempts to lift tourism in a country that has opened research opportunities in it. The plans like PRASAD, e-tourist Visa, Mobile applications for tourists are helping out to gain a new tourist footprint [3].

### B. Online travel services in India

It is expected that by the end of 2020, the count of the bookings made online is likely to reach 4 billion in India.[4] Makemytrip.com has been positioning high on the rundown of movement sites that fill in as a one-stop look for the whole group. The online booking travel related firm is a commonly known name these days & it has been collaborating with imperative players in the industry. Yatra.com[5] is the best for movement specialists, long-standing corporate customers, and easygoing hikers alike & is an excellent travel site in India.

The best costs and cashback bargains on air are assured with yatra.com.

Cleartrip online travel service provider offers internet booking for train and flight tickets for domestic and international bookings. Goibibo.com is another name which offers the best arrangements/bookings in trains, flight. From booking universal/local occasions to encouraging Foreign Exchange, Visa, Passport, protection for movement purposes. Thomas cook does this everything for the travellers. Travelguru.com, regularly appraised as a standout amongst the most went by Indian travel destinations on the web, is a fortune place of treats with regards to booking air tickets, lodging offices.

Netnography and Tourism industry: Netnography is Method used for studying culture and communities online. It is a tool to understand social interaction in digital communications [6]. Netnography is a known tool accepted as a virtual ethnography tool, is being explored along with text mining. The approach used in this paper also produces an analysis of perceptions of travellers in India using text mining techniques.

Phases in the Netnography Process [11]: Fig. 1 explains various phases in the process of netnographic analysis; below are the various phases in it:

1) *Research planning*: Research planning [3] is the initial phase in netnographic studies. This phase speaks about defining the problem, defining the research objectives, should talk about translating the research objectives into a specific set of questions.

2) *Entrée*: In this phase, online communities, blogs, groups are identified and allows an ethnographic to enter into the communities to get a better knowledge of cultures and communities.

3) *Data collection*: Data collection is a very vital phase in this type of study. In this phase, Netnographer collects data from Internet data, interview data, and field notes. Travellers' reviews downloaded using various rating & referral sites like mouthshut.com and consumeraffairs.com. These are the websites that allow the user to collect millions of reviews posted by travellers in India. The whole reviews were collected using systematics steps in the netnographic process[7].

4) *Interpretations/Data analysis*: This makes use of One of the data analysis technique known as Analytical coding, Renal data analysis consist of coding, noting, abstracting, checking and refining, generalizing and theorizing as shown in Fig. 2.

5) *The data collection* process was performed from sources like mouthshut.com [8] and Consumer Affairs.com for this study. More than 2000 reviews from these platforms have been collected and maintained in an excel sheet.

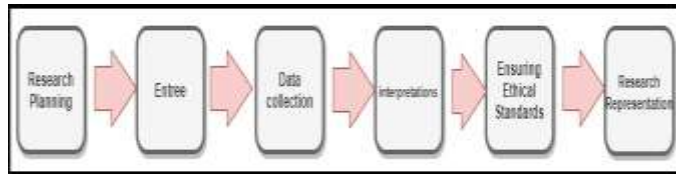


Fig. 1. Phases in Netnography Process [11].

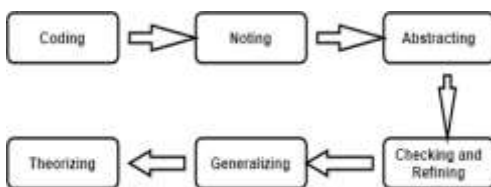


Fig. 2. Analytical Data Analysis [11].

## II. RESEARCH METHODOLOGY

The online travel agents considered in this research work includes MakeMyTrip, Goibibo, Cleartrip, redBus, and Yatra. A good number of reviews for each of these agents [12] selected are evaluated using netnographic studies. The analysis performed as netnographic analysis is used to carry out the sentiments, the levels of sentiments of travellers using

this online travel-related online services in India. Fig. 3 shows the overall rating for Cleartrip on mouthshut.com.

Fig. 4 represents the sample to review and rating for Cleartrip [7]. In addition to actual comment, Reviews contain values from 1 (low) to 5 (high) for all the website components.

Fig. 5 is a sample of the overall rating for MakeMyTrip. Referring to the Fig. 5, it is very much clear that for all five website components, the overall rating for MakeMyTrip is 2 [8].

Data collected from the identified referral, rating sites is stored in the Excel sheet. This data contains various fields as 1) Date of review, 2) Reviewer name, 3) Gender, 4) Age, 5) Location, 6) Review, 7) Source, 8) Review rating, 9) Service and support, 10) Information depth, 11) Content, 12) User-friendly, 13) Time to load.

Fig. 6 shows an Excel document of reviews collected from different sources of online travel agents.

### A. Text Pre-Processing

The reviews were collected from various online platforms & referral rating sites need to have a series of text preprocessing before referring them for analysis using Text mining. The concept of text preprocessing consists of a series of stages, which include spelling in normalization, filtering, lemmatization. The various text preprocessing tasks are being essential and this includes content cleanup, tokenization, grammatical form tagging [9].

Table I represent various phases in determining positive negative reviews the process of determining sentiments consist of some of the essential processes like splitting words, POS tagging, lemmatization, joining and then sentiment analysis.

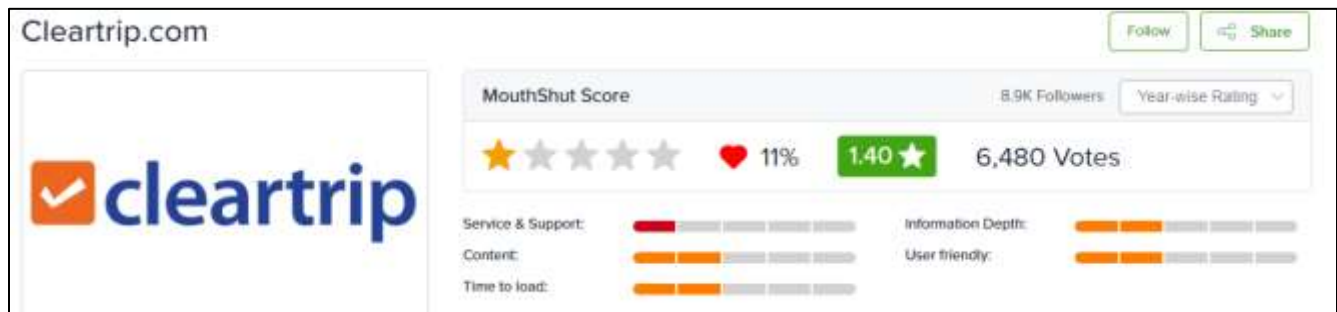


Fig. 3. The Overall Rating for Cleartrip (Source: Mouthshut.com).

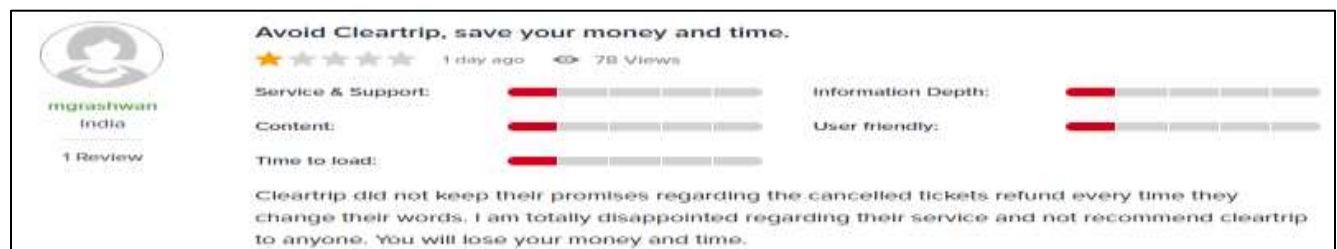


Fig. 4. Sample Review and Rating for Cleartrip (Source: Mouthshut.com).



Fig. 5. Overall Rating for MMT (Source: Mouthshut.com).

| # | Date       | Reviewer Name | Gender | Age | Location | Review  | Source        | Service       |
|---|------------|---------------|--------|-----|----------|---|---------------|---------------|
| 2 | 05-04-2018 | Sakani Man    | Female |     | India    | me in finding the best hotel with a good package which was Four Points By | mouthshut.com | Hotel/Booking |
| 3 | 05-05-2018 | Ishani Cool   | Female |     | India    | THEM CUSTOMER CARE IS VERY DULL AND OVERALL SERVICE IS VERY               | mouthshut.com | Booking       |

| Review Rating | Service & Support | Information Depth | Content | User Friendly | Time to Load | Overall Rating |
|---------------|-------------------|-------------------|---------|---------------|--------------|----------------|
| 4             | 5                 | 5                 | 5       | 5             | 4            | 5              |
| 1             | 1                 | 1                 | 1       | 1             | 1            | 1              |

Fig. 6. Different Fields in Format of Reviews Collected.

TABLE I. VARIOUS PHASES IN DETERMINING POSITIVE-NEGATIVE REVIEWS

| Process            | Output   |
|--------------------|--|
| Initial Review     | Horrible experience, hotels don't provide basic amenities and charge for everything. Customer service part is worst.   |
| Lowercase          | Horrible experience, hotels don't provide basic amenities and charge for everything. Customer service part is worst.   |
| Punctuation        | Horrible experience hotels don't provide basic amenities and charge for everything customer service part is worst  |
| Split words        | ['horrible', 'experience', 'hotels', 'don't', 'provide', 'basic', 'amenities', 'and', 'charge', 'for', 'everything', 'customer', 'service', 'part', 'is', 'worst']   |
| POS tagging        | [('horrible', 'JJ'), ('experience', 'NN'), ('hotels', 'NNS'), ('don't', 'VBP'), ('provide', 'VB'), ('basic', 'JJ'), ('amenities', 'NNS'), ('and', 'CC'), ('charge', 'NN'), ('for', 'IN'), ('everything', 'NN'), ('customer', 'NN'), ('service', 'NN'), ('part', 'NN'), ('is', 'VBZ'), ('worst', 'JJ')] |
| Lemmatization      | ['horrible', 'experience', 'hotel', 'don't', 'provide', 'basic', 'amenity', 'and', 'charge', 'for', 'everything', 'customer', 'service', 'part', 'be', 'bad']  |
| Joining            | horrible experience hotel don't provide basic amenity and charge for everything customer service part be bad   |
| Sentiment Analysis | {'negative': 0.333, 'neutral': 0.667, 'positive': 0.0, 'compound': -0.7906}  |
| Sentiment          | Positive / Negative  |

### III. RESULTS AND DISCUSSIONS

The results obtained from this research are presented in this section.

1) *Token identification*: Various subsections are being identified based on the results. After processing the textual reviews collected using text mining techniques following are the detailed list of tokens identified. The below section represents the top 25 tokens from the overall volume of reviews processed and also online agent-specific [10].

Tokens identified are used for deriving the overall perception of the travellers for the overall. Table II represents the top 25 tokens from both the categories positive and negative. The top 25 tokens based on the frequency of their occurrence has shown in the Table II For overall positive tokens care, support, help, these are the frequently used words by the reviewers while posting the reviews. Similarly, in the negative tokens category words like problems, fraud, cheat, and mistakes are some of the commonly used words by the reviews while expressing on social platforms.

Fig. 7 and 8 graphically represents the top 5 tokens in positive and negative category respectively

2) *Gender-specific (Male) Analysis of tokens*: In the gender-specific category of tokens obtained, the Table III represents the top tokens based on their frequency of occurrence in reviews that are listed. It is also identified from the Table III that when it comes to male support, care, friend, help these are the frequently used words while expressing in the form of review. And it denotes the satisfaction of the travellers the category of male.

Similarly, for or negative categories in males, the commonly used words based on the results obtained are a problem, fraud, Waste, mistake. This represents the dissatisfaction of travellers regarding online travel-related services. Tables III and IV represents the top 5 positive and top 5 negative tokens for males. Overall male analysis performed and results obtained indicate that usually while expressing on social platforms about satisfaction, travellers have used quite similar words.

TABLE II. OVERALL POSITIVE-NEGATIVE TOKENS

| Overall Positive Tokens |              |           |          | Overall Negative Tokens |               |           |          |
|-------------------------|--------------|-----------|----------|-------------------------|---------------|-----------|----------|
| Sr. No                  | Words        | Frequency | Polarity | Sr. No                  | Words         | Frequency | Polarity |
| 1                       | Care         | 489       | 0.49     | 1                       | Problem       | 244       | -0.4019  |
| 2                       | Support      | 454       | 0.40     | 2                       | Fraud         | 153       | -0.5859  |
| 3                       | Friend       | 267       | 0.49     | 3                       | Cheat         | 70        | -0.4588  |
| 4                       | Please       | 191       | 0.32     | 4                       | Mistake       | 65        | -0.34    |
| 5                       | Help         | 189       | 0.40     | 5                       | Cheater       | 60        | -0.5423  |
| 6                       | holiday      | 133       | 0.40     | 6                       | Fault         | 57        | -0.4019  |
| 7                       | thanks       | 109       | 0.44     | 7                       | Waste         | 50        | -0.4215  |
| 8                       | credit       | 90        | 0.38     | 8                       | Emergency     | 47        | -0.3818  |
| 9                       | Kind         | 88        | 0.53     | 9                       | Delay         | 45        | -0.3182  |
| 10                      | trust        | 65        | 0.51     | 10                      | Error         | 35        | -0.4019  |
| 11                      | thank        | 59        | 0.36     | 11                      | Loss          | 33        | -0.3182  |
| 12                      | solution     | 42        | 0.32     | 12                      | Complain      | 27        | -0.3612  |
| 13                      | promise      | 40        | 0.32     | 13                      | Trouble       | 25        | -0.4019  |
| 14                      | value        | 37        | 0.34     | 14                      | Penalty       | 25        | -0.4588  |
| 15                      | hope         | 30        | 0.44     | 15                      | Shock         | 22        | -0.3818  |
| 16                      | comfort      | 29        | 0.36     | 16                      | Hell          | 21        | -0.6808  |
| 17                      | profit       | 26        | 0.44     | 17                      | Fool          | 20        | -0.4404  |
| 18                      | resolve      | 25        | 0.38     | 18                      | inconvenience | 19        | -0.3612  |
| 19                      | hand         | 23        | 0.49     | 19                      | Doubt         | 19        | -0.3612  |
| 20                      | satisfaction | 22        | 0.44     | 20                      | Rude          | 16        | -0.4588  |

Overall male analysis performed and results obtained indicate that usually while expressing on social platforms about dissatisfaction, travellers have used quite similar words.

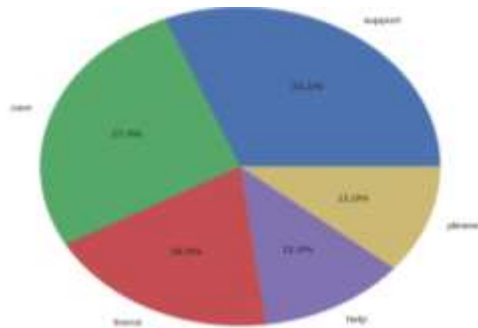


Fig. 7. Overall Positive Tokens (Top Five).

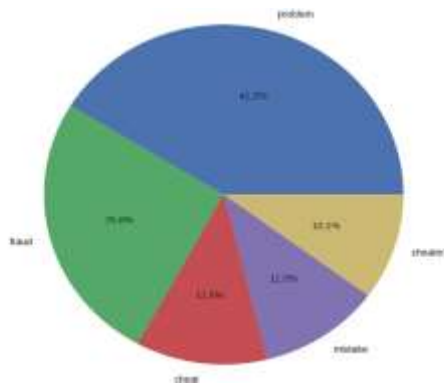


Fig. 8. Overall Negative Tokens (Top Five).

TABLE III. TOP FIVE POSITIVE MALE

| Overall Male Positive |         |           |          |
|-----------------------|---------|-----------|----------|
| Sr. No                | Words   | Frequency | polarity |
| 1                     | Support | 327       | 0.40     |
| 2                     | Care    | 284       | 0.49     |
| 3                     | Friend  | 199       | 0.49     |
| 4                     | Help    | 125       | 0.40     |
| 5                     | please  | 116       | 0.32     |

TABLE IV. TOP FIVE NEGATIVE MALE

| Overall Male Negative |           |           |          |
|-----------------------|-----------|-----------|----------|
| Sr. No                | words     | Frequency | Polarity |
| 1                     | problem   | 173       | -0.4019  |
| 2                     | Fraud     | 89        | -0.5859  |
| 3                     | Waste     | 34        | -0.4215  |
| 4                     | mistake   | 33        | -0.34    |
| 5                     | emergency | 33        | -0.3818  |

3) *Gender-specific (Female) Analysis of tokens:* This section describes various tokens used in the form of reviews while expressing on social platforms by female travellers. The analysis carried out in below Table V and VI represents the top 5 tokens based on their value of frequency. It also represents their importance while deriving positive and negative sentiments from the reviews hosted by female travellers. In a positive category, the frequently used word by female Travellers is Care, support, please, and so on. In negative categories, the most frequently identified words are problem fraud, cheat, a mistake.

Tables V and VI are the tabular representation of the top 5 reviews from gender-specific analysis of tokens in female traveller.

Overall female analysis performed and results obtained indicate that usually while expressing on social platforms about dissatisfaction, travellers have used quite similar words.

The analysis carried out in table VI represents the top 5 tokens based on their value of frequency for the positive female category.

Table VII represents a comparison matrix with the top ten tokens in each OTA. Based on the results shown, it is very much clear that satisfied travellers are frequently expressing with words like support, help, kind, comfort, and dissatisfied travellers are using words like problem, fraud, cheater, mistake while expressing on social platforms.

The comparison matrix shows the frequency of particular words when Traveller post their comments reviews on various sites. This comparison matrix also helps us understand comparative analysis between the keywords in the form of tokens for all the five online travel agents chosen in this research. Words like problem, help, and support are common in all the OTA.

4) *Keyword Analysis:* Following part of the research paper presents various keyword visualizations to understand more importance of each of the keyword plays in the entire

study. Fig. 9 explains lines of code referred from Python-based implementation which has helped to get the overall graphical representation of keywords and its analysis.

Fig. 10 represents the first Visualization of keywords for the real data values collected for the research purpose. From the results obtained, it very much clear that words like Hotel, book, ticket, service, call are the most critical words based on their appearances.

Fig. 11 and 12 represents the essential keywords in positive and negative category for the entire study. From the overall positive category, it is very much clear that words like best, awesome, trust, kind, super are the critical words describing positive sentiments in this category. In contrast, in overall negative word analysis. Words like fraud, bad, horrible, pathetic, failed is the crucial words that frequently occurred while expressing sentiments by the Travellers.

TABLE V. TOP FIVE POSITIVE FEMALE

| Overall Female Positive |         |           |          |
|-------------------------|---------|-----------|----------|
| Sr. No                  | Words   | Frequency | polarity |
| 1                       | Care    | 193       | 0.49     |
| 2                       | support | 115       | 0.40     |
| 3                       | please  | 73        | 0.32     |
| 4                       | friend  | 63        | 0.49     |
| 5                       | Help    | 62        | 0.40     |

TABLE VI. TOP FIVE POSITIVE FEMALE

| Overall Female Negative |         |           |          |
|-------------------------|---------|-----------|----------|
| Sr. No                  | words   | Frequency | Polarity |
| 1                       | problem | 68        | -0.4019  |
| 2                       | Fraud   | 60        | -0.5859  |
| 3                       | Cheat   | 34        | -0.4588  |
| 4                       | mistake | 31        | -0.34    |
| 5                       | cheater | 29        | -0.5423  |

TABLE VII. COMPARISON MATRIX FOR OTA'S AGAINST TOP 10 TOKENS

| Sr. No | Yatra    |            | redBus   |           | MMT      |           | Goibibo  |          | Cleartrip |          |
|--------|----------|------------|----------|-----------|----------|-----------|----------|----------|-----------|----------|
|        | Positive | Negative   | Positive | Negative  | Positive | Negative  | Positive | Negative | Positive  | Negative |
| 1      | care     | problem    | Care     | problem   | support  | problem   | care     | problem  | support   | fraud    |
| 2      | please   | fraud      | Support  | fraud     | care     | cheat     | support  | fraud    | care      | cheater  |
| 3      | support  | waste      | Friend   | mistake   | friend   | fraud     | friend   | cheat    | friend    | problem  |
| 4      | holiday  | cheat      | Please   | delay     | holiday  | mistake   | please   | cheater  | help      | cheat    |
| 5      | help     | emergency  | Help     | emergency | help     | error     | help     | waste    | please    | fault    |
| 6      | thanks   | mistake    | kind     | fault     | please   | emergency | thanks   | fault    | credit    | delay    |
| 7      | friend   | penalty    | thanks   | waste     | kind     | loss      | credit   | mistake  | thanks    | mistake  |
| 8      | credit   | fraudsters | comfort  | trouble   | thanks   | fault     | holiday  | fake     | thank     | cheating |
| 9      | kind     | Cheater    | thank    | cheater   | promise  | penalty   | trust    | loss     | trust     | doubt    |
| 10     | value    | Delay      | trust    | rude      | credit   | cheater   | kind     | fool     | holiday   | scam     |

```
wordcloud = WordCloud( background_color ='white')
wordcloud.generate_from_frequencies(frequencies=d1)
plt.figure()
plt.imshow(wordcloud, interpolation="bilinear")
plt.axis("off")
plt.savefig(abc+'_overall.png', dpi=600)
plt.show()
```

Fig. 9. Keyword Analysis LOC.



Fig. 10. Visualization of Overall Word Cloud.



Fig. 14. Cleartrip Dissatisfied Travellers Word Cloud.



Fig. 11. Overall Satisfied Travellers.

6) *Goibibo Keyword analysis:* Fig. 15 and 16, covers keywords based on their frequent occurrence for Goibibo. Based on the analysis carried out for Goibibo keywords excellent, best, great, trust, free are the crucial words in the positive category for Goibibo. Whereas worst, bad, fraud, hate is the commonly observed keywords in the negative category for Goibibo.



Fig. 12. Overall Dissatisfied Travellers.

7) *MMT Keyword Analysis:* MakeMyTrip (MMT) keyword analysis is expressed in Fig. 17 and 18 in graphical representation. In a positive category, best, free, kind, happy, wonderful are critical words for MMT. In contrast, when it comes to negative category pathetic, bad, worst, horrible, fraud, these are the words carrying much importance in MMT keyword analysis.



Fig. 15. Goibibo Satisfied Travellers Word Cloud.

The section below describes the online travel agent (OTA) Specific keyword analysis.

5) *Cleartrip Keyword analysis:* Fig. 13 and 14 depicts the visualization of keywords in Cleartrip. Based on the results obtained, it is very much clear that words like best, trust, great, happy are the most critical in a positive category for Cleartrip. In contrast, Fig. 13 shows bad, fraud, pathetic, fail other common words in negative categories for Cleartrip.



Fig. 13. Cleartrip Satisfied Travellers Word Cloud.



Fig. 16. Goibibo Dissatisfied Travellers Word Cloud.



Fig. 17. MMT Satisfied Travellers Word Cloud.



Fig. 22. Yatra Dissatisfied Travellers Word Cloud.



Fig. 18. MMT Dissatisfied Travellers Word Cloud.

8) *redBus Keyword Analysis*: Similar to other online travel agents, keyword analysis is performed for redBus. Fig. 19 and 20 depicts various important words in positive and negative categories for redBus. Best, happy, kind, comfortable these are the top-rated words in the positive category, and bad, worst, horrible these are some of the highly-rated words in the negative category for redBus.

9) *Yatra Keyword Analysis*: Fig. 21 and 22 is the word count analysis done on online reviews collected for Yatra online travel agent. The outcome of the analysis demonstrates the words like best, great, kind, splendid, love is the essential words in the positive category for the Yatra. In the negative category, analysis says the words like bad, fraud, worst, unprofessional, horrible are some of the words that play an essential role in Yatra.



Fig. 19. redBus Satisfied Travellers Word Cloud.



Fig. 20. redBus Dissatisfied Travellers Word Cloud.



Fig. 21. Yatra Satisfied Travellers Word Cloud.

10) *Analysis based on Rating*: Each of the review collected from the Mouthshut.com has a rating value attached to it ranging from 1 as Low to 5 as High. Following is the OTA's and reviews based on rating value percentage Fig. 23 covers lines of code for getting the statistical value of review ratings for all the reviews which are being collected and processed and relevant graphical representations are obtained with the help of similar lines of code.

The Table VIII represents a statistical analysis of reviews regarding the overall rating given by the reviewer on the scale of 1 to 5, where 1 is low, and 5 is high. The table shows a good number of reviews for each of the online travel agents as review rating low, which means dissatisfaction of travellers is low. At the same time the volume of reviews classified under 5 represents the high level of satisfaction of travellers using an online travel agent and there relevant services MakeMyTrip, Goibibo and Cleartrip are having good volume of reviews under the highest rating.

```

sum = []
percent = []
for i in range(5):
    sumval = 0
    for j in range(5):
        sumval += len(ratingList[j][i])
    sum.append(sumval)
for i in range(5):
    percent = []
    for j in range(5):
        percent.append(round(((len(ratingList[j][i]) * 100) /
sum[i] ), 1 ))
    worksheet.write(i, j, round(((len(ratingList[j][i]) * 100) /
sum[i] ), 1 ))
    worksheet.write(i, 5, sum[i])
percent.append(sum[i])
print(percent)
workbook.close()

```

Fig. 23. Rating based Analysis LOC.



TABLE VIII. OTA'S REVIEW RATING

| Review Rating |           | 1   | 2  | 3  | 4  | 5  | Total Reviews |
|---------------|-----------|-----|----|----|----|----|---------------|
| OTA's         | Cleartrip | 236 | 13 | 12 | 33 | 57 | 351           |
|               | Goibibo   | 321 | 33 | 36 | 54 | 75 | 519           |
|               | MMT       | 281 | 18 | 48 | 90 | 94 | 531           |
|               | redBus    | 292 | 38 | 45 | 46 | 48 | 469           |
|               | Yatra     | 172 | 17 | 15 | 40 | 21 | 265           |

Table IX represents the statistical analysis of review ratings for all the OTA in percentage collected and processed for each of these online travel agents. From the table, it is very much clear that MakeMyTrip is having the highest percentage of review rating value with 5. It means the satisfaction level of travellers using various services offered by MakeMyTrip is also high in comparison to other online travel agents.

Cleartrip has 16.2% reviews based on the overall rating given to the review this percentage and it's second after MakeMyTrip. Goibibo is third in the list with the highest percentage under 5. It is very much clear from table values that Yatra is having low preference under the top-level, which means satisfaction is very low for the Yatra.

TABLE IX. OTA'S REVIEW RATING PERCENTAGE

| Review Rating |           | 1    | 2   | 3   | 4    | 5    | Total volume |
|---------------|-----------|------|-----|-----|------|------|--------------|
| OTA           | Cleartrip | 67.2 | 3.7 | 3.4 | 9.4  | 16.2 | 351          |
|               | Goibibo   | 61.8 | 6.4 | 6.9 | 10.4 | 14.5 | 519          |
|               | MMT       | 52.9 | 3.4 | 9   | 16.9 | 17.7 | 531          |
|               | redBus    | 62.3 | 8.1 | 9.6 | 9.8  | 10.2 | 469          |
|               | Yatra     | 64.9 | 6.4 | 5.7 | 15.1 | 7.9  | 265          |

Looking at Table IX for Yatra, 64.9 percent of Travellers have given Low review rating, represent the percentage of the satisfaction of Travellers using services offered by Yatra.

IV. ARCHITECTURAL VIEW OF MODEL

In this research work, netnography and Text mining techniques are used as an integrated approach. The following are the various stages in the architectural view of the basic model in this study.

Fig. 24 represents the architectural view of the platform used in the text mining process. As mentioned in the diagram, the following are the steps that help get a visualized representation of specific results. The first stage in the architectural view is inputting of the review. In this research study, the reviews or posts used are collected from various online platforms using Netnographic guidelines. Collected reviews are maintained in a repository. The second stage is the preprocessing stage. In this stage, preprocessing of the text reviews collected is done using various techniques of splitting the words, converting the words, checking for the missing values is performed. The next stage is cleaning the reviews; in this process of cleaning of the reviews, it is very much essential that some of the processes which help the text to

refine and get a meaningful text for further processing. Sentiment score stage, the sentiments of each word, and then for the entire review are obtained. In this stage, using the Vader sentiment lexicon dictionary. The sentiments are derived as positive, negative, neutral—compound sentiment score help to classify the review into positive or negative. In the review classification stage, based on the compound sentiment scores with more than 0.05 compounded courses are classified as positive, and the sentiment scores between -0.05 and 0.0 are considered neutral review. The review has a compounded score of less than -0.05 classified as negative reviews. In the visualization stage, based on the review, the various analysis has been carried out in the graphical representation has been generated in the form of the word cloud, pie charts.

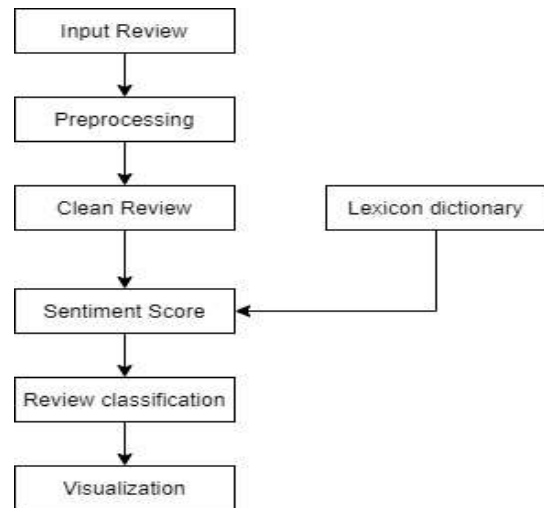


Fig. 24. Architectural view of Approach used for Visualization.

V. STATISTICAL ANALYSIS OF RESULT

In this section, the statistical analysis results obtained using the Chi-square test are discussed. A Chi-square test for various parameters is performed to understand and test the relationship between various variables based on the results of the test acceptance or rejection of the hypothesis is being chosen. Performed Chi-square test on the following scenarios:

1) Scenario 1(OTA vs. Gender): Here, the Chi-square statistical test is used to test whether two variables are independent or not. Chi-square hypothesis testing is used to understand whether there is any relationship between the online travel agent and the gender of travellers.

a) Null hypothesis (H0): No association/relationship between gender of traveller and OTA.

b) The alternate hypothesis (H1): Gender of traveller and OTA has an association/relationship

Table X(a) and X(b) shows the calculations of test statistics.

P-value obtained is .0000009307. With the degree of freedom 4 and at 5% level of significance, the critical value/tabular value is 9.49. The calculated chi-square value(33.54) is greater than the critical or table value(9.49).

TABLE X. (A) TEST STATISTICS FOR OTA VS. GENDER

| Gender<br>OTA | Male     |          |       | Female   |          |       | Total    |       |
|---------------|----------|----------|-------|----------|----------|-------|----------|-------|
|               | Observed | Expected | %     | Observed | Expected | %     | Observed | %     |
| Cleartrip     | 242      | 233.38   | 16.74 | 110      | 118.62   | 14.97 | 352      | 16.14 |
| Goibibo       | 429      | 385.87   | 29.67 | 153      | 196.13   | 20.82 | 582      | 26.69 |
| MMT           | 325      | 341.44   | 22.48 | 190      | 173.56   | 25.85 | 515      | 23.61 |
| redBus        | 306      | 310.95   | 21.16 | 163      | 158.05   | 22.18 | 469      | 21.5  |
| Yatra         | 144      | 174.37   | 9.96  | 119      | 88.63    | 16.19 | 263      | 12.06 |
| <b>Total</b>  | 1446     |          | 100   | 735      |          | 100   | 2181     | 100   |

(b) Test Statistics for OTA vs. Gender

| Observed                    | Expected | (obs-exp)^2 | (obs-exp)^2/ Exp |
|-----------------------------|----------|-------------|------------------|
| 242                         | 233.38   | 74.30       | 0.32             |
| 429                         | 385.87   | 1860.20     | 4.82             |
| 325                         | 341.44   | 270.27      | 0.79             |
| 306                         | 310.95   | 24.50       | 0.08             |
| 144                         | 174.37   | 922.34      | 5.29             |
| 110                         | 118.62   | 74.30       | 0.63             |
| 153                         | 196.13   | 1860.20     | 9.48             |
| 190                         | 173.56   | 270.27      | 1.56             |
| 163                         | 158.05   | 24.50       | 0.16             |
| 119                         | 88.63    | 922.34      | 10.41            |
| Calculated Chi-Square Value |          |             | 33.54            |

There is enough of the statistical evidence to reject the null hypothesis and to accept the fact that there is an association or relationship between various OTAs and the gender of travellers.

Since  $33.52 > 9.48$  or our P-value  $< 0.05$

The alternate hypothesis is accepted. Thus there is an association between gender and people using OTA.

Fig. 25 is the graphical representation of observed values for males and females for all the online travel agents against the total volume of reviews collected.

2) *OTA vs. Positive Negative Sentiments*: In this section, the Chi-square statistical test performed to understand the relationship between positive negative sentiments of travellers and various online travel-related service providers (OTA).

a) *Null hypothesis (H0)*: No association between Sentiment and people using OTA.

b) *Alternate Hypothesis (H1)*: Association between Sentiment and people using OTA.

From Table XI, P-value obtained using the chi-square test is 0.01, chi-square value from P-value is 12.78; critical or tabulated chi-square value at the degree of freedom 4 and level of significance 5% is 9.49; since  $12.78 > 9.48$  An alternate hypothesis is accepted; it means an association or relationship between sentiments and people using OTA. Fig. 26 is a graphical representation of about table values representing the

relationship between OTA Sentiment against the total number of reviews for the online travel agents.

3) *Statistical Analysis of Positive negative ratio Vs. Gender*: In the below section, a detailed analysis of the relationship between positive-negative sentiment ratio and gender. Based on the data shown in the Table, P-value using the Chi-square statistical test is 0.00. The following are the null and alternate hypotheses framed to identify whether there is any relationship between the positive-negative sentiment ratio and gender.

a) *Null hypothesis (H0)*: No association between Sentiment and Gender

b) *Alternate Hypothesis (H1)*: Association between Sentiment and Gender

Referring to Table XII, The Chi-square value calculated using P-value is 12.96. And with a 5% level of significance and 4 degrees of freedom, the tabular value for chi-square is 9.48. Since the calculated value, 12.96 is more than the tabular value for chi-square 9.48. There is enough evidence to reject the null hypothesis and accept the alternate hypothesis, which says there is a relationship between Sentiment and gender. Since  $12.95 > 9.48$ , Alternate hypothesis is accepted.

Fig. 27 represents a detailed classification of the total volume of reviews and sentiments as a positive and negative gender-wise. The finding from the analyses says that the male and female total percentage of negative review sentiment in male is less than female.

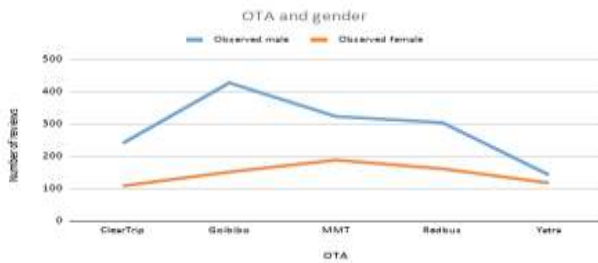


Fig. 25. OTA Specific Gender Values.

TABLE XI. TEST STATISTICS FOR OTA VS. POSITIVE NEGATIVE SENTIMENTS

| Gender OTA   | Positive    |          |            | Negative    |          |            | Total Observed |
|--------------|-------------|----------|------------|-------------|----------|------------|----------------|
|              | Observed    | Expected | %          | Observed    | Expected | %          |                |
| Cleartrip    | 165         | 178.5    | 14.92      | 187         | 173.5    | 17.4       | 352            |
| Goibibo      | 303         | 295.14   | 27.4       | 279         | 286.86   | 25.95      | 582            |
| MMT          | 291         | 261.16   | 26.3       | 224         | 253.84   | 20.84      | 515            |
| redBus       | 225         | 237.83   | 20.34      | 244         | 231.17   | 22.7       | 469            |
| Yatra        | 122         | 133.37   | 11.03      | 141         | 129.63   | 13.12      | 263            |
| <b>Total</b> | <b>1106</b> |          | <b>100</b> | <b>1075</b> |          | <b>100</b> | <b>2181</b>    |

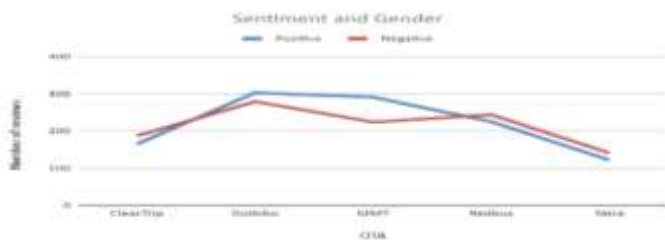


Fig. 26. Sentiments and Gender Relationship for OTA.

TABLE XII. RELATIONSHIP BETWEEN SENTIMENTS AND GENDER

| Gender Sentiment | Observed    |            |             |            | Expected |        | Total       |
|------------------|-------------|------------|-------------|------------|----------|--------|-------------|
|                  | Male        | Male %     | Female      | Female %   | Male     | Female |             |
| Positive         | 773         | 69.89      | 673         | 62.60      | 733.28   | 712.72 | 1446        |
| Negative         | 333         | 30.11      | 402         | 37.40      | 372.72   | 362.28 | 735         |
| <b>Total</b>     | <b>1106</b> | <b>100</b> | <b>1075</b> | <b>100</b> |          |        | <b>2181</b> |

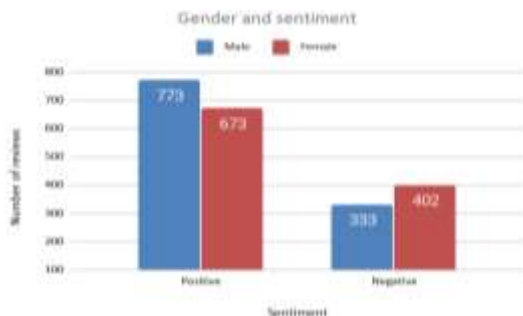


Fig. 27. Classification Reviews and Sentiments.

## VI. SERVICE SPECIFIC ANALYSIS

Online travel agents used in this research provides a various set of services to the users, customers, travellers. These services include hotel booking, flight booking, bus booking. In this section, the results obtained regarding these services are evaluated, and reasonable interpretations are derived. Fig. 28 shows lines of code help to do the analysis with reference to sentiments for each of the OTA specific to services identified.

```

xpos = np.arange(len(OTA))
plt.bar(xpos-0.2, posi, width=0.4, label="Positive")
plt.bar(xpos+0.2, negi, width=0.4, label="Negative")
plt.xticks(xpos, OTA)
plt.ylabel("Number of "+word+" reviews")
plt.title("Sentiment of "+word+" reviews")
plt.legend()
plt.savefig(word+"_OTA")
    
```

Fig. 28. Service Specific Sentiment Analysis LOC.

The graphical representation below Fig. 29 clearly indicates for bus-related bookings our services Travellers have preferred redBus.

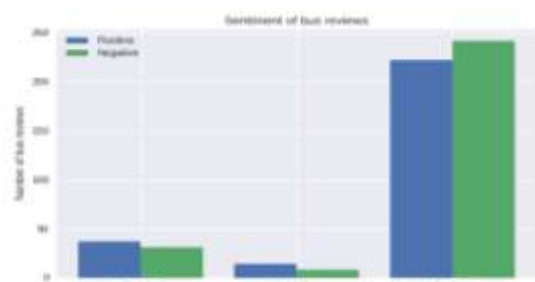


Fig. 29. Sentiments of Bus Reviews for OTA.

After redBus Travellers preferred Goibibo and then Makemytrip, the finding here is for bus-related services, redBus is on top of all.

Flight-related services are offered by the following online travel agents chosen for the study, Cleartrip, Goibibo, MakeMyTrip, and Yatra. The graphical representation Fig. 30 shows that for the flight-related services, the volume of negative reviews and sentiments are more with Cleartrip than other service providers. For Goibibo, the positive sentiments are little more than the negative sentiments which we can interpret in a manner that Travellers or customers are a little happier regarding flight-related services offered by Goibibo. Similarly, for MakeMyTrip, positive sentiments are more than negative sentiments for all the possible reviews processed for flight-related services which also represents the satisfaction of travellers is more than dissatisfaction for Makemytrip. For Yatra, the positive sentiments are less than the negative sentiments. Also, the volume of reviews for Yatra regarding flight-related services is low; it represents the satisfaction is

lower than the dissatisfaction and preference for flight-related services is not towards Yatra.

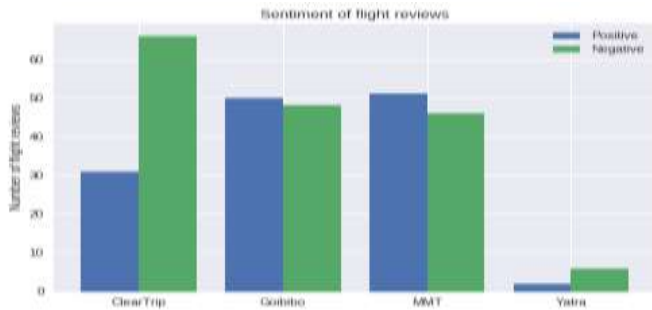


Fig. 30. The Sentiment of Flight Reviews for OTA's.

The following are the online travel agents that offer hotel booking related services online Cleartrip, Goibibo, MakeMyTrip, Yatra. Fig. 31 clearly explains the results of the analysis done regarding Hotel related services offered by these online travel agents in India. When it comes to Cleartrip Hotel related services, Fig. 31 represents the negative sentiments are more than positive it means when it comes to hotel-related services, online dissatisfaction is more observed in customers or travellers about Cleartrip. Goibibo looks at the top choice for Hotel related services online. Again based on the analysis, it is proved that positive sentiments are more than negative sentiments. It means the satisfaction of travellers using Hotel related services offered by Goibibo is high. MakeMyTrip also follows the line of Goibibo when it comes to hotel-related services. It has become the second preference for Hotel related services after Goibibo. It is also clearly indicated that Yatra is rarely preferred by travellers when it comes to online Hotel related services.

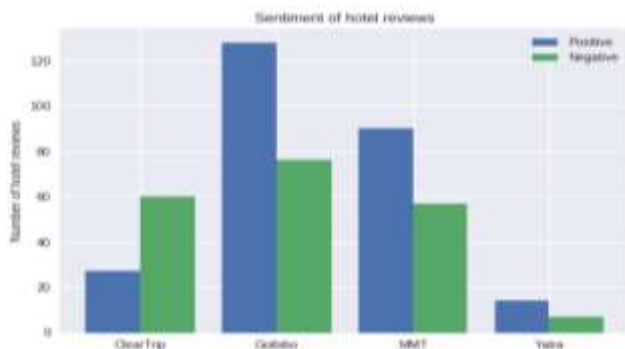


Fig. 31. The Sentiment of Hotel Reviews for OTA's.

## VII. CONCLUSION

Customer is King for a successful business, customer satisfaction is considered to be a prime attribute. This research helps in consumer decision making while choosing various travel-related services and also businesses to choose their strategies to improve on services offered based on listening to the consumers.

Netnography and Text mining techniques used to perform analysis and processing of reviews and comments collected from various online platforms. The review collection has followed all the guidelines suggested by netnographic studies for five online travel agents i.e. MakeMyTrip (MMT), Yatra,

Cleartrip, Goibibo, redBus analysis is performed using text mining as an integrated approach.

Based on results obtained regarding service-specific analysis, it is possible to conclude that for flight booking, related services traveller's satisfaction is more with MMT, Goibibo. Consumers have preferred MakeMyTrip as the first option and Goibibo as Second and followed by Yatra. For Hotel booking related services satisfaction of travellers using services from MMT, Goibibo is high compared to others. For bus booking related services, users are more happy and satisfied with redBus then Goibibo and MMT.

This research work has used various approaches concerning the understanding of satisfaction, dissatisfaction, or the perception of travellers using selective travel-related online service providers in India to build a competitive advantage to customers and company too. Online travel service providers in India can make appropriate decisions according to the results obtained from the study regarding customer perception and create a competitive advantage. Also, consumers can choose a particular online travel-related service provider based on the result of the study for improved services and traveller satisfaction.

## REFERENCES

- [1] C. Science and M. Studies, "A Study on Tourist Satisfaction towards Hotel Related Services in Gujarat Tourism," vol. 4, no. 6, pp. 105–110, 2016.
- [2] "Internet users in India to cross 500 mn in 2016: Prasad," May 2016.
- [3] D. Mane and D. Bhatnagar, "Integrating Netnographic Analysis and Text Mining for Assessing Satisfaction of Travellers Visiting to India - A Review of Literature," 2020, pp. 564–572.
- [4] "Top Online Travel Agencies (OTAs) in India - THE 'PERFECT' HOTEL." [Online]. Available: <http://whiteskyhospitality.co.uk/top-online-travel-agencies-otas-in-india/>. [Accessed: 28-Jul-2020].
- [5] J. Chakravarthi and V. Gopal, "Comparison of traditional and online travel services: A concept note," IUP J. Bus. Strateg., vol. 9, no. 1, pp. 45–59, 2012.
- [6] D. Mane., "Text Mining to Understand Major Keywords Explaining Sentiments of Travelers Using Travel Related Online," books.google.com.
- [7] "CLEARTRIP.COM Reviews, Feedback, Complaint, Experience, Customer Care Number - MouthShut.com." [Online]. Available: <https://www.mouthshut.com/websites/Cleartrip-com-reviews-925062909-srch>. [Accessed: 28-Aug-2020].
- [8] "MAKEMYTRIP.COM Reviews, Feedback, Complaint, Experience, Customer Care Number - MouthShut.com." [Online]. Available: <https://www.mouthshut.com/websites/MakeMyTrip-com-reviews-925031929-srch>. [Accessed: 28-Aug-2020].
- [9] P. Gupta, R. Tiwari, and N. Robert, "Sentiment Analysis and Text Summarization of Online Reviews: A Survey," 2016 Int. Conf. Commun. Signal Process., pp. 241–245, 2016.
- [10] Prameswari P., Surjandari I., and Laoh E.. Opinion Mining from Online Reviews in Bali Tourist Area. 226–230. (2017).
- [11] Kozinets R., Netnography: doing ethnographic research online. Sage Publication, London (2010b).
- [12] Top 10 travel Companies in India|TopEcommerceStartups., Retrieved from <http://topecommercestartups.com/top-10-travel-companies-in-india/Travel&#amp;>; (2014).

# Multi-Dimensional Fraud Detection Metrics in Business Processes and their Application

Badr Omair<sup>1</sup>, Ahmad Alturki<sup>2</sup>

Faculty of Computer and Information Sciences  
King Saud University, Riyadh, Kingdom of Saudi Arabia

**Abstract**—Occupational fraud is defined as the deliberate misuse of one’s occupation for personal enrichment. It poses a significant challenge for organizations and governments. Estimates indicate that the funds involved in occupational fraud cases investigated across 125 countries between 2018 and 2019 exceeded US\$3.6 billion. Process-based fraud (PBF) is a form of occupational fraud that is perpetrated inside business processes. Business processes underlie the logic of the work that organizations undertake, and they are used to execute an organization’s strategies to achieve organizational goals. Business processes should be examined for potential fraud risks to ensure that businesses achieve their objectives. While it is impossible to prevent fraud entirely, it must be detected. However, PBF detection metrics are not well developed at present. They are scattered, unstandardized, not validated, and, in some cases, absent. This study aimed to develop a comprehensive PBF detection metric by leveraging and operationalizing a taxonomy of fraud detection metrics for business processes as an underlying theory. 41 PBF detection metrics were deduced from the taxonomy using design science research. To evaluate their utility, the application of the metrics was undertaken using illustrative scenarios, and a real example of the implementation of the metrics was provided. The developed metrics form a complete, classified, validated, and standardized list of PBF detection metrics, which include all the necessary PBF detection dimensions. It is expected that the stakeholders involved in PBF detection will use the metrics established in this work in their practice to increase the effectiveness of the PBF detection process.

**Keywords**—Business process fraud; fraud detection; fraud indicators; fraud measures; fraud metrics; PBF; red flags

## I. INTRODUCTION

Fraud refers to an action that is designed to deceive others. Fraud results in a loss for the victim and gain for the perpetrator [1]. The Association of Certified Fraud Examiners (ACFE)<sup>1</sup> defines occupational fraud as the “use of one’s occupation for personal enrichment through the deliberate misuse or misapplication of the employing organization’s resources or assets” [2, p. 86]. Organizations and individuals alike can be financially or physically affected by fraud [3].

Fraud can either be internal, when it is committed by someone inside an organization, or external, when it originates from outside an organization [4]. In this research, the focus is on internal or occupational fraud.

Fraud is becoming a globally prevalent threat [5]. It is estimated that the overall loss resulting from 2,504 cases of occupational fraud that were investigated between January 2018 and September 2019 exceeded US\$3.6 billion across 125 countries [2]. The ACFE estimates that organizations lose approximately 5% of their revenues to fraud each year [2]. The wave of financial scandals that has been sweeping the world in the current century has also heightened the awareness of the need to manage fraud risk [6].

Process-based fraud (PBF) is a form of fraud that occurs in business processes. It can be identified by measuring the deviation from the process model [7]. However, deviation in the business process model is not always regarded as fraud; in order to confirm that fraud has taken place, a domain expert must investigate the matter.

Business process refers to a collection of related events, activities, decision points, actors, and objects that lead to an outcome that is valuable to at least one customer [8]. Business processes are core assets of organizations [8], and they are essential in the implementation of organizational strategy [9]. Business processes should be examined to detect any associated potential fraud risks that may threaten the achievement of business objectives [10]. However, at present, PBF detection metrics are not well addressed [11]. They are incomplete, overlapping, scattered, and not standardized [11]. Furthermore, the increase in fraud in recent years reflects the persistent nature of the issue [12]. Therefore, as it is impossible to prevent PBF completely, detecting it when it occurs is essentially.

This manuscript aims to develop comprehensive metrics that cover all the components necessary for the effective detection of PBF. The developed metrics will contribute to the effective detection of PBF as they provide a comprehensive, validated, and standardized list of PBF detection metrics.

First, the metrics are deduced from the taxonomy of fraud detection metrics for business processes [13]. The taxonomy serves as the underlying theory using design science research (DSR). The use of this taxonomy provides a complete understanding of PBF detection, coverage of all PBF detection elements, and a checklist of best practices that define PBF detection metrics [13]. Second, an illustrative scenario, as an evaluation method [14], is provided for each of the developed metrics in order to validate their utility. Ultimately, an implementation that uses the process mining technique is proposed to demonstrate the technical application of the metrics.

<sup>1</sup> <https://www.acfe.com>

The remaining contents of this paper are organized as follows: Section II provides the background of the topic; Section III explains the methodology followed in the current work; Section IV proposes the complete PBF detection metrics; Section V provides a real example of the implementation of the metrics; Section VI shows and discusses the results; and, finally, in Section VII, the conclusions and direction in which the work in this field may progress in the future are presented.

## II. BACKGROUND

Implementing fraud detection and fraud prevention systems is essential for effective fraud risk management [15]. Fraud prevention consists of measures to avoid or reduce fraud. In addition, in fraud detection, measures that help identify fraud when it occurs are used [15]. Since preventing every instance of fraud is impossible, continuous application of fraud detection techniques is necessary to protect against any instances that were not prevented [3].

Fraud detection techniques can be placed into one of three categories [16]. First, the misuse-based detection technique uses a predefined list (i.e., known patterns) of possible fraud schemes to detect fraud. It is an expert fraud detection system that uses predefined metrics. Its advantage is a low false alarm rate, but it cannot detect instances of fraud that follow new patterns [16]. Second, the anomaly-based technique can be implemented using machine learning techniques, which leads to the detection of any suspicious behavior that deviates from standard behavior [17], [18]. It does not require a predefined list of fraud schemes, and it can detect new cases of fraud. However, it suffers from a high false alarm rate [19]. Third, the hybrid technique attempts to combine the previous two techniques to overcome their limitations [16].

Successful fraud detection must include an examination of business processes to identify the potential origins of fraud [20]. Business processes are the core of business process management (BPM), which is a management discipline that uses business processes to implement organizational strategy [9]. It is a management discipline that requires continued focus, and often, significant changes in management style [9].

PBF detection metrics form the intersection between fraud risk management and BPM, as reflected in the bidirectional arrow mentioned in Fig. 1. The use of such metrics is common in fraud risk assessment and process monitoring and control<sup>2</sup> which are elements of fraud risk management and BPM, respectively.

Fraud detection can be achieved using a taxonomy to predefine initial fraud schemes [1], [21]. A taxonomy is a set of dimensions, each consisting of a set of mutually exclusive and collectively exhaustive characteristics [22]. A taxonomy of fraud detection metrics for business processes was proposed

<sup>2</sup> Performance measures are usually identified during the process analysis phase of BPM. In some cases, they are identified during the process identification phase [8]. Moreover, business process measures can be classified as measures for business process models and execution [63]. Since fraud detection is the goal, this study focuses on measures executed in the process monitoring and control phase to determine how well the executed processes work with regard to the chosen measures [9].

in [13], as depicted in Fig. 2. The taxonomy provides a holistic view of fraud detection in business processes. It consists of the dimensions examined in the following subsections.

### A. Fraud Domain

This dimension covers the application domain of fraud detection. Knowing the fraud domain is crucial in the detection of fraud because it allows an understanding to be gained of the problem domain [23]. In addition, specific fraud, which is particular to certain domains, exists, and these cases require special handling. This dimension contains two characteristics:

- *General*: Describes all metrics that can be used in any application domain.
- *Specific*: Covers a particular application domain, such as finance.

### B. Fraud Data Scheme

This dimension covers all the potential fraud schemes in the data. Fraud data schemes provide a list of possible data schemes used for committing fraud, which means that understanding them is critical for detection. This dimension contains the following data schemes:

- *Anomalous*: Covers any data that can be characterized as ambiguous or exceptional (e.g., too long, too short, excessive, and outliers).
- *Discrepant*: Describes inconsistent data (e.g., the conflict between input and output, and between past and current).
- *Missing*: Covers insufficient and absent data.
- *Wrong*: Covers incorrect data (e.g., inaccurate, non-conforming, fictitious, error, and outdated).

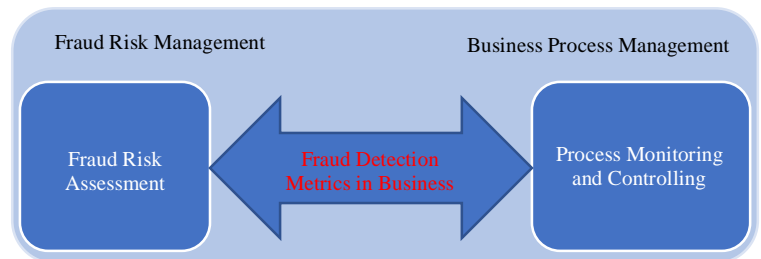


Fig. 1. Execution Scope of PBF Detection Metrics.

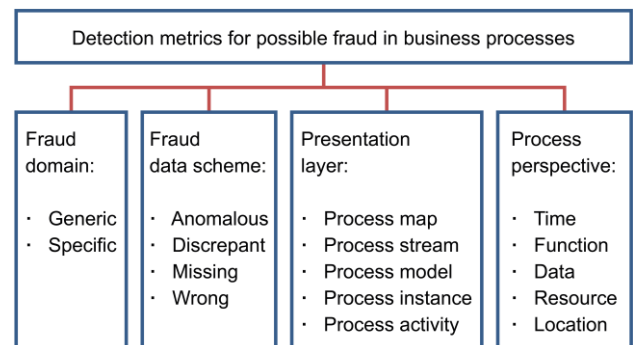


Fig. 2. Taxonomy of Fraud Detection metrics in Business Processes. Source: [13].

### C. Presentation Layer

This dimension aims to examine all layers of the business processes, as illustrated in Fig. 3.<sup>3</sup> The layers are essential for detecting fraud because every layer can give specific auditing information [24]. Additionally, some fraud cases do not become apparent by looking at a single layer. The dimension contains the following characteristic layers:

- *Process map*: Gives an overview of all business processes and determines their relationships. The process map also contains aggregated data on all business processes in the organization. It is useful for planning fraud detection in business processes.
- *Process stream*: Offers a greater level of detail compared to the process map. It helps set the scope by focusing on a collection of processes that form a specific (and usually vital) business cycle, such as the purchase-to-pay cycle. This layer allows fraud examiners to aggregate data on a particular business cycle.
- *Process model*: Represents a single business process, such as the payment process. It provides more detail on the structure of the process, its controls, activities, and actors. This layer contains aggregated data on many instances of a specific business process.
- *Process instance*: Depicts the details of one particular instance of a process model. It contains concrete data on one specific business process instance, such as payment instance number 123.
- *Process activity*: This is the lowest layer in the *presentation layer* dimension. It can be considered an element of the *process instance* layer with a particular focus. It gives concrete data with more detail on a specific activity in a particular process instance, such as approval activity.

### D. Process Perspective

This dimension looks at business process from various angles because, for successful fraud detection, it is necessary to examine all aspects of business process [20]. This dimension contains the following characteristics:

- *Time*: This perspective regards business process's time (e.g., throughput time, actual processing time, waiting time, and deadlines).
- *Function*: This perspective is concerned with the implementation of the activities in business process (e.g., work frequency, work sequence, work decision, process steps, and process control flow).
- *Data*: The data perspective covers all the data that are entered, consumed, and delivered by business process (e.g., process objects).
- *Resource*: This perspective involves all the actors that interact with business process, including customers,

software, business role, business units, suppliers, and employees.

- *Location*: This perspective is concerned with the location of business process's execution.

The results of the literature review on PBF detection metrics<sup>4</sup> are summarized in Fig. 4 in the form of a literature map [11]. The literature map illustrates the topics relevant to fraud detection metrics in business processes, as well as the frequency of their recurrence in the literature. Omair and Alturki [11] demonstrated that, at present, the explicitly defined PBF detection metrics, which are listed in Table I, do not adequately address the essential conceptual perspectives of business process.

Combined metrics and process mining can improve fraud detection [25]. Process mining is a methodology that aims to discover, monitor, and improve real processes by analyzing their event logs [26]. It connects model-based process analysis (e.g., simulation) and data-oriented analysis techniques (e.g., data mining) [27]. Process mining associates the actual processes with their data and the process models [28].

Process mining has been successfully applied to detect fraud [23], [29]–[31]. It can reveal fraudulent transactions that cannot be detected using traditional audit methods [29], [32], [33]. Relying on measurements of throughput processing (not just measurements of the input-output relation), process mining can identify a problem's root cause. This involves identifying the process model, and, subsequently, the performance of the process [34].

Using process mining to detect fraud has many advantages. Since event logs are automatically logged in most existing systems [35], it is possible to save time and effort, and to improve detection accuracy by taking real and complete data as opposed to samples [36]. Also, reading from event logs ensures independence from human intervention, which guarantees unaltered and error-free data [37]. According to the ACFE report [2], the median time for detecting fraud is 14 months. During the interval between occurrence and detection, the most significant financial losses tend to occur. However, using online process mining solutions can change this reality [38].

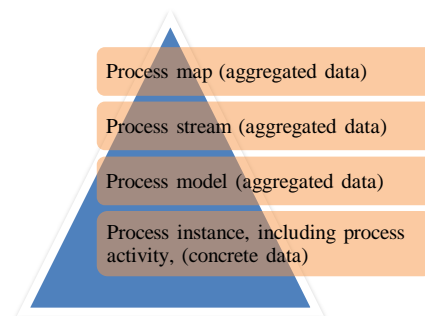


Fig. 3. Presentation Layers of Business Processes.

<sup>3</sup> For more information, see [13], [24].

<sup>4</sup> For the complete literature review and analysis, see [11].

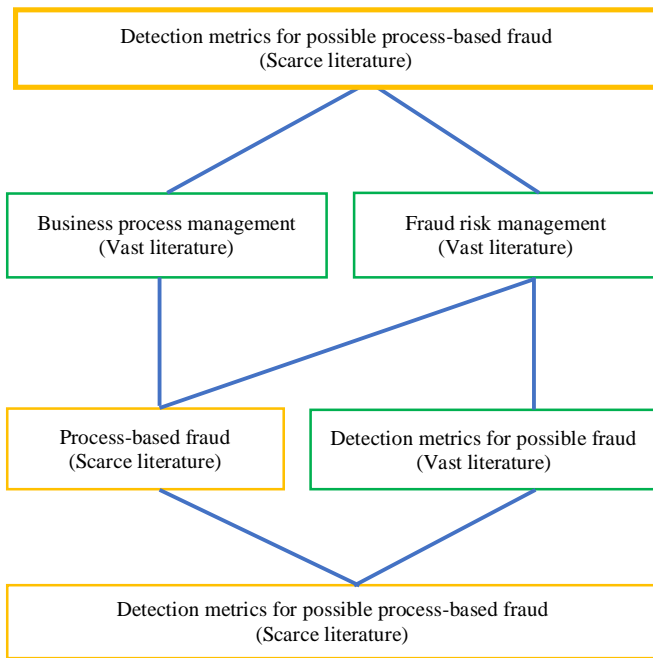


Fig. 4. Literature Map. Source: [11].

TABLE I. FRAUD DETECTION METRICS IN BUSINESS PROCESSES. ADAPTED FROM [24]

| ID | Metric name           | Explanation  | Reference                  |
|----|-----------------------|--|----------------------------|
| 1  | Skipped activity      | Not executing an activity that is prescribed in the standard operating procedure (SOP). The skipped activity is either a routine activity or a decision activity [42].   | [7], [30], [31], [42]–[45] |
| 2  | Wrong resources       | The activity is performed by an actor who is not defined in the SOP.   | [7], [30], [31], [42]–[46] |
| 3  | Wrong duty            | The same actor executes different activities, which should require different privileges. This includes “wrong duty sequence” in the sequence activity, “wrong duty decision” in the decision activity, and “combined wrong duty”, a combination of wrong duty sequence and wrong decision sequence [42]. | [7], [30], [31], [42]–[45] |
| 4  | Wrong pattern         | Deviation from the standard sequence prescribed in the SOP.  | [7], [30], [31], [42]–[46] |
| 5  | Wrong decision        | Decision activity execution is a deviation from standard decision execution, as stated in the SOP.   | [7], [30], [31], [42]–[45] |
| 6  | Wrong throughput time | The activity execution time deviates from the standard time, as stated in the SOP. It includes “wrong throughput time min” and “wrong throughput time max” [42].   | [7], [30], [31], [42]–[45] |
| 7  | Parallel event        | Nonparallel events are performed simultaneously.   | [7], [30]                  |
| 8  | Originator behavior   | The actor’s behavior while executing the activity is anomalous.  | [7], [30], [31]            |

Process mining anomaly techniques include control flow analysis, role resource analysis, throughput time analysis, and decision point analysis [39]. The study undertaken by [4], which proposed a process mining method for PBF detection, suggested the concept “1 + 5 + 1”, which includes (1) log preparation; (5) (a) log analysis, (b) performance analysis, (c) social analysis, (d) conformance analysis, (e) process analysis; and (1) refocusing and iteration. A combination of the red flag approach (i.e., metrics approach) and process mining were proposed in [25] to reduce the false positive rate in detecting fraud. The method connects the red flag approach with process mining by using the red flag to present unusual behavior, whereas process mining involves visualizing the business process flow. In [40], a validated method, based on the most accepted lifecycle model for the implementation of the process mining project [41], was proposed for an application in auditing information systems. It used process mining as an expert system engine to address the limitations of other auditing methods involved in fraud detection, including sampling, due to questionable effectiveness as they lack automation and have a narrow scope.

### III. METHODOLOGY

In her remarkable and exceptional work, Gregor [47] explained information systems (IS) theories in terms of five types: analytic theory, explaining theory, prediction theory, explaining and prediction theory, and design and action theory<sup>5</sup> Taxonomy is a taxonomic theory and can be classified as an analysis theory [47]. Analysis theories define or classify specific dimensions or characteristics of individuals, groups, situations, or events by describing the shared features found in discrete observations [47]. These theories answer *what* questions, and they are used as a foundation for developing more advanced theories, as shown in Fig. 5 [47], [48].

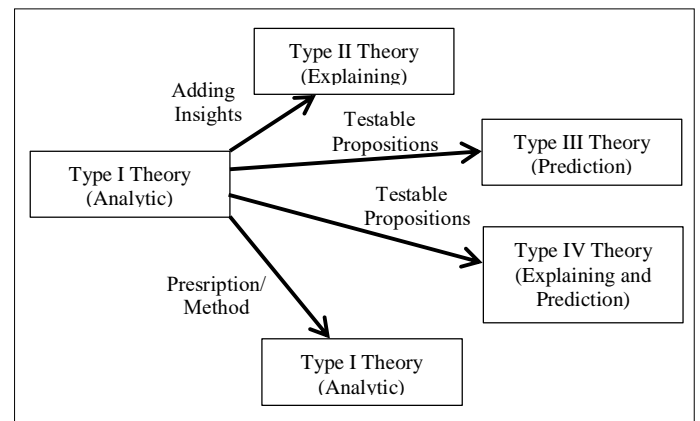


Fig. 5. Evolution of Analytic Theories into other Types. Source: [48].

The DSR methodology can be used to conduct research when the desired goal is an artifact or a recommendation [49]. DSR artifacts are classified into constructs, models, methods, and instantiations [50]. The developed PBF detection metrics are subsumed under the *method* artifact type [51]. This study aims to design an artificial (i.e., human-made) artifact (i.e., PBF detection metric), which fits well within the DSR

<sup>5</sup> For more information, see [47].



paradigm [52]. Furthermore, the pragmatic viewpoint of DSR, which confirms the inability to separate utility from reality [49], is suitable for the nature of the activity of PBF detection.

Following the DSR paradigm, the taxonomic theory [13] was used in this research as the foundation for deriving PBF detection metrics. The taxonomy [13] was developed using DSR’s build/evaluate cycle [52], which led to the definition of the building blocks of PBF detection metrics by implementing the method of Nickerson *et al.* [22]. Since taxonomy can be used as a foundation to produce new knowledge [22], [47], [48], [53], the taxonomy of fraud detection metrics for business processes [13] was used deductively to develop the PBF detection metrics (i.e., the taxonomy’s objects). Adapting [54], the following steps were taken to develop the metrics:

- Define the measured entity in the study, namely, business process.
- Specify the attributes of the defined entity (i.e., business process), which are already developed by the taxonomy (i.e., the taxonomy’s dimensions and characteristics) [13].
- Define the metrics by matching the attributes of the defined entity.

Theoretical validation of the developed metrics can be achieved through the use of a validated taxonomic theory [13]. In addition, in order to evaluate the utility of every developed metric, an illustrative scenario was used [16]. Lastly, an implementation was provided to explain the metrics technical application.

#### IV. PBF DETECTION METRICS

Using the taxonomy of fraud detection metrics for business processes as the underlying theory [13], PBF detection metrics can be derived by matching the characteristics of the taxonomy’s dimensions. Selecting the matched characteristics depends on the application domain, project scope, and the case situation. However, general PBF detection metrics can be developed by matching the selected characteristics from the process perspectives, presentation layers, and fraud data schemes dimensions.<sup>6,7</sup> Table II shows the derived list of PBF detection metrics, including the metric’s ID, name, description, and the illustrative scenario. The generally derived PBF detection metrics covered all the dimensions of PBF detection (i.e., full-dimensional metrics), as stated in the taxonomy of fraud detection metrics for business processes [13].

TABLE II. GENERAL SAMPLES OF PBF DETECTION METRICS

| ID | Metric name         | Description   | Illustrative scenario  |
|----|---------------------|---|--|
| 1  | Wrong activity time | Indicates whether the process activity’s time is incorrect. | The execution time of the approval activity in invoice XYZ is not valid. |

<sup>6</sup> Other metrics can be similarly developed by matching the selected characteristics that should be specified for every project.

<sup>7</sup> The selected characteristic of the fraud domain dimension is *general*. This is because the scope of the developed metrics does not focus on a specific fraud domain.

|    |                             |   |  |
|----|-----------------------------|---|--|
| 2  | Wrong instance time         | Shows whether the process instance’s time is incorrect.                       | The waiting time between activity A and activity B in an invoiced instance exceeds the allowed time.                                     |
| 3  | Wrong stream time           | Indicates whether the process stream’s time is incorrect.                     | The waiting time between the raising of invoice XYZ and its payment as processes in the purchase-to-pay stream exceeds the allowed time. |
| 4  | Discrepant instance time    | Shows whether the process instance’s time causes conflict.                    | The throughput time of an activity is longer than the throughput time of the instance that includes the activity.                        |
| 5  | Discrepant stream time      | Reveals whether the process stream’s time causes conflict.                    | The execution time of invoice XYZ and its payment as processes in the purchase-to-pay stream are identical.                              |
| 6  | Anomalous activity time     | Indicates whether the process activity’s time is abnormal.                    | The execution of the approval activity in invoice XYZ occurred outside of the working hours.   |
| 7  | Anomalous instance time     | Indicates whether the process instance’s time is abnormal.                    | The throughput time of a payment instance is too short.  |
| 8  | Anomalous model time        | Shows whether the process model’s time is abnormal.                           | The execution time of all payment instances for supplier XYZ are all at 8 P.M.   |
| 9  | Anomalous stream time       | Indicates whether the process stream’s time is abnormal.                      | The waiting time between receiving and inspection as processes in the purchase-to-pay stream is very long.                               |
| 10 | Anomalous map time          | Indicates whether the process map’s time at the <i>map layer</i> is abnormal. | The total execution time of all the organization’s processes is too short.   |
| 11 | Wrong activity function     | Indicates whether the process activity’s work is incorrect.                   | The decision was incorrectly made in activity XYZ.   |
| 12 | Wrong instance function     | Shows whether the process instance’s work is incorrect.                       | A payment instance must not be executed because the vendor’s work is not yet finished.   |
| 13 | Wrong stream function       | Reveals whether the process stream’s work is incorrect.                       | A payment process was executed before the receiving process in the purchase-to-pay stream.   |
| 14 | Missing activity function   | Indicates whether the necessary process activity’s work is missing.           | The approval activity in invoice XYZ is missing.   |
| 15 | Missing instance function   | Demonstrates whether the necessary process instance’s work is missing.        | The inspection instance to show that item XYZ was checked is missing.  |
| 16 | Anomalous activity function | Indicates whether the process activity’s work is unusual.                     | The decision made in activity XYZ was unexpected.  |
| 17 | Anomalous instance function | Shows whether the process instance’s work is unusual.                         | Unnecessary activities (i.e., excessive work) are performed in executing an invoice instance.  |
| 18 | Anomalous model function    | Indicates whether the process model’s work is unusual.                        | The number of refund instances of customer XYZ is unusual.   |
| 19 | Anomalous                   | Shows whether the   | Purchase-to-pay processes  |

|    |                             |  |  |
|----|-----------------------------|--|--|
|    | stream function             | process stream's work is unusual.  | for supplier XYZ have always had a non-standard process flow without justifications.   |
| 20 | Anomalous map function      | Indicates whether the process map's work at the <i>map layer</i> is unusual.           | Cancellation of 25% of the organization's processes.   |
| 21 | Discrepant stream function  | Shows if the process stream's work causes conflict.                                    | The payment instances are more than the invoice instances as processes in the order-to-cash stream; however, they should be the same.      |
| 22 | Wrong activity data         | Indicates whether the data produced or consumed by the process activity are incorrect. | The attached document in activity XYZ at invoice A is invalid.   |
| 23 | Missing activity data       | Indicates whether the data produced or consumed by the process activity are missing.   | The signature data in activity XYZ of invoice A is missing.  |
| 24 | Discrepant activity data    | Shows whether the data produced or consumed by the process activity are inconsistent.  | The attached form in the activity XYZ at invoice A has a signature date that follows the activity date.                                    |
| 25 | Discrepant instance data    | Shows whether the data produced or consumed by the process instance are inconsistent.  | In an invoice instance, the input data of activity B does not match the output data of activity A, though they should be equal.            |
| 26 | Discrepant stream data      | Indicates whether the data produced or consumed by a process stream are inconsistent.  | The total amount of orders and the total cash received as processes in the order-to-cash stream should be equal but they differ.           |
| 27 | Anomalous activity data     | Shows whether the data produced or consumed by the process activity are suspicious.    | The activity XYZ has unnecessary recorded data (maybe to complicate the auditing process).   |
| 28 | Anomalous instance data     | Indicates whether the data produced or consumed by a process instance are suspicious.  | The attached document in the activities A and B of a process instance are in different formats. Even the document should not be different. |
| 29 | Anomalous model data        | Shows whether the data produced or consumed by the process model are questionable.     | The inspection instances of all the items from supplier XYZ always have a lengthy inspection report.                                       |
| 30 | Wrong activity resource     | Indicates whether the process activity's resource is incorrect.                        | An employee not authorized to perform activity XYZ in an invoice instance performed it.  |
| 31 | Wrong instance resource     | Shows whether the process instance's resource is incorrect.                            | Issue and review of invoice XYZ performed by the same employee (violates the separation of duties law).                                    |
| 32 | Missing activity resource   | Indicates if the process activity's resource is missing.                               | An anonymous person performed activity XYZ in an invoice instance.   |
| 33 | Anomalous activity resource | Shows whether the process activity's resource is suspicious.                           | Activity XYZ in a payment instance, usually executed by employee X is executed by employee   |

|    |                             |  |   |
|----|-----------------------------|--|---|
|    |                             |  | Y instead.  |
| 34 | Anomalous instance resource | Shows whether the process instance's resource is suspicious.     | The same employee did most of the activities in a receiving instance XYZ.   |
| 35 | Anomalous model resource    | Indicates whether the process model's resource is suspicious.    | The same employee approved all the payments for supplier XYZ.   |
| 36 | Wrong activity location     | Shows whether the process activity's location is incorrect.      | The activity XYZ of a receiving inventory instance was executed outside the approved receiving area.                        |
| 37 | Wrong instance location     | Reveals whether the process instance's location is incorrect.    | Two activities to be executed at the same location for an invoice instance performed at different locations.                |
| 38 | Anomalous activity location | Indicates whether the process activity's location is suspicious. | The execution location of activity XYZ in a payment instance was very distant.  |
| 39 | Anomalous model location    | Shows whether the process model's location is suspicious.        | All the large payments are made only at one location.   |
| 40 | Anomalous stream location   | Indicates whether the process stream's location is suspicious.   | Two processes that are usually executed in the same place in the order-to-cash stream were executed at different locations. |
| 41 | Missing activity location   | Shows whether the process activity's location is missing.        | The execution location of activity XYZ in a payment instance is not specified.  |

## V. IMPLEMENTATION

Based on [25], [40], [41], as well as the taxonomy developed in [13], a method can be proposed for implementing PBF detection metrics. The method uses data and process mining to ensure an effective PBF detection process. Both techniques are used to detect fraud in business processes [45], [55]. Although data mining and process mining share many features, the key difference is that data mining aims to discover previously unknown and interesting patterns in the datasets, while process mining focuses on finding process relationships [28]. Thus, data mining techniques for detecting fraud are usually unsuitable for analyzing the behavior of control flow in a business process [39]. However, process mining can be used to assess the control flow of a business process [56] and to analyze process performance, event sequence, and process roles [57]. Still, process mining focuses on the control flow of transactions [56] and not on process content (e.g., transaction value). Therefore, data mining and process mining are both needed.

Real data [58] on purchase-to-pay process events in a multinational paints and coatings company were used for implementation.<sup>8</sup> The implementation method is illustrated in Fig. 6 and described in the following steps:

<sup>8</sup> To reduce data noise, the data were filtered according to document type, item category, and timeframe to include "standard PO," "three-way matching," "invoice before GR," and 2018 (quarter 2).

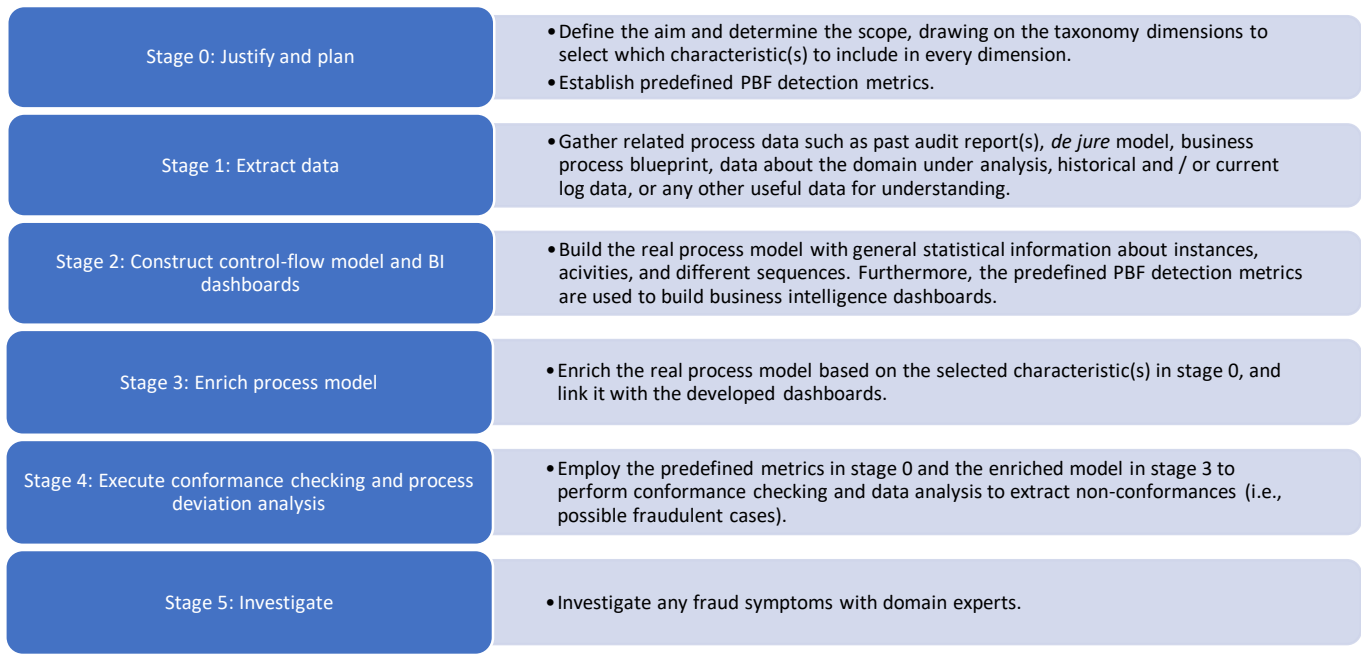


Fig. 6. Implementation Steps. Adapted from [25], [40], [41].

**Stage 0:** At this stage, the scope and aims should be defined after establishing a thorough understanding of the application domain. This includes understanding the business process, identifying the theoretical existence of fraud schemes, cataloging all potential fraud methods and red flags<sup>9</sup>, defining the general multi-dimensional metrics by using the taxonomy, and defining specific multi-dimensional metrics for the selected fraud schemes and methods. Every metric may include a metric formula, data source, metric description, data update frequency, metric unit, threshold or compared value, related fraud scheme, and fraud method or red flag.

In this implementation, the aim was to detect fraud in the purchase-to-pay process by examining execution deviations. The scope was determined based on the following dimensions and characteristics of the taxonomy of fraud detection metrics for business processes [13]:

- **Fraud domain:** In this implementation, the purchase-to-pay business process was selected. Thus, {specific: finance and general} were chosen as the fraud areas for the implementation because general PBF detection metrics are also used.
- **Presentation layer(s):** {process stream, model, instance, and activity} were selected to satisfy the aim. However, the process stream layer was not included in the implementation due to missing data.
- **Process perspective(s):** {time, function, data, and resource} were selected. Location perspective data are not available. However, depending on the case situation and data availability, it may be useful to include all process perspectives.

- **Fraud data scheme(s):** To specify critical data schemes that can effectively detect fraud in this implementation, {anomalous, discrepant, missing, and wrong} were selected. The selection of the fraud data scheme characteristics was based on the case situation and the quality of existing data. However, if possible, it is always useful to include all fraud data schemes.

The selected dimensions, along with their characteristics, ought to assist in developing the predefined metrics. Fraud examiners can also add more useful metrics based on their experience. In this implementation, the generic and specific metrics defined in Appendix A are used based on the case situation and the existing data.<sup>10</sup> The specific multi-dimensional metrics for the fraud schemes and fraud methods are defined based on the common fraud schemes appearing in the fraud tree [10].<sup>11</sup> The fraud tree was selected for the following reasons: (1) it represents a comprehensive classification of the most common occupational financial fraud schemes; and (2) it is developed by a standards organization (ACFE).

**Stage 1:** At this stage, all the useful process data for detecting PBF should be collected. Examples of data that should be collected are the past audit reports, process events log, and process model, as depicted in Fig. 7 [59]. This model is referred to as the *de jure* model, which represents the desired, ideal, or required process.

<sup>9</sup> Red flags are signs of potentially fraudulent behavior [62].

<sup>10</sup> Sound knowledge of business rules is valuable in defining effective metrics.

<sup>11</sup> For more information about the fraud tree, see [10].

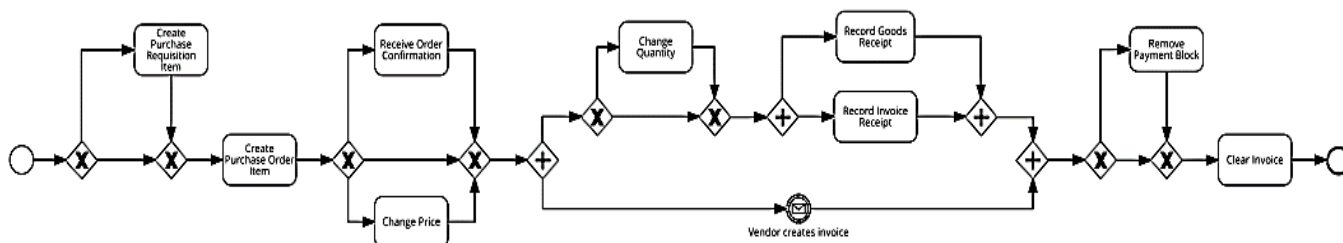


Fig. 7. The De Jure Model for the Purchase-to-Pay Process. Source: [59].

**Stage 2:** Using the process mining discovery technique, the *de facto* model with general statistical information was constructed as shown in Fig. 8. The *de facto* model describes reality with potential violations [60]. It was implemented using the Celonis process mining software.<sup>12</sup> It is possible for the auditor to analyze differences between the *de jure* and *de facto* models in order to detect fraud [33].

Moreover, the predefined metrics were represented on dashboards, as shown in Appendix B. In this case, the Celonis process mining software was also used.

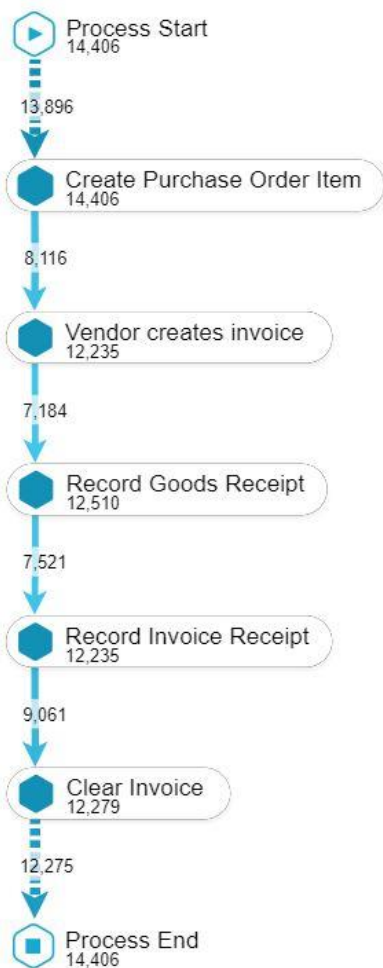


Fig. 8. The De Facto Model, with Case Frequency, for the Purchase-to-Pay Process.

<sup>12</sup> See [www.celonis.com](http://www.celonis.com)

**Stage 3:** This stage involves enriching the *de facto* model based on the process perspective characteristics selected in stage 0, as shown in Appendix B. In addition, the *de facto* model is linked to the dashboards that are used to represent the predefined metrics using Celonis process mining software.

**Stage 4:** Conformance checking and process deviation analysis should be applied to combine misuse-based techniques and anomaly-based techniques. The misuse-based technique is implemented by creating dashboards that leverage business intelligence (BI) techniques for the predefined metrics, while the anomaly-based technique is implemented using the conformance checking technique.

Conformance checking is used to compare the business process with its SOP [30]. This is relevant to auditing [40] because it can detect, locate, and explain the deviation from the behavior expected in business process [56]. It helps detect the occurrence of event skipping and enables analysis of the flow of the business process [30]. Using conformance checking to classify standard and non-standard business process variants can assist in detecting potential risks [33].

A process variant is a single path (i.e., routing) that is followed by at least one business process instance [33]. All business process instances that follow the same path are grouped into the same variant [33]. Thus, it is possible to examine process variants to find out all business process instances that are in non-standard paths [33]. In turn, each process variant can be prioritized using the metrics, thereby reducing the rate of false positives in detecting fraud [25]. Reducing false positives saves time and cost [61].

**Stage 5:** In this stage, the fraud symptoms should be investigated with domain experts to confirm the presence or absence of fraud [25].

## VI. RESULTS AND DISCUSSION

Using the enriched model in stage 3, the conformance checking procedure was applied to extract non-conformances that form potentially fraudulent cases. The findings of the conformance checking revealed that there were 431 process flow variants (control-flow perspective). The number of variants is usually large because the process should be flexible to meet all business needs. Thus, the use of metrics as filters is essential to save time and effort, and to discover new signs of fraud.

Using the enriched model assists in fraud detection without the influence of the fraud examiner [40]. Moreover, using the predefined metrics in stage 0 ensures the accuracy and comprehensiveness of fraud detection. This is because the

predefined metrics can be used to detect fraud in the content perspective (not just the control-flow perspective) of the business process.

The combination of visual analytics and process mining can help to identify data integrity issues such as missing, non-conforming, or anomalous activities undertaken by a privileged user, or those with suspiciously short execution times [56]. Furthermore, applying the metrics using process mining reduces the number of false positives in fraud detection [25]. Thus, conformance checking and process deviation analysis are used to detect PBF [62].

In Appendices A and B, the implementation screens and results are provided. Each implementation screen serves as a link between the process flow view and the data view to present a complete view. The results show that 13 metrics produced results that should be investigated.

This implementation shows that the developed metrics can be used in the following ways: (1) directly, thereby conserving time and effort; (2) as a template, thereby facilitating the definition of other metrics and ensuring consistency among PBF detection stakeholders; and (3) to determine the implementation scope. Additionally, the developed metrics are process-oriented metrics that can measure throughput processing, as opposed to just measuring process input–output relations. This helps to detect and predict fraud, with its root cause, in its initial stages.

## VII. CONCLUSIONS

This study sought to develop a comprehensive list of fraud detection metrics for business processes. A taxonomy of fraud detection metrics for business processes was used as a “building” theory to generate all possible metrics for detecting fraud in business processes. Compared to the 8 existing PBF detection metrics, 41 comprehensive metrics were developed, classified, and demonstrated. These metrics cover each of the PBF detection dimensions that are not entirely (e.g., presentation layer) or partially incorporated into existing PBF detection metrics. Additionally, their applications were demonstrated by using illustrative scenarios. Finally, their technical implementation was explained by providing an implementation that offers an accurate and comprehensive view for PBF examiners.

The study’s contributions to the literature are twofold. First, the study offers improved DSR artifacts (i.e., the developed metrics and their implementation method), which can enhance the ability to detect PBF. Second, the study enriched the construction of the taxonomic theory [13] (i.e., by leveraging the taxonomy for a purpose beyond analysis). This is a step toward developing advanced theories such as design and action theory. The study also is relevant due to its practical contribution in improving PBF detection in the workplace. PBF stakeholders can improve their practices by using the developed PBF detection metrics to bolster their effectiveness.

The limited availability of data on fraud is one of the limitations of this study. This relates to the fact fraud is a sensitive topic in public discussion, and so it is not an issue spoken about openly. However, the data issued by standard-

setting organizations such as the Committee of Sponsoring Organizations (COSO)<sup>13</sup> and the ACFE can mitigate this limitation to a certain degree. Nevertheless, the data from these organizations are mainly from the finance domain. In addition to these limitations, reviewing the metrics results with domain experts (i.e., the investigation step) is needed to confirm fraud cases. However, the scope here is specified to detect possible PBF.

Extending and validating the metrics in other domains (e.g., the telecommunications sector) is suggested as a possible direction for future research. In addition, case studies within organizations, which prioritize the use of the metrics in their specific context, are suggested. Linking each metric to a full list of possible deviation patterns is another worthwhile research opportunity. For example, the *wrong instance function* is a suitable metric that can be linked with deviation patterns such as looping, swapping, and inserting activities in the process model.

## ACKNOWLEDGMENT

The authors would like to thank the deanship of scientific research at King Saud University for funding and supporting this research through the DSR Graduate Students’ Research Support initiative.

## REFERENCES

- [1] Cotton, S. Johnigan, and L. Givarz, *Fraud risk management guide*. COSO, 2016.
- [2] “Report to the nations: Global study on occupational fraud and abuse,” 2020.
- [3] D. Al-Jumeily, A. Hussain, A. MacDermott, G. Seeckts, and J. Lunn, “Methods and techniques to support the development of fraud detection system,” in *IWSSIP*, 2015, pp. 224–227.
- [4] J. J. Stoop, “Process mining and fraud detection - A case study on the theoretical and practical value of using process mining for the detection of fraudulent behavior in the procurement process,” Twente University, 2012.
- [5] F. Sinaga and R. Sarno, “Business process anomaly detection using multi-level class association rule learning,” *IPTEK J. Proc. Ser.*, vol. 2, no. 1, 2016.
- [6] D. S. Kerr, “The importance of the CobiT framework IT processes for effective internal control over financial reporting in organizations: An internationale surveys,” *Inf. Manag.*, vol. 50, no. 7, pp. 590–597, 2013.
- [7] S. Huda, R. Sarno, and T. Ahmad, “Increasing accuracy of process-based fraud detection using a behavior model,” *Int. J. Softw. Eng. Its Appl.*, vol. 10, no. 5, pp. 175–188, May 2016.
- [8] M. Dumas, M. La Rosa, J. Mendling, and H. A. Reijers, *Fundamentals of business process management*. New York, NY, USA: Springer, 2013.
- [9] J. Jeston, *Business process management practical guidelines to successful implementations*. Taylor and Francis, 2017.
- [10] “Fraud tree,” ACFE. [Online]. Available: <http://www.acfe.com/fraud-tree.aspx>. [Accessed: 10-May-2020].
- [11] B. Omair and A. Alturki, “A systematic literature review of fraud detection metrics in business processes,” *IEEE Access*, vol. 8, no. 1, pp. 26893–26903, Feb. 2020.
- [12] S. Bhattacharyya, S. Jha, K. Tharakunnel, and J. C. Westland, “Data mining for credit card fraud: a comparative study,” *Decis. Support Syst. Sci.*, vol. 50, no. 3, pp. 602–613, Feb. 2011.
- [13] B. Omair and A. Alturki, “Taxonomy of fraud detection metrics for business processes,” *IEEE Access*, vol. 8, pp. 71364–71377, 2020.

<sup>13</sup> The COSO of the Treadway Commission is a joint initiative to combat corporate fraud. <https://www.coso.org/>

- [14] K. Peffers, M. Rothenberger, T. Tuunanen, and R. Vaezi, "Design science research evaluation," in International Conference on Design Science Research in Information Systems, 2012, pp. 398–410.
- [15] B. Baesens, V. Van Vlasselaer, and W. Verbeke, *Fraud analytics using descriptive, predictive, and social network techniques: A Guide to data science for fraud detection*. Hoboken, NJ, USA: John Wiley and Sons, 2015.
- [16] A. Abdallah, M. A. Maarof, and A. Zainal, "Fraud detection system: A survey," *J. Netw. Comput. Appl.*, vol. 68, pp. 90–113, Jun. 2016.
- [17] V. Jyothsna, "A review of anomaly based intrusion detection systems," *Int. J. Comput. Appl.*, vol. 28, no. 7, pp. 975–8887, Sep. 2011.
- [18] J. Akhilomen, "Data mining application for cyber credit-card fraud detection system," in Industrial Conference on Data Mining, 2013, pp. 218–228.
- [19] K. Mule and M. Kulkarni, "Credit card fraud detection using hidden Markov model (HMM)," *Int. J. Innov. Technol. Adapt. Manag.*, vol. 1, no. 6, Aug. 2014.
- [20] R. Nisbet, G. Miner, and K. Yale, "Fraud detection," in *Handbook of Statistical Analysis and Data Mining Applications*, Amsterdam, Netherlands: Elsevier, 2018, pp. 289–302.
- [21] T. W. Singleton and A. J. Singleton, *Fraud risk assessment*, vol. 160. John Wiley and Sons, 2011.
- [22] R. C. Nickerson, U. Varshney, and J. Muntermann, "A method for taxonomy development and its application in information systems," *Eur. J. Inf. Syst.*, vol. 22, no. 3, pp. 336–359, May 2013.
- [23] J. West and M. Bhattacharya, "An investigation on experimental issues in financial fraud mining," in ICIEA, 2016, vol. 80, pp. 1796–1801.
- [24] M. Werner, "Process model representation layers for financial audits," *Proc. Annu. Hawaii Int. Conf. Syst. Sci.*, vol. 2016-March, pp. 5338–5347, 2016.
- [25] G. Baader and H. Kremer, "Reducing false positives in fraud detection: Combining the red flag approach with process mining," *Int. J. Account. Inf. Syst.*, vol. 31, no. March, pp. 1–16, Dec. 2018.
- [26] G. Vossen, "The Process Mining Manifesto—An interview with Wil van der Aalst," *Inf. Syst.*, vol. 37, no. 3, pp. 288–290, May 2012.
- [27] H. A. Reijers, I. Vanderfeesten, and W. M. P. van der Aalst, "The effectiveness of workflow management systems: A longitudinal study," *Int. J. Inf. Manage.*, vol. 36, no. 1, pp. 126–141, Feb. 2016.
- [28] S.-M.-R. Beheshti et al., *Process Analytics*. Cham: Springer International Publishing, 2016.
- [29] W.-S. Yang and S.-Y. Hwang, "A process-mining framework for the detection of healthcare fraud and abuse," *Expert Syst. Appl.*, vol. 31, no. 1, pp. 56–68, Jul. 2006.
- [30] S. Huda, R. Sarno, and T. Ahmad, "Fuzzy MADM approach for rating of process-based fraud," *J. ICT Res. Appl.*, vol. 9, no. 2, pp. 111–128, Nov. 2015.
- [31] R. Sarno, R. D. Dewandono, T. Ahmad, M. F. Naufal, and F. Sinaga, "Hybrid association rule learning and process mining for fraud detection," *IAENG Int. J. Comput. Sci.*, vol. 42, no. 2, pp. 59–72, Apr. 2015.
- [32] M. Jans, M. G. Alles, and M. A. Vasarhelyi, "A Field study on the use of process mining of event logs as an analytical procedure in auditing," *Account. Rev.*, vol. 89, no. 5, pp. 1751–1773, Sep. 2014.
- [33] T. Chiu, "Exploring New Audit Evidence: the Application of Process Mining in Auditing," Rutgers, The State University of New Jersey, 2018.
- [34] M. zur Muehlen and R. Shapiro, "Business Process Analytics," in *Handbook on Business Process Management 2*, Second., Berlin, Heidelberg: Springer Berlin Heidelberg, 2015, pp. 243–263.
- [35] T. Chiu, Y. Wang, and M. A. Vasarhelyi, "A framework of applying process mining for fraud scheme detection," *SSRN Electron. J.*, Jun. 2017.
- [36] M. Leyer, D. Heckl, and J. Moormann, "Process Performance Measurement," in *Handbook on Business Process Management 2*, Berlin, Heidelberg: Springer Berlin Heidelberg, 2015, pp. 227–241.
- [37] M. Jans, M. Alles, and M. Vasarhelyi, "The case for process mining in auditing: Sources of value added and areas of application," *Int. J. Account. Inf. Syst.*, vol. 14, no. 1, pp. 1–20, 2013.
- [38] C. dos S. Garcia et al., "Process mining techniques and applications – A systematic mapping study," *Expert Syst. Appl.*, vol. 133, pp. 260–295, Nov. 2019.
- [39] R. Sarno, F. Sinaga, and K. R. Sungkono, "Anomaly detection in business processes using process mining and fuzzy association rule learning," *J. Big Data*, vol. 7, no. 1, p. 5, Dec. 2020.
- [40] P. Zerbino, D. Aloini, R. Dulmin, and V. Mininno, "Process-mining-enabled audit of information systems: Methodology and an application," *Expert Syst. Appl.*, vol. 110, pp. 80–92, Nov. 2018.
- [41] W. van der Aalst et al., "Process mining manifesto," in *Lecture Notes in Business Information Processing*, vol. 99 LNBIP, no. PART 1, Berlin, Heidelberg: Springer, Berlin, Heidelberg, 2012, pp. 169–194.
- [42] R. Sarno and F. P. Sinaga, "Business process anomaly detection using ontology-based process modelling and Multi-Level Class Association Rule Learning," in IC3INA, 2015, pp. 12–17.
- [43] E. S. Pane, A. D. Wibawa, and M. H. Purnomo, "Event log-based fraud rating using interval type-2 fuzzy sets in fuzzy AHP," in *IEEE Region 10 Conference (TENCON)*, 2016, pp. 1965–1968.
- [44] S. Huda, T. Ahmad, R. Sarno, and H. A. Santoso, "Identification of process-based fraud patterns in credit application," in *ICOICT*, 2014, pp. 84–89.
- [45] D. Rahmawati, R. Sarno, C. Fatchah, and D. Sunaryono, "Fraud detection on event log of bank financial credit business process using Hidden Markov Model algorithm," in *3rd ICSITech*, 2017, pp. 35–40.
- [46] H. A. Hartanto, R. Sarno, and N. F. Ariyani, "Linked warning criterion on ontology-based key performance indicators," in *ISemantic*, 2016, pp. 211–216.
- [47] S. Gregor, "The nature of theory in information systems," *MIS Q.*, vol. 30, no. 3, pp. 611–642, Sep. 2006.
- [48] J. Muntermann, R. Nickerson, and U. Varshney, "Towards the development of a taxonomic theory," in *21st AMCIS*, 2015, no. Gregor 2006, pp. 1–15.
- [49] A. Dresch, D. P. Lacerda, and J. A. V. Antunes Jr, *Design science research*. Cham: Springer International Publishing, 2015.
- [50] S. T. March, "Design and natural science research on information technology," *Decis. Support Syst.*, vol. 15, no. 4, pp. 251–266, Dec. 2003.
- [51] O. M. Sangupamba, N. Prat, and I. Comyn-Wattiau, "Business intelligence and big data in the cloud: opportunities for design-science researchers," in *International Conference on Conceptual Modeling*, 2014, pp. 75–84.
- [52] A. Hevner, S. March, and J. Park, "Design science in information systems research," *MIS Q. Manag. Inf. Syst.*, vol. 28, no. 1, pp. 75–105, Mar. 2004.
- [53] B. Omair and A. Alturki, "An improved method for taxonomy development in information systems," *Int. J. Adv. Comput. Sci. Appl.*, vol. 11, no. 4, pp. 535–540, 2020.
- [54] N. E. Fenton and S. L. Pfleeger, *Software metrics: A rigorous and practical approach*, vol. 2. 1997.
- [55] A. Boenner, "Bayer: Process mining supports digital transformation in internal audit," in *Process Mining in Action*, Cham: Springer International Publishing, 2020, pp. 159–168.
- [56] G. Moggia and Z. Varga, "Connecting data and processes in audit — some considerations about the use of process mining," *European Court of Auditors*, 2020.
- [57] M. Jans, J. M. van der Werf, N. Lybaert, and K. Vanhoof, "A business process mining application for internal transaction fraud mitigation," *Expert Syst. Appl.*, vol. 38, no. 10, pp. 13351–13359, Sep. 2011.
- [58] B. F. (Boudewijn) Van Dongen, "BPI Challenge 2019." 4TU.Centre for Research Data, 2019.
- [59] K. Diba, S. Remy, and L. Pufahl, "Compliance and Performance Analysis of Procurement Processes Using Process Mining," in *International Conference on Process Mining*, 2019.

- [60] W. M. P. van der Aalst, K. M. van Hee, J. M. van der Werf, and M. Verdonk, "Auditing 2.0: Using process mining to support tomorrow's auditor," Computer (Long Beach, Calif.), vol. 43, no. 3, pp. 90–93, Mar. 2010.
- [61] J. Luell, "Employee fraud detection under real world conditions," University of Zurich, 2010.
- [62] G. Baader and H. Krcmar, "Reducing false positives in fraud detection: Combining the red flag approach with process mining," Int. J. Account. Inf. Syst., vol. 31, pp. 1–16, Dec. 2018.
- [63] L. Sánchez González, F. García Rubio, F. Ruiz González, and M. Piattini Velthuis, "Measurement in business processes: A systematic review," Bus. Process Manag. J., vol. 16, no. 1, pp. 114–134, Feb. 2010.

APPENDIX A

|                            |   |                           |          |
|----------------------------|---|---------------------------|----------|
| <b>Metric name</b>         | WAT_Generic   | <b>Threshold</b>          | 0        |
| <b>Fraud data scheme</b>   | Wrong   | <b>Presentation layer</b> | Activity |
| <b>Process perspective</b> | Time  | <b>Fraud domain</b>       | Generic  |
| <b>Metric description</b>  | Counts the activities with execution time not in 2018 (Q2).   |                           |          |
| <b>Metric formula</b>      | <pre>SUM(CASE WHEN YEAR("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."COMPLETE_TIMESTAMP") &lt;&gt; 2018 THEN 1.0         WHEN QUARTER("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."COMPLETE_TIMESTAMP") &lt;&gt; 2 THEN 1.0         ELSE 0.0 END)</pre>   |                           |          |
| <b>Result</b>              | 0   |                           |          |
| <b>Metric name</b>         | AAT_Generic   | <b>Threshold</b>          | 0        |
| <b>Fraud data scheme</b>   | Anomalous   | <b>Presentation layer</b> | Activity |
| <b>Process perspective</b> | Time  | <b>Fraud domain</b>       | Generic  |
| <b>Metric description</b>  | Counts the activities with execution time outside of normal working hours (between 8 PM and 6 AM).  |                           |          |
| <b>Metric formula</b>      | <pre>SUM(CASE WHEN HOURS("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."COMPLETE_TIMESTAMP") &gt;= 20 THEN 1         WHEN HOURS("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."COMPLETE_TIMESTAMP") &lt;= 6 THEN 1 ELSE 0 END)</pre>  |                           |          |
| <b>Result</b>              | 15254   |                           |          |
| <b>Metric name</b>         | AIT_Generic   | <b>Threshold</b>          | 0        |
| <b>Fraud data scheme</b>   | Anomalous   | <b>Presentation layer</b> | Instance |
| <b>Process perspective</b> | Time  | <b>Fraud domain</b>       | Generic  |
| <b>Metric description</b>  | Monitors instances throughput time that is less than 2 days.  |                           |          |
| <b>Metric formula</b>      | <pre>CASE WHEN AVG(CALC_THROUGHPUT(CASE_START TO CASE_END,         REMAP_TIMESTAMPS("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."COMPLETE_TIMESTAMP", DAYS))) &lt; 2 THEN 1 ELSE 0 END</pre>  |                           |          |
| <b>Result</b>              | 920   |                           |          |
| <b>Metric name</b>         | AMT_Generic   | <b>Threshold</b>          | 0        |
| <b>Fraud data scheme</b>   | Anomalous   | <b>Presentation layer</b> | Model    |
| <b>Process perspective</b> | Time  | <b>Fraud domain</b>       | Generic  |
| <b>Metric description</b>  | The instance throughput time is less than or greater than 43 days (average instance throughput time) by 50%   |                           |          |
| <b>Metric formula</b>      | <pre>CASE WHEN AVG(CALC_THROUGHPUT(CASE_START TO CASE_END, REMAP_TIMESTAMPS("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."COMPLETE_TIMESTAMP", DAYS))) &lt;= 43 * 0.5 THEN 1         WHEN AVG(CALC_THROUGHPUT(CASE_START TO CASE_END, REMAP_TIMESTAMPS("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."COMPLETE_TIMESTAMP", DAYS))) &gt;= 43 * 1.5 THEN 1         ELSE 0 END</pre>  |                           |          |
| <b>Result</b>              | 5187  |                           |          |
| <b>Metric name</b>         | WIF_Generic   | <b>Threshold</b>          | 0        |
| <b>Fraud data scheme</b>   | Wrong   | <b>Presentation layer</b> | Instance |
| <b>Process perspective</b> | Function  | <b>Fraud domain</b>       | Generic  |
| <b>Metric description</b>  | Monitors the wrong work sequence (i.e., "Create Purchase Order Item" activity occurred after "Receive Order Confirmation")  |                           |          |
| <b>Metric formula</b>      | <pre>AVG(CASE WHEN PU_COUNT("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES",         "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY",         "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY" = 'Create Purchase Order Item') = 0 THEN 0         WHEN PU_COUNT("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES",         "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY",         "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY" = 'Receive Order Confirmation') = 0 THEN 0         WHEN DATEDIFF(mi, PU_FIRST("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES",         "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."COMPLETE_TIMESTAMP",         "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY" = 'Create Purchase Order Item'),         PU_FIRST("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES",         "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."COMPLETE_TIMESTAMP",         "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY" = 'Receive Order Confirmation')) &lt; 0 THEN 1 ELSE 0 END)</pre> |                           |          |
| <b>Result</b>              | 0   |                           |          |

|                            |  |                           |          |
|----------------------------|--|---------------------------|----------|
| <b>Metric name</b>         | MIF_Generic  | <b>Threshold</b>          | 0        |
| <b>Fraud data scheme</b>   | Missing  | <b>Presentation layer</b> | Instance |
| <b>Process perspective</b> | Function   | <b>Fraud domain</b>       | Generic  |
| <b>Metric description</b>  | Finds an instance where the “Purchasing Document” (PO number) is null, which is because every event should be connected to a PO number in the events log |                           |          |
| <b>Metric formula</b>      | <code>KPI("Filtered count", ISNULL("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."(CASE) PURCHASING DOCUMENT") = 1 )</code>                          |                           |          |
| <b>Result</b>              | 0  |                           |          |

|                            |   |                           |         |
|----------------------------|---|---------------------------|---------|
| <b>Metric name</b>         | AMF_Generic   | <b>Threshold</b>          | 0       |
| <b>Fraud data scheme</b>   | Anomalous   | <b>Presentation layer</b> | Model   |
| <b>Process perspective</b> | Function  | <b>Fraud domain</b>       | Generic |
| <b>Metric description</b>  | Finds the less frequent activity “Record Subsequent Invoice” (which occurred only once in the events log) |                           |         |
| <b>Metric formula</b>      | <code>AVG (MATCH_ACTIVITIES(NODE[ 'Record Subsequent Invoice' ] ))</code>                                 |                           |         |
| <b>Result</b>              | 1   |                           |         |

|                            |  |                           |          |
|----------------------------|--|---------------------------|----------|
| <b>Metric name</b>         | AIF_Generic  | <b>Threshold</b>          | 0        |
| <b>Fraud data scheme</b>   | Anomalous  | <b>Presentation layer</b> | Instance |
| <b>Process perspective</b> | Function   | <b>Fraud domain</b>       | Generic  |
| <b>Metric description</b>  | Counts instances that do not end with “Clear Invoice” activity   |                           |          |
| <b>Metric formula</b>      | <code>AVG (CASE WHEN PU_LAST("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY") &lt;&gt; 'Clear Invoice' THEN 1 ELSE 0 END)</code> |                           |          |
| <b>Result</b>              | 2132   |                           |          |

|                            |  |                           |          |
|----------------------------|--|---------------------------|----------|
| <b>Metric name</b>         | MAD_Generic  | <b>Threshold</b>          | 0        |
| <b>Fraud data scheme</b>   | Missing  | <b>Presentation layer</b> | Activity |
| <b>Process perspective</b> | Data   | <b>Fraud domain</b>       | Generic  |
| <b>Metric description</b>  | Counts activities with missing price (null).   |                           |          |
| <b>Metric formula</b>      | <code>SUM(CASE WHEN ISNULL("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."CUMULATIVE NET WORTH (EUR)") = 1 THEN 1 ELSE 0 END)</code> |                           |          |
| <b>Result</b>              | 0  |                           |          |

|                            |   |                           |          |
|----------------------------|---|---------------------------|----------|
| <b>Metric name</b>         | DID_Generic   | <b>Threshold</b>          | 0        |
| <b>Fraud data scheme</b>   | Discrepant  | <b>Presentation layer</b> | Instance |
| <b>Process perspective</b> | Data  | <b>Fraud domain</b>       | Generic  |
| <b>Metric description</b>  | Counts instances where the number of “Create Purchase Order Item” is not equal to that of “Create Purchase Requisition Item”  |                           |          |
| <b>Metric formula</b>      | <code>AVG(CASE WHEN PU_COUNT("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY" = 'Create Purchase Requisition Item') = 0 THEN 0 WHEN PU_SUM("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."CUMULATIVE NET WORTH (EUR)", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY" = 'Create Purchase Requisition Item') - PU_SUM("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."CUMULATIVE NET WORTH (EUR)", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY" = 'Create Purchase Order Item') &lt;&gt; 0 THEN 1 ELSE 0 END)</code> |                           |          |
| <b>Result</b>              | 0   |                           |          |

|                            |   |                           |          |
|----------------------------|---|---------------------------|----------|
| <b>Metric name</b>         | DIF_Generic   | <b>Threshold</b>          | 0        |
| <b>Fraud data scheme</b>   | Discrepant  | <b>Presentation layer</b> | Instance |
| <b>Process perspective</b> | Function  | <b>Fraud domain</b>       | Generic  |
| <b>Metric description</b>  | The metric checks if the execution times of (first) “Record Invoice Receipt” and (first) “Clear Invoice” are the same, and it also checks whether “Create Purchase Order Item” activity occurred at the same time as “Receive Order Confirmation” |                           |          |



|                       |   |
|-----------------------|---|
| <b>Metric formula</b> | <pre>AVG(CASE WHEN PU_COUNT("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY" = 'Record Invoice Receipt') = 0 THEN 0 WHEN PU_FIRST("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."COMPLETE_TIMESTAMP", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY" = 'Record Invoice Receipt') = PU_FIRST("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."COMPLETE_TIMESTAMP", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY" = 'Clear Invoice') THEN 1 ELSE 0 END) + AVG(CASE WHEN PU_COUNT("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY" = 'Create Purchase Order Item') = 0 THEN 0 WHEN PU_COUNT("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY" = 'Receive Order Confirmation') = 0 THEN 0 WHEN DATEDIFF(mi, PU_FIRST("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."COMPLETE_TIMESTAMP", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY" = 'Create Purchase Order Item'), PU_FIRST("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."COMPLETE_TIMESTAMP", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY" = 'Receive Order Confirmation')) = 0 THEN 1 ELSE 0 END)</pre> |
| <b>Result</b>         | 7   |

|                            |  |                           |          |
|----------------------------|--|---------------------------|----------|
| <b>Metric name</b>         | AAD_Generic  | <b>Threshold</b>          | 0        |
| <b>Fraud data scheme</b>   | Anomalous  | <b>Presentation layer</b> | Activity |
| <b>Process perspective</b> | Data   | <b>Fraud domain</b>       | Generic  |
| <b>Metric description</b>  | Monitors activities where the price of the purchased item is equal to 1 euro   |                           |          |
| <b>Metric formula</b>      | <pre>SUM(CASE WHEN "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."CUMULATIVE NET WORTH (EUR)" = 1 THEN 1 ELSE 0 END)</pre> |                           |          |
| <b>Result</b>              | 13   |                           |          |

|                            |  |                           |         |
|----------------------------|--|---------------------------|---------|
| <b>Metric name</b>         | AMD_Generic  | <b>Threshold</b>          | 0       |
| <b>Fraud data scheme</b>   | Anomalous  | <b>Presentation layer</b> | Model   |
| <b>Process perspective</b> | Data   | <b>Fraud domain</b>       | Generic |
| <b>Metric description</b>  | Counts instances that have a price less than or greater than 383 euros (average price) by 50%  |                           |         |
| <b>Metric formula</b>      | <pre>CASE WHEN AVG("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."CUMULATIVE NET WORTH (EUR)") &lt;= 383 * 0.5 THEN 1 WHEN AVG("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."CUMULATIVE NET WORTH (EUR)") &gt;= 383 * 1.5 THEN 1 ELSE 0 END</pre> |                           |         |
| <b>Result</b>              | 10783  |                           |         |

|                            |   |                           |          |
|----------------------------|---|---------------------------|----------|
| <b>Metric name</b>         | WAR_Generic   | <b>Threshold</b>          | 0        |
| <b>Fraud data scheme</b>   | Wrong   | <b>Presentation layer</b> | Activity |
| <b>Process perspective</b> | Resource  | <b>Fraud domain</b>       | Generic  |
| <b>Metric description</b>  | Monitors the resource of "Vendor Create Invoice", which should be a vendor value (i.e., NONE in the events log)   |                           |          |
| <b>Metric formula</b>      | <pre>SUM(CASE WHEN "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY" &lt;&gt; 'Vendor creates invoice' THEN 0 WHEN "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."RESOURCE" &lt;&gt; 'NONE' THEN 1 ELSE 0 END)</pre> |                           |          |
| <b>Result</b>              | 0   |                           |          |

|                            |   |                           |          |
|----------------------------|---|---------------------------|----------|
| <b>Metric name</b>         | WIR_Generic   | <b>Threshold</b>          | 0        |
| <b>Fraud data scheme</b>   | Wrong   | <b>Presentation layer</b> | Instance |
| <b>Process perspective</b> | Resource  | <b>Fraud domain</b>       | Generic  |
| <b>Metric description</b>  | Checks the violation of the segregation of duties rule, where the resource of "Record Invoice Receipt" should not be the same resource as "Clear Invoice"   |                           |          |
| <b>Metric formula</b>      | <pre>AVG(CASE WHEN PU_COUNT("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY" = 'Record Invoice Receipt') = 0 THEN 0 WHEN PU_FIRST("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."RESOURCE", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY" = 'Record Invoice Receipt') = PU_FIRST("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."RESOURCE", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY" = 'Clear Invoice') THEN 1 ELSE 0 END)</pre> |                           |          |
| <b>Result</b>              | 13  |                           |          |

|                          |             |                           |          |
|--------------------------|-------------|---------------------------|----------|
| <b>Metric name</b>       | MAR_Generic | <b>Threshold</b>          | 0        |
| <b>Fraud data scheme</b> | Missing     | <b>Presentation layer</b> | Activity |

|                            |  |                     |         |
|----------------------------|--|---------------------|---------|
| <b>Process perspective</b> | Resource   | <b>Fraud domain</b> | Generic |
| <b>Metric description</b>  | Counts activities with null resource   |                     |         |
| <b>Metric formula</b>      | <code>SUM(CASE WHEN ISNULL("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."RESOURCE") = 1 THEN 1 ELSE 0 END)</code> |                     |         |
| <b>Result</b>              | 0  |                     |         |

|                            |  |                           |          |
|----------------------------|--|---------------------------|----------|
| <b>Metric name</b>         | AIR_Generic  | <b>Threshold</b>          | 0        |
| <b>Fraud data scheme</b>   | Anomalous  | <b>Presentation layer</b> | Instance |
| <b>Process perspective</b> | Resource   | <b>Fraud domain</b>       | Generic  |
| <b>Metric description</b>  | Checks whether a resource undertook more than one activity in a complete instance  |                           |          |
| <b>Metric formula</b>      | <code>AVG(CASE WHEN PU_COUNT("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY" = 'Clear Invoice') = 0 THEN 0 WHEN PU_COUNT_DISTINCT("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."RESOURCE") &lt;&gt; PU_COUNT_DISTINCT("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY") THEN 1 ELSE 0 END)</code> |                           |          |
| <b>Result</b>              | 963  |                           |          |

|                            |  |                           |         |
|----------------------------|--|---------------------------|---------|
| <b>Metric name</b>         | AMR_Generic  | <b>Threshold</b>          | 0       |
| <b>Fraud data scheme</b>   | Anomalous  | <b>Presentation layer</b> | Model   |
| <b>Process perspective</b> | Resource   | <b>Fraud domain</b>       | Generic |
| <b>Metric description</b>  | Monitors to determine whether employee frequency is suspicious (e.g., appears only once)                           |                           |         |
| <b>Metric formula</b>      | <code>COUNT_TABLE("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES")</code><br>(based on resource dimension) |                           |         |
| <b>Result</b>              | 0  |                           |         |

|                            |   |                           |   |
|----------------------------|---|---------------------------|---|
| <b>Metric name</b>         | MIF_BillAndHold   | <b>Threshold</b>          | 0   |
| <b>Fraud data scheme</b>   | Missing   | <b>Presentation layer</b> | Instance  |
| <b>Process perspective</b> | Function  | <b>Fraud domain</b>       | Finance,<br>fictitious expenses,<br>bill-and-hold |
| <b>Metric description</b>  | By using "Missing Instance Function", it is possible to define this specific metric, which checks whether the "Clear Invoice" activity is missing, while "Record Invoice Receipt" exists  |                           |   |
| <b>Metric formula</b>      | <code>AVG (CASE WHEN PU_COUNT("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY" = 'Record Invoice Receipt') = 0 THEN 0 WHEN PU_COUNT("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY" = 'Clear Invoice') = 0 THEN 1 ELSE 0 END)</code> |                           |   |
| <b>Result</b>              | 4   |                           |   |

|                            |   |                           |   |
|----------------------------|---|---------------------------|---|
| <b>Metric name</b>         | WAD_OmiOfExp  | <b>Threshold</b>          | 0   |
| <b>Fraud data scheme</b>   | Wrong   | <b>Presentation layer</b> | Activity  |
| <b>Process perspective</b> | Data  | <b>Fraud domain</b>       | Finance,<br>concealed liabilities<br>and expenses,<br>omission of<br>expenses |
| <b>Metric description</b>  | By using "wrong activity data", this specific metric can be defined, which checks to see whether the activity price is equal to zero                      |                           |   |
| <b>Metric formula</b>      | <code>CASE WHEN KPI("Filtered count", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."CUMULATIVE NET WORTH (EUR)" = 0) &gt; 0 THEN 1 ELSE 0 END</code> |                           |   |

|                            |  |                           |  |
|----------------------------|--|---------------------------|--|
| <b>Result</b>              | 237  |                           |  |
| <b>Metric name</b>         | AMR_ExcOfSup   | <b>Threshold</b>          | 0  |
| <b>Fraud data scheme</b>   | Anomalous  | <b>Presentation layer</b> | Model  |
| <b>Process perspective</b> | Resource   | <b>Fraud domain</b>       | Finance,<br>economic extortion,<br>exclusion of<br>specific supplier |
| <b>Metric description</b>  | Checks if the supplier frequency (e.g., #PO and PO value) is sharply decreasing over time                |                           |  |
| <b>Metric formula</b>      | <b>COUNT_TABLE("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES")</b><br>(based on time dimension) |                           |  |
| <b>Result</b>              | 0  |                           |  |

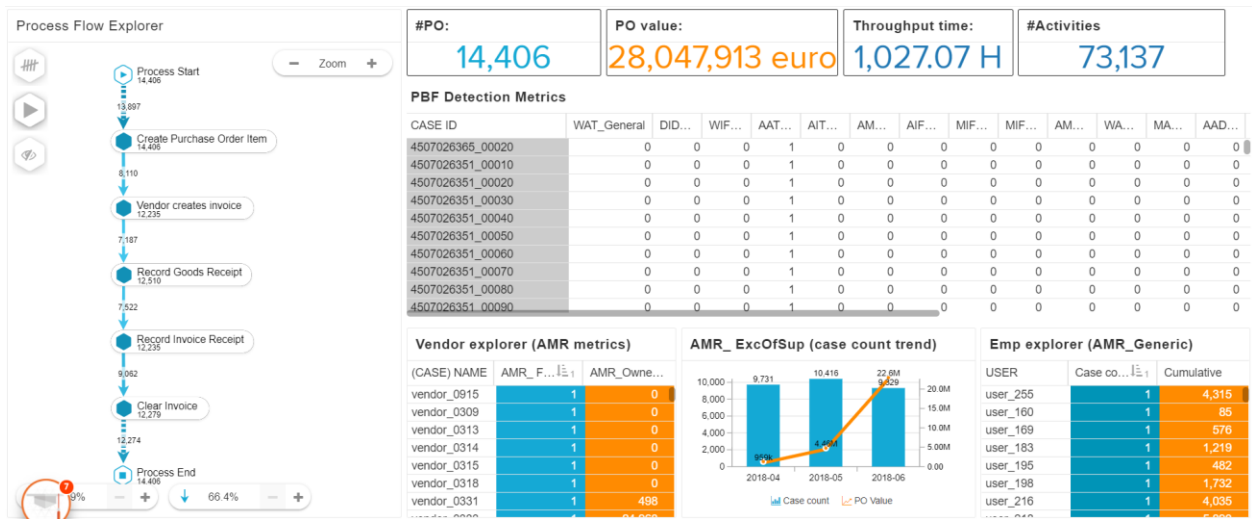
|                            |  |                           |   |
|----------------------------|--|---------------------------|---|
| <b>Metric name</b>         | AMR_OwnershipOfSup   | <b>Threshold</b>          | N/A   |
| <b>Fraud data scheme</b>   | Anomalous  | <b>Presentation layer</b> | Model   |
| <b>Process perspective</b> | Resource   | <b>Fraud domain</b>       | Finance, conflict<br>of interest,<br>ownership of<br>supplier |
| <b>Metric description</b>  | Checks if the supplier frequency (e.g., #PO and PO value) is sharply increasing over time                  |                           |   |
| <b>Metric formula</b>      | <b>COUNT_TABLE("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES")</b><br>(based on vendor dimension) |                           |   |
| <b>Result</b>              | 0  |                           |   |

|                            |   |                           |   |
|----------------------------|---|---------------------------|---|
| <b>Metric name</b>         | AMR_FictitiousSup   | <b>Threshold</b>          | 0, 2  |
| <b>Fraud data scheme</b>   | Anomalous   | <b>Presentation layer</b> | Model   |
| <b>Process perspective</b> | Resource  | <b>Fraud domain</b>       | Finance,<br>fraudulent<br>disbursements of<br>cash, fictitious<br>suppliers |
| <b>Metric description</b>  | By using "Anomalous model resource", this specific metric can be defined, which checks for a supplier that appears only once (showing supplier frequency) |                           |   |
| <b>Metric formula</b>      | <b>COUNT_TABLE("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES")</b><br>(based on vendor dimension)  |                           |   |
| <b>Result</b>              | 155   |                           |   |

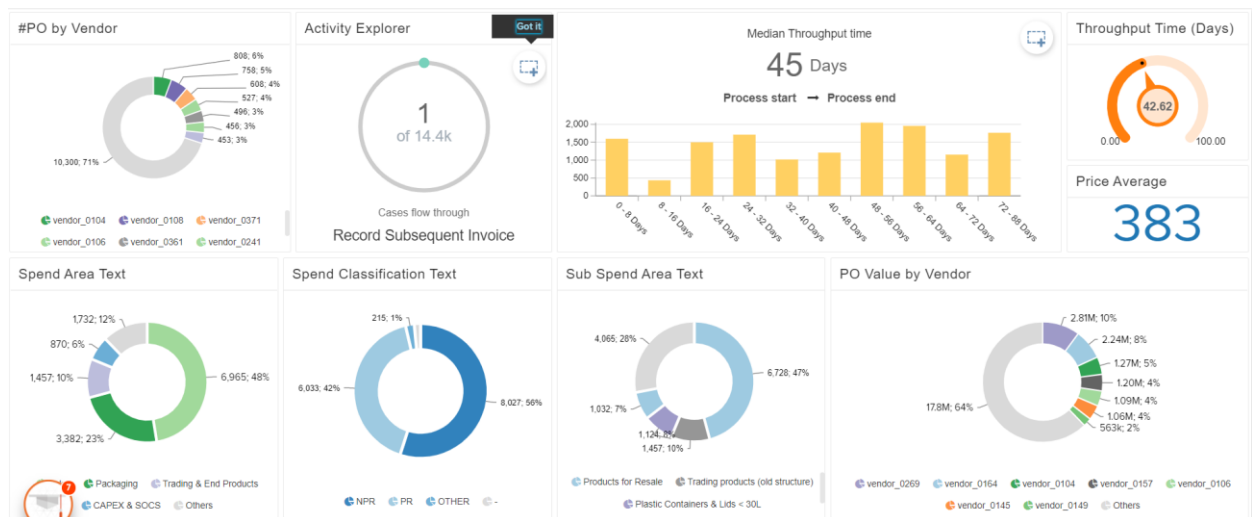
|                            |   |                           |  |
|----------------------------|---|---------------------------|--|
| <b>Metric name</b>         | DIR_PhantomSup  | <b>Threshold</b>          | 0  |
| <b>Fraud data scheme</b>   | Discrepant  | <b>Presentation layer</b> | Instance   |
| <b>Process perspective</b> | Resource  | <b>Fraud domain</b>       | Finance,<br>fictitious expenses,<br>phantom supplier |
| <b>Metric description</b>  | Checks if vendor in "Record Invoice Receipt" is different than that in "Vendor Create Invoice"  |                           |  |
| <b>Metric formula</b>      | <pre> AVG(CASE WHEN PU_COUNT("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY" = 'Record Invoice Receipt') = 0 THEN 0 WHEN PU_FIRST("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."(CASE) NAME", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY" = 'Record Invoice Receipt') &lt;&gt; PU_FIRST("Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy_CASES", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."(CASE) NAME", "Q2_Disco_export_TAX_xlsx_Q2_Disco_export_-_Copy"."ACTIVITY" = 'Vendor creates invoice') THEN 1 ELSE 0 END) </pre> |                           |  |

|                            |   |                           |  |
|----------------------------|---|---------------------------|--|
| <b>Result</b>              | 0   |                           |  |
| <b>Metric name</b>         | AIF_FakeInv   | <b>Threshold</b>          | 0  |
| <b>Fraud data scheme</b>   | Anomalous   | <b>Presentation layer</b> | Instance   |
| <b>Process perspective</b> | Function  | <b>Fraud domain</b>       | Finance,<br>fictitious expenses,<br>fake invoice |
| <b>Metric description</b>  | By using "Anomalous instance function", this specific metric can be defined, which checks activity frequency for "Cancel Invoice Receipt" to determine whether it occurs more than once. This is because fraud may be undertaken by creating fake invoices (e.g., to increase expenses for any reason), which are canceled at a later date. |                           |  |
| <b>Metric formula</b>      | <code>CASE WHEN KPI("Ratio", MATCH_ACTIVITIES(NODE_ANY['Cancel Invoice Receipt'] ) = 1) &gt; 1 THEN 1 ELSE 0 END</code>   |                           |  |
| <b>Result</b>              | 0   |                           |  |

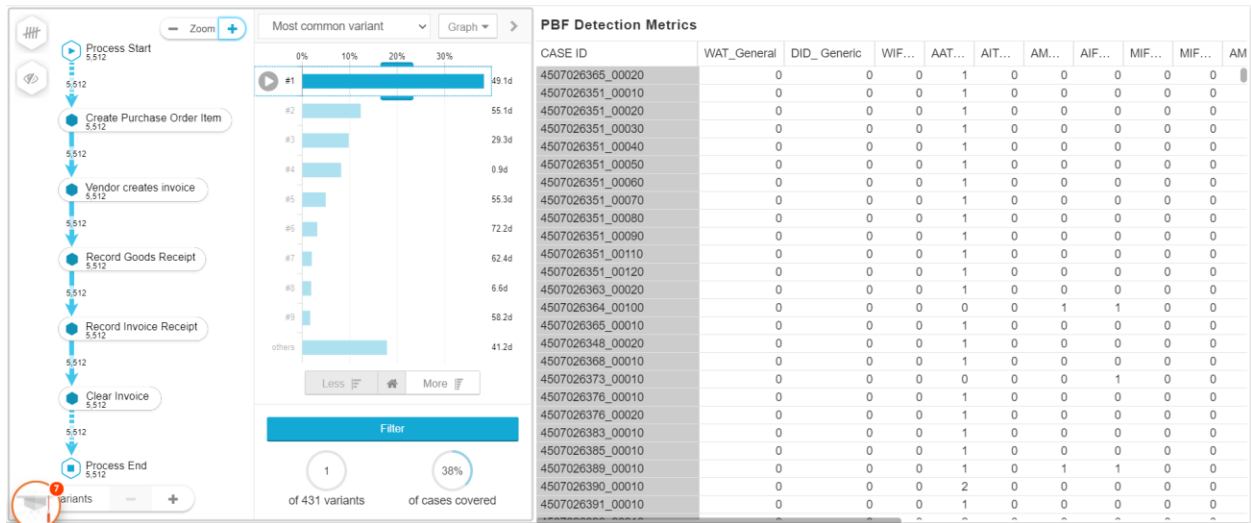
APPENDIX B



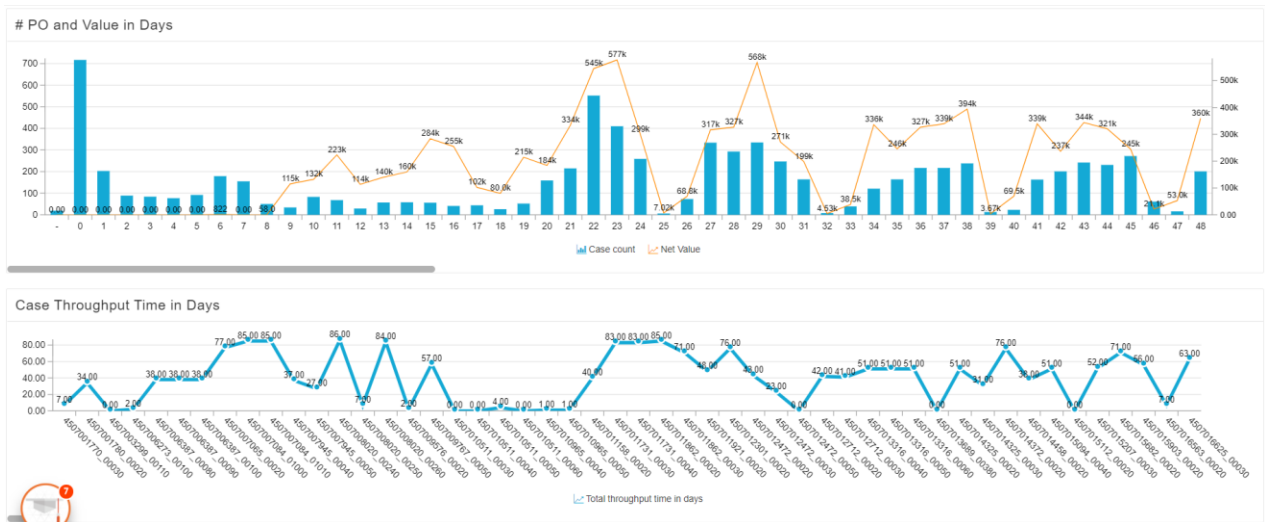
Screen 1. Process view Linked with Data view.



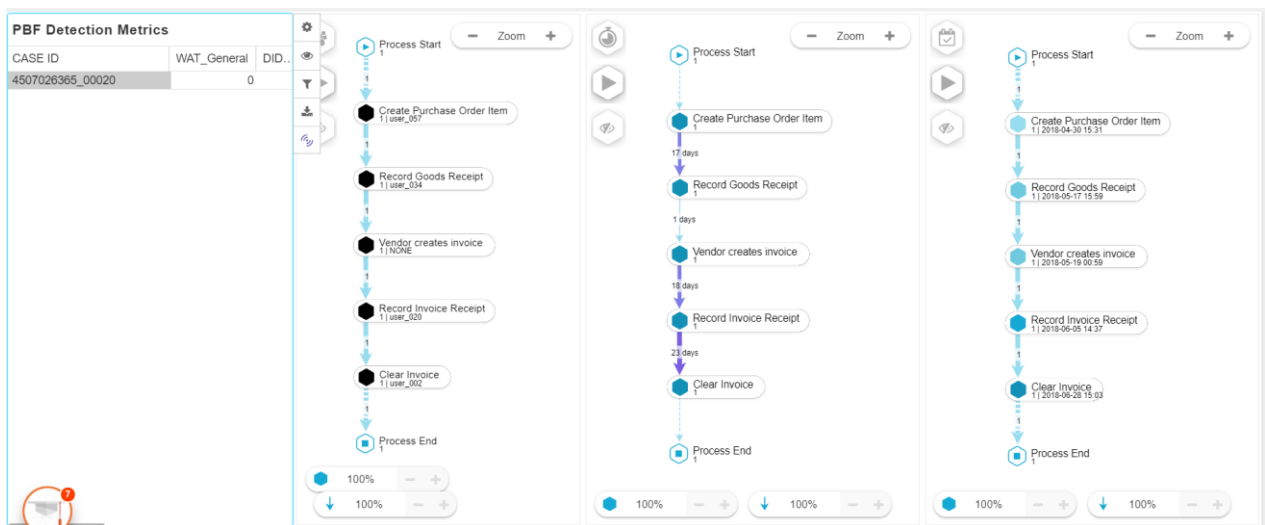
Screen 2. BI Dashboards for Analyzing Process Content.



Screen 3. Process Variant Explorer with Predefined Metrics.



Screen 4. Trend Analysis Dashboards.



Screen 5. An Enriched Process Model for Case id 4507026365\_00020, Showing Activity Frequency, Username, Average Throughput Time, and Execution Timestamp.

# Video Processing for Animation at Key Points of Movement in the Mimosa Pudica

Rodolfo Romero-Herrera<sup>1</sup>

Ciencias e Ingeniería de la Computación  
ESCOM IPN, Ciudad de México

Laura Mendez-Segundo<sup>2</sup>

Ingeniería en sistemas Computacionales  
ESCOM IPN, Ciudad de México

**Abstract**—The processing of an image of a moving plant is inadequate, for this reason, digital video processing must be incorporated, which allows the behavior of an algorithm to be analyzed over time. A method is presented that takes images of a plant with autonomous movement filmed on video; the frames are digitally processed and the information is used to generate animations. Our representation of the structure is derived from an analysis of the image where the plant is deformed; the projections of the movement of the plant are recovered from the video frames and are used as a basis to generate videograms in an animation based on key points taken from an image; Harris and Brisk algorithms are applied. The main plant used is the Mimosa Pudica. Once the frames have been obtained, correlation is proposed as a mechanism to find movement. The techniques are equally useful for any other moving plant such as carnivores or sunflowers.

**Keywords**—Harris; Brisk; correlation; ROI (Region of Interest); Canny. Sobel; Mimosa Pudica; movement

## I. INTRODUCTION

Most people think of plants as inactive [1]. Plants move in different ways and for different reasons [2]. Numerous investigations have led to a deeper understanding of the physiology and biomechanics of these living beings; however, they are not yet fully understood [3]. Because their movements occur in multiple ways [4].

The leaves of Mimosa pudica and many other legumes are characterized by their motor organs that allow the leaves to carry out sleep movements [5] [6]. These actions appear to be regulated by electrical signals and chemical properties [7]. The movements go through four stages: open state, closed state, locked state, semi-open state [8] [9]. The investigations try to discover if these plants have an electrical component [10]; based on the characteristic of cyclic voltage current, where the possibility of memory must manifest itself [11]; in this way, mathematical and electrical-chemical models have been proposed [12].

The study using video shots when the plant is impacted by rain or some other substance can reveal how the system is activated [13]. However, Mimosa pudica rapidly closes its leaves in response to mechanical stimulation [14]; in such a way that the properties of the electrical signals generated can be studied [7], but a video turns out to be a better option because, although they lack muscle, plants have developed a remarkable variety of mechanisms to create movement. Video can show how the plant uses mechanical instability to

accelerate its movements [15]. In the video you do not necessarily see what the eyes observe; Frames can reflect the effect of applying techniques such as the based on classifiers, neural networks, corner detectors, etc. [16] [17].

Many structural and functional properties that plants possess have great potential to stimulate new concepts and innovative ideas in the field of biomimetic engineering. Key from biology can be used to create efficient and optimized structures [5]. The study of plant movement has generated applications such as robots [18] [19], energy harvesting systems [20], development of simple teaching laboratories to illustrate the dynamic qualities of plant movement using smartphones [21]; the realization of tactile sensors inspired by the schismatic movement of plants [22]. The use of animation techniques and their algorithms can lead to a greater number of applications, especially in the area of artificial intelligence [23]. Therefore, the impact of the system with animation has potential aspects. Since the development of a system that animates the movement of the Mimosa pudica allows the analysis in time intervals.

The referenced articles (1 to 23) perform an analysis considering a mechanism that isolates the plant considering an image, such a situation avoids a method that allows reviewing the behavior of the plant in a sequence of time, and with the disadvantage of not being under a natural environment, this a problem is solved in this article through image processing techniques and the use of key points generated from methods such as Harris and Brisk presented in animation.

Another obvious problem is that the processing of the referenced works is applied to isolated plants, which affects the natural behavior of the plant. To avoid the problem, a way to isolate the problem is by using ROI (Region of Interest).

The main objective is to develop a procedure that allows generating animations from keypoints detected in videos of the mimosa pudica plant; The system first performs motion detection and subsequently performs a statistical analysis of the sequence of frames in an animation, without modifying the natural environment in which the Mimosa pudica is found.

The present work aims to carry out video processing applicable to plants with autonomous movement, through the digital treatment of frames of the Mimosa Pudica and to generate animations of the collected data, which allow its evolution to be visualized. In addition to producing a technique based on key points for the analysis of movements. In this way,

the hypothesis that the movement can be detected utilizing the correlation coefficient of frames with key points is supported.

The main motivation is the development of the system that will allow the intelligent behavior of plants to be analyzed in their natural environment, since it will permit the generation of systems based on the movements of the plants, starting with sensors, reaching the generation of artificial intelligence paradigms. based on the results delivered by the system or methodology presented.

## II. METHODS

### A. Incremental Life Cycle

The research methodology used is based on the Incremental Life Cycle of a system. The model is built by increasing functionalities. It was carried out by modules that fulfill different functions within the investigation. The increments allow the capacities of a system to be gradually improved [24].

In the diagram of Fig. 1, the blocks of the procedure are observed. The pre-processing adapts the video frames through filters, enhancement of certain parameters, conversion to gray levels, etc; the part of the interest of the Mimosa Púdica plant is located through ROI (Region of Interest). Following the incremental model, key points are located, by the Harris or Brisk method and an analysis of key points or Landmarks is performed [25]. Once finished, the animation is generated on which the amount of movement is determined. Animations can be generated in various stages; which improves the behavior analysis of the techniques used.

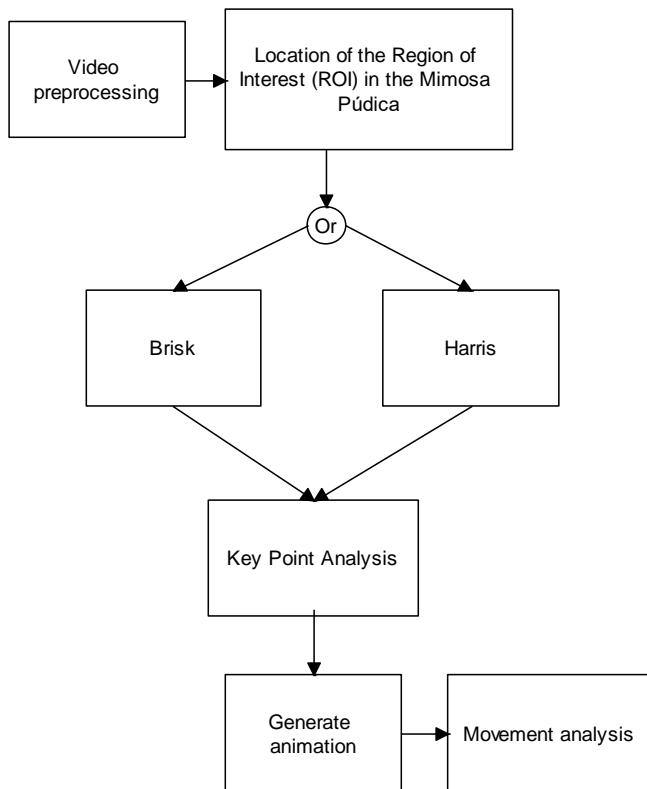


Fig. 1. Affective Pattern Recognition and Learning Systems.

### B. Region of Interest

A region of interest (ROI) is a part of an image that you want to filter or operate on in some way. To determine the area you can use many shapes, such as circles, ellipses, polygons, rectangles; hand-drawn shapes were chosen. In the present project, it was used to create a binary over-face image [26].

### C. Sobel

The Sobel operator handles a matrix where the central row or column of the filter is given greater weight [27]. Said matrix is defined in equation (1).

$$H_x^s = \begin{pmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{pmatrix}$$

and

$$H_y^s = \begin{pmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ 1 & 2 & 1 \end{pmatrix} \quad (1)$$

### D. Canny

The method applies several filters in different directions and resolutions, which are combined. The method has three objectives: Minimize the number of false, Improve edge localization, deliver a border of one-pixel width. The technique is based on the gradient and the second derivative of Laplacian [27].

### E. BRISK Corner Detector (Binary Robust Invariable Scalable Key Points)

The detector can be used to locate corners with multiple scales. The goal is to decompose an image into regions of local interest; which reduces complexity while exploiting appearance properties. The hotspot detector finds regions in the image that stands out even when the viewer of point of observation changes. See Fig. 2. BRISK has three phases, the detection of the characteristic, the composition of the descriptor, and finally the pairing of the key points [28].

The technique is applied during the sampling of the intensity of the image at the standard  $p_i$  point, and a Gaussian smoothing with standard deviation  $\sigma_i$  proportional to the distance between the points of the respective circle [22]. Also, the pattern is positioned and scaled according to the key point  $k$ , in one of the  $(N(N-1))/2$  pairs of sampled points  $p_i, p_j$ . The smoothed intensity values at these points are  $I(p_i, \sigma_i)$  and  $I(p_j, \sigma_j)$  respectively, and are used to estimate the gradient  $g(p_i, p_j)$ . See equation 2.

$$g(p_i, p_j) = (p_j, p_i) \frac{I(p_i, \sigma_i) - I(p_j, \sigma_j)}{\|p_j - p_i\|^2} \quad (2)$$

The distance thresholds are  $\delta_{\max} = 9.75t$  and  $\delta_{\min} = 13.67t$  with a  $t(k)$  scale. Iterating through the even points in  $L$ , the direction over all the characteristics of the pattern of key point  $k$  is estimated, using equation 3.

$$g = \begin{pmatrix} g_x \\ g_y \end{pmatrix} = \frac{1}{L} \sum_{(p_i, p_j) \in L} g(p_i, p_j) \quad (3)$$

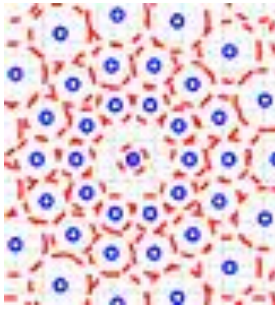


Fig. 2. Sampling Pattern for the BRISK Descriptor [28].

### F. Harris Detector

A corner can be interpreted as the junction of two edges, where one edge is a sudden change in the brightness of the image. The Harris algorithm detects corner points, regardless of rotation and change in gray level [29]. The response is used to select isolated corner pixels, and reduce the pixels at the edge. The Harris matrix is a symmetric matrix similar to a covariance matrix. The main diagonal is made up of the two averages of the square gradients "y". See equation 4. The off-diagonal elements are the averages of the cross product of the gradient  $\langle G_{xy} \rangle$ .

$$A_{Harris} = \begin{bmatrix} \langle G_x^2 \rangle & \langle G_{xy} \rangle \\ \langle G_{xy} \rangle & \langle G_y^2 \rangle \end{bmatrix} \quad (4)$$

First, the measure of the response of the corner  $R$  is considered. The contours of the constant  $R$  are shown by thin lines.  $R$  is positive in the corner region, negative in the border regions, and small in the flat region. The  $R$  values increase the contrast and magnitude of the response. The flat region is specified by  $Tr$ , which falls below some selected threshold.  $R$  can be determined by equation 5.

$$R = \det(A_{Harris}) - k Tr^2(A_{Harris}) \quad (5)$$

where  $k$  is a constant typically with a value of 0.04;  $R$  can also be expressed with gradients:

$$R = (\langle G_x^2 \rangle \langle G_y^2 \rangle - \langle G_{xy} \rangle^2) - k(\langle G_x^2 \rangle + \langle G_y^2 \rangle)^2 \quad (6)$$

So when the response is greater than a predefined threshold, a corner is detected. See equation (7).

$$R > k_{thresh}$$

$$(\langle G_x^2 \rangle \langle G_y^2 \rangle - \langle G_{xy} \rangle^2) - k(\langle G_x^2 \rangle + \langle G_y^2 \rangle)^2 > k_{thresh} \quad (7)$$

In this way, a pixel in the corner region is selected (positive response) if its response is a local maximum of 8 axes. Pixels in the edge region are considered boundaries if their responses are local minima and negative in the  $x$  or  $y$  directions; depending on, if the magnitude of the first gradient is greater in the  $x$  or  $y$ -direction. The result is continuous thin edges that generally end in corner regions.

### G. Animations

An environmental animation algorithm is a laboratory for the dynamic investigation of programs. Multiple graphs of an algorithm in action can be presented, exposing program

properties that in other cases would be difficult to understand or notice.

To animate an algorithm, the frame of reference must be prepared. The framework works as a test, connecting the algorithm to an input generator, and monitoring both input and internal events, displaying the action in one or more views. A videogram is an image of an abstract state of the algorithm or its data, handled by events [30]. Frames show relevant parameters and variables, of the current state. An algorithm in the video becomes clearer. The idea is to help humans quickly get a clear picture of what is happening. Properties can be intuited that can be verified with more formal methods [31].

Steps to create a frame animation:

- 1) Run a simulation or generate data.
- 2) Draw / Render the stage in time  $t_k$ .
- 3) Take a snapshot of the scene
- 4) Advance time  $t_k$  to  $t_k + 15$ . Save movie

### H. Correlation

The covariance of two random variables  $X$  and  $\gamma$ , with a joint probability density function  $f(x, y)$ , is defined as:

$$Cov(X, Y) = \sigma_{x,y} = E[(X - \mu_x)(Y - \mu_y)] \quad (8)$$

Where  $\mu$  is the mean and  $E$  is the mathematical expectation. The correlation coefficient is given by:

$$Corr(X, \gamma) = \rho_{x,\gamma} = \frac{Cov(X, \gamma)}{\sigma_x \sigma_\gamma} = \frac{\sigma_{x,\gamma}}{\sigma_x \sigma_\gamma} \quad (9)$$

Where  $\sigma_x > 0$  y  $\sigma_\gamma > 0$

Correlation is a measure of the linear relationship between two random variables. If the joint distribution of two variables has a correlation coefficient, then  $-1 \leq \rho_{x,\gamma} \leq 1$ . When  $\rho_{x,\gamma} = 1$ , then  $X$  and  $\gamma$  are perfectly related positively. The conclusion is that the possible values of  $X$  and  $\gamma$  lie on a line with a positive slope. On the other hand, when  $\rho_{x,\gamma} = -1$  then the situation is opposite:  $X$  and  $\gamma$  are perfectly negatively correlated. If  $X$  and  $\gamma$  are independent, then  $\rho_{x,\gamma} = 0$  [32].

## III. RESULTS

### A. Active Contours vs. Region of Interest

One of the main problems for the analysis of the mimosa pudica is that is immersed within an environment whose image segmentation makes a problem is not solved; since if the segmentation is done by color, the algorithm is confused when finds similarities between the other sheets and backgrounds with a similar color. Thus, the first option that was considered was a semiautomatic segmentation. The original snake model of the active contour proposed by Kass [33], is represented as a parameterized curve  $v[s] = [x(s), y(s)]$ ,  $s \in [0, 1]$  that moves through a spatial domain and that seeks to minimize the following energy functional.

$$E_{snake}^* = \int_0^1 E_{snake}(v(s)) ds \quad (10)$$

Active Contours seems like a good option, but the results show the opposite, as made known in Fig. 3 and 4. The situation is difficult to analyze in the image in Fig. 5 because



even for one person it is difficult to separate the branch to observe. In Fig. 4, the cut leaf can be detected since the selected area (in yellow line) is modified due to the algorithm of the snake. This fact impairs the analysis since part of the movement of the plant is truncated. The intrusion of other objects is again due to the similarity of patterns with the rest of the leaves and the surrounding environment, together with the fact that the leaf to be analyzed is moving. That is, the area of energy covered by the integral of equation (10) is modified, which is easily observable in animation with frames. For this reason is better to use a region of interest by freehand limiting.

The binary mask created allows us to isolate the region to be analyzed using ROI. The area remains throughout the cycle of movements of the plant. Therefore, the edges should be chosen with excess, but without going too far from the extremes to be analyzed. Trying to cover the entire movement. As seen in Fig. 5. In this way, the problem of analyzing the plant in its natural environment is solved without affecting the conditions that we wish to review. In Fig. 5 can be seen that the frame covers only the area to be analyzed without invading other regions, so this method gives better results than the snake algorithm. Fig. 5 shows the intermediate Frame of the animation sequence.

If the key point detection algorithms are applied to the pudic Mimosa without using ROI, the results shown in Tables I and II are obtained. Several videos were processed with shots of the Mimosa; the results were similar to those presented in Tables I and II. Table I shows the Mimosa with the leaves practically horizontally. It is not possible to observe the details that happen over time in a single image of the added value of the animation. Using the Brisk algorithm, motion detection is easy, as more circles appear as the sheets are folded. Finally, the Harris algorithm detects corner points that change position, but they are not enough. The case presented in Table I is the worst-case evaluated.

Table II shows a case where the Mimosa leaf is facing the camera. For the Brisk algorithm, although letter-shaped noise is detected, such fact is easy to observe the accumulation of detections at points where the sheet is folded. Finally, when applying the Harris algorithm, the stem and green points are detected at the tips of the leaves, reducing the distance between points as the movement evolves; Tables I and II show the utility of these methods used; however, a single image does not allow us to observe the behavior of the plant concerning the passage of time; therefore, an animation is essential. However, the video generation consumes several computing resources; therefore, only the key points should be used rather than the entire image.

The primary results would indicate that is feasible to analyze the movement of the plant just by using the key point detectors. However, is not the case, since the problem arises when the branch to be analyzed is found in clusters of mimosa; Separating the plant with a mechanism modifies the normal

conditions and therefore changes the behavior of the plant. ROI is used to solve the problem. The results are shown in Table III.

Before the detection of key points, digital filters are applied and we have Fig. 4. The Sobel and Canny techniques improve the detection of key points, which is best observed with animation and that would be difficult to notice with a single image or with a set of them, especially if is considered that the conditions are changing over time. This fact is observed in Table III for the cases presented in the second 1, 2.5, and 4 when applying Brisk to the video.

Results are better when applying Canny Edge Detector. Table IV shows the results when applying the Harris detector with Canny. The second 1 is the case when the movement starts, the second 2.5 is the middle of the movement and the second 4 is the end of the movement. You can see the similarities and differences in the detected points.



Fig. 3. Original Image [34].



Fig. 4. Snake Method.



Fig. 5. ROI Method.

TABLE I. RESULTS OF MIMOSA PROCESSING USING DIFFERENT ALGORITHMS

|  |  |  |  |                |
|--|--|--|--|----------------|
|  |  |  |  | Original Image |
|  |  |  |  | Brisk          |
|  |  |  |  | Harris         |

TABLE II. RESULTS OF MIMOSA PROCESSING USING DIFFERENT ALGORITHMS, CASE 2

|  |  |  |  |                 |
|--|--|--|--|-----------------|
|  |  |  |  | Imagen Original |
|  |  |  |  | Brisk           |
|  |  |  |  | Harris          |

TABLE III. KEY POINTS WERE DETECTED

| Tiempo                 | 1 seg | 2.5 seg | 4 seg |
|------------------------|-------|---------|-------|
| Sobel<br>Fotograma<br> |       |         |       |
| Canny<br>Fotograma<br> |       |         |       |

TABLE IV. KEY POINTS WITH CANNY EDGE DETECTION

| Tiempo        | 1 seg | 2.5 seg | 4 seg |
|---------------|-------|---------|-------|
| Fotograma<br> |       |         |       |

## B. Statistical and Correlation Analysis

### Brisk

The detection of key points by the Brisk technique was applied to videos of Mimosa Púdica. The first video frame to be analyzed is shown in Fig. 6. The Region of Interest exposed in Fig. 7 will be processed.

The area shown by Fig. 7 is chosen, and from this selection, the frames that form the animation are generated. Thanks to the generation of the animation are possible to perform the statistical analysis during the time interval in which the movement of the plant lasts.

Animation with Brisk key points was obtained; the final frame is shown in Fig. 8. A single 3-second video can generate more than 300 images. In the same way, as with the Harris method is observed that the amount of information to be processed allows the saving of computing resources since the whole image or its contours do not have to be processed, the key points are enough.



Fig. 6. Image to Process.

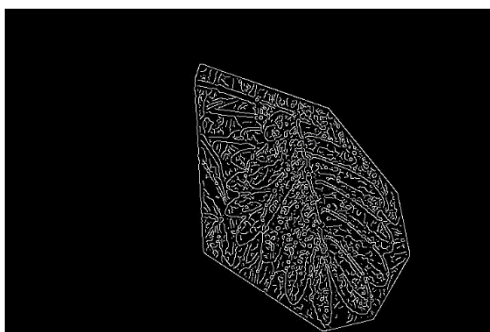


Fig. 7. ROI using Brisk Technique.

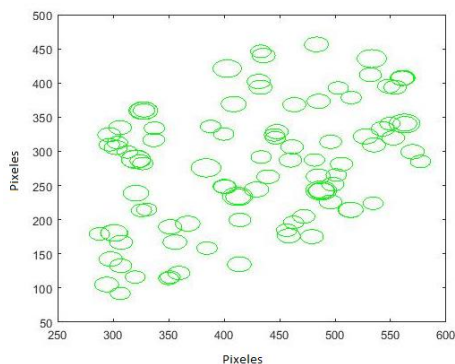


Fig. 8. The Final Frame of the Animation.

Fig. 9 shows the result of the application of the basic statistics of the key points. Obtaining the correlation coefficient. Changes that reflect the movement of the plant are observed in the graph.

Table V presents the analysis of the graph. A value equal to 1 or -1 is never observed, so there is a linear dependence between images, but obviously, the movements are not sudden. The maximum value is 0.62 for frame 90. The minimum has a value of 0.31 because there is no movement.

### Harris

The video processing was also performed by applying the Harris corner detector to a square area of the leaf frame of the Mimosa Púdica. As shown in Fig. 10. Only the contours generated by Canny are shown in Fig. 10.

Fig. 11 shows the key points of the final frame. With frames of 400 x 400 pixels. The amount of information processed is much less than if the original image were processed, so computing resources are saved; Since we have the position of the key points, the statistics of the time interval that the video lasts are taken.

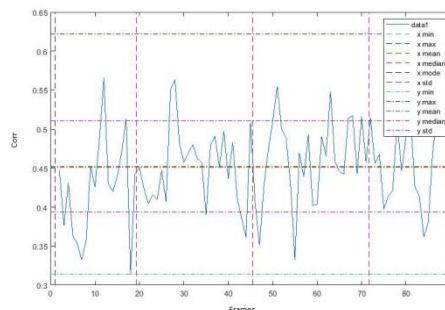


Fig. 9. Correlation Plot.

TABLE V. BASIC STATISTICS OF CORRELATION

|       | Fotograma | Correlación |
|-------|-----------|-------------|
| min   | 1         | 0.3135      |
| max   | 90        | 0.6217      |
| Mean  | 45.5      | 0.4524      |
| Mode  | 1         | 0.3135      |
| std   | 26.12     | 0.05869     |
| range | 89        | 0.3082      |

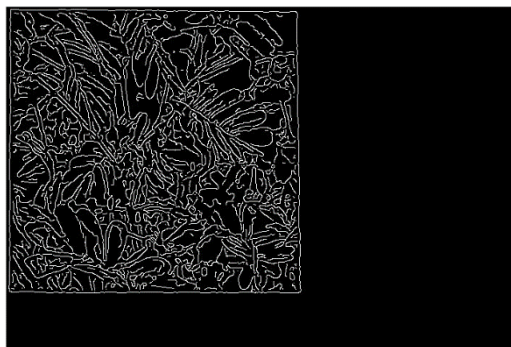


Fig. 10. ROI by applying Harris Corner Detector.

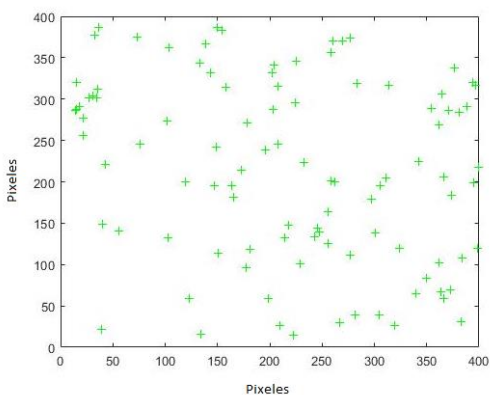


Fig. 11. Harris Key Point.

The graph in Fig. 12 shows the basic statistic points between frames. As there is a greater number of *Mimosa pudica* plants, more chaotic movements are shown where the correlation coefficient can range from -0.14 to 0.28 as shown in Table VI.

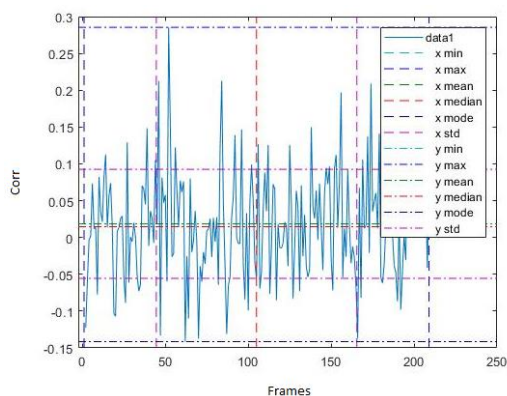


Fig. 12. Graph of the basic Statistics of the Correlation.

TABLE VI. BASIC STATISTICS OF THE CORRELATION

|       | X     | Y       |
|-------|-------|---------|
| min   | 1     | -0.1415 |
| max   | 209   | 0.2854  |
| Mean  | 105   | 0.0185  |
| Mode  | 1     | -0.1415 |
| std   | 60.48 | 0.7399  |
| range | 208   | 0.4269  |

#### IV. CONCLUSIONS

The processing of an image of a moving plant is inadequate, for this reason, digital video processing must be incorporated, and the generation of animations facilitate the review of the behavior of an algorithm over time, solving the problem of having a large magnitude of time-consuming images compared to animation. An analysis was proposed employing a non-autonomous segmentation (Region of interest) that avoids that when the conditions change over time, the segmentation of the information is out of control. The referenced articles isolate the plant from its natural

environment, with the proposal it is feasible to carry out the analysis within the natural environment of the plant without using mechanisms or systems that affect the behavior or the analysis.

We comply to develop a procedure that allows generating animations from Harris or Brisk key points of the pudic mimosa plant, where motion detection was first achieved, and later a statistical analysis of the sequence of frames in animation is performed. The use of key points as data to be analyzed is another contribution, since it greatly reduces the computation used, unlike the referenced articles.

As for the key point detection algorithms, we can take the information and perform simulations and/or animations based on the movement of the plants. The ROI application manages to reduce the amount of information to be analyzed and improves the processing speed of landmarks.

Regarding the analysis carried out as a research tool, the application of the detection of key points can yield interesting results that can lead to applications of biomimetic impact, the development of sensors, and the use of movement mechanics for the generation of small energy sources.

The correlation coefficient presents us with an alternative for the detection and analysis of the movement of the *Mimosa Pudica*. As the correlation coefficient increases, the amount of movement decreases.

Regarding animations, the relative execution of different algorithms becomes clearer and improves their potential. The idea is to help humans quickly acquire a clear picture of what is happening. Due to the exploration of the algorithm in a visual way, you can intuit properties that are verified with methods more formal.

As future work, using the implemented algorithms, the behavior of the *Mimosa pudica* plants will be observed for the generation of paradigms in artificial intelligence based on the behavior of the plants.

#### ACKNOWLEDGMENT

The authors acknowledge the support received for the development of the research project to the IPN (Instituto Politécnico Nacional).

#### REFERENCES

- [1] Jae Young Kim, Young-Joon Park, June-Hee Lee & Chung-Mo Park Developmental polarity shapes thermo-induced nastic movements in plants, *Plant Signaling & Behavior*, 14:8, DOI: 10.1080/15592324.2019.1617609, 2019.
- [2] L. Wagner et al, "The plant leaf movement analyzer (PALMA): a simple tool for the analysis of periodic cotyledon and leaf movement in *Arabidopsis thaliana*," *Plant Methods*, vol. 13, (1), pp. 2, 2017.
- [3] S. Poppinga et al, "Biomechanical analysis of prey capture in the carnivorous Southern bladderwort (*Utricularia australis*)," *Scientific Reports*, vol. 7, (1), pp. 1776-10, 2017.
- [4] Guo, Qiaohang, et al. "Fast nastic motion of plants and bioinspired structures." *Journal of the Royal Society Interface* 12.110, 2015: 20150598.
- [5] S. Sugito et al, "An analysis of mimosa pudica leaves movement by using LoggerPro software," in 2016, DOI: 10.1088/1742-6596/739/1/012121.

- [6] G. Muhammad et al, "Mimosa pudica L., a High-Value Medicinal Plant as a Source of Bioactives for Pharmaceuticals," *Comprehensive Reviews in Food Science and Food Safety*, vol. 15, (2), pp. 303-315, 2016.
- [7] ] Z. N. Ismarrubie et al, "Bio-Mechanism Response of Mimosa Pudica against External Stimulation," *Advanced Materials Research*, vol. 1125, pp. 588-592, 2015.
- [8] Volkov, Alexander G. "Signaling in electrical networks of the Venus flytrap (*Dionaea muscipula* Ellis)." *Bioelectrochemist* 125", pp: 25-32, 2019.
- [9] Venus Flytrap: How an Excitable, Carnivorous Plant Works por Hedrich, Rainer; Neher, Erwin *Trends in Plant Science*, 03/2018, Volumen 23, Número 3.
- [10] A. G. Volkov et al, "Memristors in plants," *Plant Signaling & Behavior*, vol. 9, (3), pp. e28152, 2014.
- [11] S. Poppinga et al, "Biomechanical analysis of prey capture in the carnivorous Southern bladderwort (*Utricularia australis*)," *Scientific Reports*, vol. 7, (1), pp. 1776-10, 2017.
- [12] ] J. Gim and C. Ahn, "Design and Analysis of Osmosis-based Artificial Muscle," *Journal of Bionic Engineering*, vol. 16, (1), pp. 56-65, 2019.
- [13] U. Bauer et al, "Mechanism for rapid passive-dynamic prey capture in a pitcher plant," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 112, (43), pp. 13384-13389, 2015.
- [14] K. Tagawa, M. Watanabe and T. Yahara, "A sensitive flower: mechanical stimulation induces rapid flower closure in *Drosera* spp. (*Droseraceae*)," *Plant Species Biology*, vol. 33, (2), pp. 153-157, 2018.
- [15] J. Edwards et al, "The Role of Water in Fast Plant Movements," *Integrative and Comparative Biology*, vol. 59, (6), pp. 1525-1534, 2019.
- [16] E. Yigit et al, "A study on visual features of leaves in plant identification using artificial intelligence techniques," *Computers and Electronics in Agriculture*, vol. 156, pp. 369-377, 2019.
- [17] Brattland, Vegard, et al. "Image processing of leaf movements in *Mimosa pudica*." *Scandinavian Conference on Image Analysis*. Springer, Cham, 2017.
- [18] L. Shi et al, "Development and evaluation of a Venus flytrap-inspired microrobot," *Microsystem Technologies*, vol. 22, (8), pp. 1949-1958, 2016.
- [19] Chang, Chung-Liang, and Jin-Long Shie. "Design and implementation of a bionic mimosa robot with delicate leaf swing behavior." *Micromachines* , 42-62, 2015.
- [20] B. Ali et al, "Fuzzy Logic based Energy Harvesting with the Movement of Plants Branches and Leaves," *Pakistan Journal of Agricultural Sciences*, vol. 53, (2), pp. 449-454, 2016.
- [21] E. D. Brenner, "Smartphones for Teaching Plant Movement," *The American Biology Teacher*, vol. 79, (9), pp. 740- 745, 2017.
- [22] BASIR, Siti Nora, et al. design concept of a new bio-inspired tactile sensor based on main pulvinus motor organ cells distribution of *Mimosa Pudica* plant. En 2014 International Symposium on Micro-NanoMechatronics and Human Science (MHS). IEEE, 2014. p. 1-6.
- [23] Hu, Tao, et al. "Generating Video Animation from Single Still Image in Social Media Based on Intelligent Computing." *Journal of Visual Communication and Image Representation* : 102812, 2020.
- [24] Ramos, Daniel, et al. *Curso de Ingeniería de Software: 2ª Edición*. IT Campus Academy, 2017.
- [25] Y. Wu and Q. Ji, "Facial Landmark Detection: A Literature Survey," *International Journal of Computer Vision*, vol. 127, (2), pp. 1-28, 2018;2019.
- [26] Greenham, Kathleen, et al. "TRiP: Tracking Rhythms in Plants, an automated leaf movement analysis program for circadian period estimation." *Plant methods* 11.1 , 2015: 33.
- [27] Chacón Murguía, Mario I. Percepción visual-Aplicada a la robótica. Alfaomega Grupo Editor, 2016.
- [28] Leutenegger, S., M. Chli and R. Siegwart. "BRISK: Binary Robust Invariant Scalable Keypoints", *Proceedings of the IEEE International Conference, ICCV*, 2011.
- [29] Rodolfo Romero, Francisco Gallegos, José Elias Romero, "Video Analysis with Faces using Harris Detector", *Journal of Advanced Computer Science and Applications (IJACSA)*, Volume 10 Issue 8, 2019.
- [30] Zahra, Beenish. "Algorithm and Technique for Animation." *LGURJCSIT* 1.1 .pp: 22-36, 2017.
- [31] Végh, Ladislav, and Veronika Stoffová. "Algorithm Animations for Teaching and Learning the Main Ideas of Basic Sortings." *Informatics in Education* 16.1 pp: 121-140, 2017.
- [32] Martínez, W. L., and A. R. Martínez. "Computational Statistics Handbook with MATLAB, 124-132." ,2016.
- [33] ] M. Deriche, A. Amin and M. Qureshi, "Color image segmentation by combining the convex active contour and the Chan Vese model," *Pattern Analysis and Applications : PAA*, vol. 22, (2), pp. 343-357, 2017;2019.

# Population based Optimized and Condensed Fuzzy Deep Belief Network for Credit Card Fraudulent Detection

Jisha M.V<sup>1</sup>

Research Scholar  
Dept. of Computer Science  
Nehru arts and Science College  
Coimbatore, Tamilnadu, India

D. Vimal Kumar<sup>2</sup>

Associate Professor  
Dept. of Computer Science  
Nehru arts and Science College  
Coimbatore, Tamilnadu, India

**Abstract**—In this information era, with the advancement in technology, there is a high risk due to financial fraud which is a continually increasing menace during online transactions. Credit card fraudulent identification is a toughest challenge because of two important issues, as the profile of the credit card user's behavior changes constantly and credit card datasets are skewed. The factors which greatly affects the credit card fraudulent transaction detection are primarily based on data sampling models, features involved in feature selection and detection approaches implied. To overwhelm these issues, instead of using certainty theory, this paper encapsulates with three different empowered models are deployed for intellectual way of fraudulent transaction detection. In this work uncertainty theory of intuitionistic fuzzy theorem to determine the significant features which will influence the detection process effectively. Maximized relevancy among dependent and independent features of credit card dataset are determined using grade of membership and non-membership information of each features. The intuitionistic fuzzy mutual information with the knowledge of entropy it selects the features with highest information score as significant feature subset. This proposed model devised Fuzzy Deep Belief Network enriched with Sea Turtle Foraging for credit card fraudulent detection (EFDBN-STFA). The fuzzy deep belief network greatly handles the complex pattern of credit card transactions with its deep knowledge and stacked restricted Boltzmann machine the pattern of dataset is analyzed. The weights assigned to the hidden nodes are fine-tuned by the sea turtle foraging using its fitness measure and thus it improves the detection accuracy of the FDBN. Simulation results proved the efficacy of EFDBN-STFA on two different credit card datasets with its gained ability of handling hesitation factor and optimization using metaheuristic approach, it achieves higher detection rate with reduced false alarms compared to other existing detection models.

**Keywords**—Credit card fraudulent; uncertainty; intuitionistic fuzzy; fuzzy deep belief network; sea turtle foraging

## I. INTRODUCTION

In modern days, usage of internet for commercial transactions started increasing exponentially because of its availability and flexibility. Usage of credit cards for online or offline transactions is very useful for business peoples [1]. But credit card fraud becomes significant issue in financial sectors, banks and card issuers. This credit card fraud detection

becomes an important and interesting topic of research for the scientific community. To handle such voluminous transaction system, it needs high sophisticated security system to analyze transactions and detect fraud transactions more quickly. This necessitates with the advantage of modern technology utilizing machine learning, mining and artificial intelligence influence's the process of credit card detection more accurately. Detection of fraud has become an important activity which involves in reducing the impact of fraudulent transactions [2]. But credit card fraud detection is very challenging from the perspective of learning process due to its nature of class imbalance.

Generally, classification models are deployed to examine all the approved transactions and alert the utmost suspicious ones. The alerts are investigated by the professionals and intimate the cardholders to discover whether it is genuine or fraudulent transactions for each altered transaction [3]. This provides feedback to the classification system during the training phase to update itself and detect the fraud detection more accurately. Meanwhile, it is very essential to highlight the major difference among user behavior and fraud analysis. The fraud detection models extract the signature of fault tricks pattern and it greatly assist during the testing process. The ultimate goal of this work is to reduce the false detection of fraudulent detection in a more precise way. This is achieved by developing a Population based Optimized and Condensed Fuzzy Deep Belief Network for Credit Card Fraudulent Detection.

The rest of the research is organized as follows: Section I gives the importance of detection of frauds, Section II highlights the related works done before. Section III describes the methodology of the proposed work. Section IV represents the conceptual framework used and Section V gives the results and its discussion. Followed by the conclusion and references.

## II. RELATED WORK

This section discusses about some of the existing work which involves in credit card fraudulent detection using machine learning algorithms and mining approaches.

Vaishnavi et al. [4] in their work anticipated a novel approach for fraud detection on streaming transaction data,

which analysis past history of customers transaction details. The behaviour patterns of transaction are extracted and card holders with same patterns are clustered depending on their amount of transaction. They used sliding window concept and transactions are aggregate to determine fraudulent and genuine transaction.

Imane et al. [5] developed a comparative model which comprised of various approaches of machine learning models are deployed for credit card fraud detection. The authors mainly focused on investigating neural network performance. They stated that this study aims to guide the researches to choose best approaches for credit card fraud detection.

Wen-Fang et al. [6] presented an outlier mining model to accurately forecast fraudulent credit card transaction. Distance summing algorithm is used to emulate variation among normal and fraud detection. The outlier approach is mainly used to detect anomalous transactions.

Maniraj et al. [7] devised a recognition model to check whether a new incoming transaction normal or fraudulent. They performed preprocessing and analyzed PCA converted credit card transaction data. They deployed isolation forest model and local outlier detection for classifying multiple type of anomaly detection.

Andrea et al. [8] in their work contributed three different approaches to discover fraudulent transactions. They handled class imbalance problem by designing a novel learning policy. They worked with real time dataset with the concept of drifting and verification.

Navneet et al. [9] investigates fraudulent transactions in banking sectors and analyzed vulnerabilities during online transaction. This work explores possible ways to prevent fraud transaction by developing graph database and finding the patterns of fraud transaction.

Salvatore et al. [10] in their work stated that understanding the purpose of meta learning policies will greatly influence during the process of fraud catching rate and deduction rate. They used skewed distribution to work with balanced data during training that results in better classification.

Dheepa and Dhanapal [11] developed a behavioral classification model using support vector machines. The features are extracted to determine the significant behavioral transaction patterns. If there any conflict occurs then it is examined as suspicious and this is considered to discover the frauds.

Chuang et al. [12] presented a mining model, which uses web services on online bank transaction. The banks which are involved in these scheme shares their knowledge about fraud patterns in a heterogenous environment and with the distributed system it further improves the ability of fraud detection with less financial loss.

Tao Guo et al. [13] developed a neural network model to discover customer's behavior pattern. The significant task is to discover any deviation from the usual transaction pattern. This is achieved by training the neural network with dataset and their confident value is computed. Those credit card

transactions with less confident value is treated as fraudulent transaction.

Suvasini Panigrahi et al. [14] in their work designed a fusion model which comprised of four components like Dempster-Shafer, rule-based filter, transaction history and Bayesian model. Rule filter is used to determine fraudulent transactions. DST is used to compute belief value of each transaction based on its evidence value.

### III. METHODOLOGY OF POPULATION BASED OPTIMIZED AND CONDENSED FUZZY DEEP BELIEF NETWORK FOR CREDIT CARD FRAUDULENT DETECTION

This proposed work aims to overcome the ambiguity, vagueness and uncertainty in prediction of credit card fraudulent detection shown in "Fig. 1". This work used two different datasets where one is collected from Kaggle repository and another dataset is collected with a case study of a specific bank. The existing models use the neural networks, support vector machine, random forest and other conventional classification models to perform this fraudulent detection process. Most of them fail to concentrate on handling vagueness and inconsistency which often arise in the real time dataset when there are transactions which cannot be finitely defined either as fraudulent or normal transaction. To overcome this problem the proposed model works in two stages, in order to reduce the redundancy among features involved in fraudulent detection and increase the relevancy among feature and class. This is achieved by adapting intuitionistic fuzzy mutual information as feature subset selection, whose ultimate goal is to choose the most significant attributes involved in process of prediction. The pattern of credit card transactions is analyzed in depth by using fuzzy Deep belief network which is fine tuned by introducing sea turtle foraging algorithm which optimizes the assignment of weights in FDBN.

#### A. Dataset Description

This work used two different credit card datasets for fraudulent transaction detection. The first dataset comprised of 284,807 transactions [15]. The input variables are in PCA transformation and they are denoted as V1, V2, ... V28 vector values. Other variables are time, amount is not transformed and class is a feature which is considered as a target variable. The second credit card dataset comprised of 30,000 transactions with 25 features including the class variable [16]. The features of this dataset are limit value, gender, education, marital status, age, payment and billing details. This proposed work uses these two credit card datasets to discover the fraudulent transaction.

#### B. Normalized Intuitionistic Fuzzy Mutual Information

A searching procedure which selects a subset of features which greatly influence the classification process is known as feature subset generation. This method applies a subset evaluation function to assess the current subset, if the present subset performs better than the previous subset, the current subset is replaced with the previous set [17]. This is a cyclic process which repeats the subset generation and evaluation until termination condition is met. The termination condition relies on both evaluation and generation function, the former

case the iteration terminates when the insertion or deletion of an attribute doesn't produce a better subset. In later case until a predefined number of attributes are selected or specified number of iterations are done.

This research work develops a normalized intuitionistic fuzzy feature subset selection scheme, which starts with an empty feature subset E. Consecutively, each feature is selection in such a way that it maximized the criteria of evaluation and add the relevant feature to E. The selected feature subset is evaluated based on the minimum redundancy and maximum relevancy principle. Each features relevancy is measure using Intuitionistic Fuzzy mutual information (IFMI) between feature  $fr_i$  and the class variable  $fr_{cl}$ . The feature's redundancy is evaluated by finding IFMI between  $fr_i$  and the subset of previously selected features which are in the E list is computed. To overcome the biased nature of multivalued features this work uses Normalized Intuitionistic Fuzzy mutual information (NIFMI). The NIFMI among two features  $fr_s$  and  $fr_t$  is computed by finding the ratio between the intuitionistic fuzzy mutual information  $IFMI(fr_s; fr_t)$  of the two attributes and the minimum entropies of those two attributes  $(H(fr_s):H(fr_t))$ . Likewise, Normalized Intuitionistic Fuzzy mutual information is defined as.

$$NIFMI(fr_s, fr_t) = \frac{IFMI(fr_s, fr_t)}{\text{Min}\{IFMI(fr_s, fr_t)\}} \quad (1)$$

Let F be the list of attributes of a dataset, which is represented as  $F = \{fr_1, fr_2, fr_3, fr_4, \dots, fr_n\}$  where n denotes number of attributes. Let us defined that C and D are two Intuitionistic fuzzy sets defined on the fuzzy sets Y. The Intuitionistic fuzzy membership value of  $R_{th}$  feature for  $i^{th}$  class represented as  $\mu_{i,r}$ , degree of non-membership value is  $\vartheta_{i,r}$  and its degree of hesitation is represented as  $\pi_{i,r}$ .

The membership value  $\mu_{i,r}$  is computed as

$$\mu_{i,r} = \left( \frac{\|\bar{fr}_i - fr_r\| \sigma}{d + \epsilon} \right)^{\frac{-2}{q-1}} \quad (2)$$

The non-membership value  $\vartheta_{i,r}$  is computed as

$$\vartheta_{i,r} = \frac{1 - \mu_{i,r}}{1 + \tau \mu_{i,r}} \quad (3)$$

The indeterminacy value  $\pi_{i,r}$

$$\pi_{i,r} = 1 - \mu_{i,r} - \vartheta_{i,r} \quad (4)$$

Where q is the intuitionistic fuzzification coefficient,  $\epsilon > 0$  is used to evade distinctiveness,  $\sigma$  denotes standard deviation while performing distance calculation [18].  $\bar{fr}_i$  signifies attributes mean value which belongs the class i. The radius of data d is denoted as  $d = \max(\|\bar{fr}_i - fr_r\| \sigma)$ . The intuitionistic fuzzy entropy (IFE) of the fuzzy sets C and D is computed as follows:

$$IFE(C) = -\frac{1}{n} \sum_{fr \in F} [\mu_C(fr) \log \mu_C(fr) + \vartheta_C(fr) \log \vartheta_C(fr) - (1 - \pi_C(fr)) \log(1 - \pi_C(fr))] \quad (5)$$

$$IFE(D) = -\frac{1}{n} \sum_{fr \in F} [\mu_D(fr) \log \mu_D(fr) + \vartheta_D(fr) \log \vartheta_D(fr) - (1 - \pi_D(fr)) \log(1 - \pi_D(fr))] - \pi_D(fr) \quad (6)$$

$$IFE(C \cup D) = -\frac{1}{n} \sum_{fr \in F} [\mu_C(fr) \vee \mu_D(fr)] \log [\mu_C(fr) \vee \mu_D(fr)] + [\vartheta_C(fr) \vee \vartheta_D(fr)] \log [\vartheta_C(fr) \vee \vartheta_D(fr)] - [1 - \pi_C(fr) \vee \pi_D(fr)] \log [1 - \pi_C(fr) \vee \pi_D(fr)] - [\pi_C(fr) \vee \pi_D(fr)] \quad (7)$$

$$IFMI(C, D) = IFE(C) + IFE(D) - IFE(C, D) \quad (8)$$

Normalized Intuitionistic Fuzzy Mutual Information based feature selection with maximized relevancy and minimized redundancy of an Intuitionistic Fuzzy dataset C with the class feature cl is computed as in "equation (9)".

$$NIFMI(C, f_{cl}) = \frac{IFMI(C, f_{cl})}{\text{Min}\{IFMI(C, f_{cl})\}} \quad (9)$$

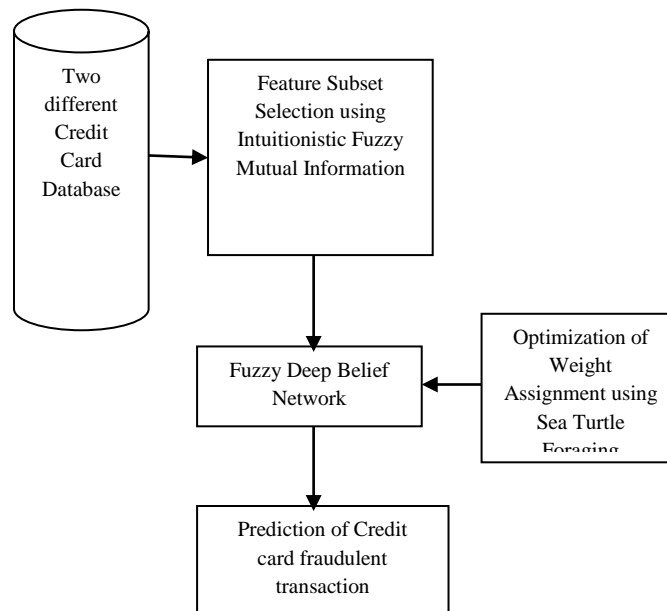


Fig. 1. Overall Workflow of the Proposed Population based Optimized Fuzzy Deep Belief Network for Credit Card Fraudulent Detection.



C. Fuzzy Deep Belief Network

Prior to assigning the input to the DBN network the domain value has to be converted to the fuzzy representation using the formula:

$$\mu_{CDS}(x) = \left\{ \begin{array}{ll} 0 & (x < a) \text{ or } (x > d) \\ \frac{x-a}{b-a} & a \leq x \leq b \\ 1 & b \leq x \leq c \\ \frac{d-x}{d-c} & c \leq x \leq d \end{array} \right\} \quad (10)$$

where  $\mu_{CDS}(x)$  refers to the membership of each instances of credit card dataset towards normal transaction as in “Fig. 2”.

A sort of Deep neural network which comprised of multiple layers of belief network known as deep belief network. In this model each layer is a Restricted Boltzmann Machine (RBM) which are stacked to each other and constructs deep belief network. DBN consist of two various types of networks they are belief network and restricted Boltzmann machine [19]. A Belief network is comprised of layers of stochastic binary units whose connections are weighted. This network is acyclic graph which permits to observe the kind of data the belief network believes. It adjusts the weights of the states between these units so that the network can produce appropriate result. The binary units in belief networks have either the state 0 or 1.

The initial process of DBN is to learn a layer of features of the visible units with contrastive divergence method. Next, to treat the activations of previously trained features as visible units and learn features of features in a second layer. At last, the entire DBN is trained when the final hidden layer finished its learning process. The greedy learning approach is used for training the DBN, because while training RBM with CD for each layer it falls under local optimum and the next stacked RBM layers takes those trained optimal values and look more local optimum. Finally, all the layers are consistently involving for local optima it gets its global optimum.

As shown in the “Fig. 3”, Restricted Boltzmann machine is a recurrent neural network which consist of binary units and undirected edges among units. The probability distributions on visible and hidden units are termed with its function of energy. The functions are formulated as follows:

$$\text{Prob}(Vs, Hd) = \frac{1}{Z} \exp(-Eg(Vs, Hd)) \quad (11)$$

$$Z = \sum_{Vs, Hd} \exp(-Eg(Vs, Hd)) \quad (12)$$

$$Eg(Vs, Hd) = -\sum_i va_i Vs_i - \sum_j hb_j Hd_j - \sum_i \sum_j Vs_i wt_{i,j} Hd_j \quad (13)$$

where Vs refers to the visible node and Hd refers to Hidden nodes, va denotes visible bias, hb signifies the hidden bias of the final layer and finally wt refers to the weight value between the previous layer and the present layer.

Like logistic regression, the conditional probabilities  $\text{Prob}(Vs_i = 1 | Hd)$  and  $\text{Prob}(Hd_j = 1 | Vs)$  and when a hidden vector  $Hd(Hd_1, \dots, Hd_j, \dots, Hd_m)$  is known, the activation probability of the ith visible unit can be computed as follows:

$$\text{Prob}(Vs_i = 1 | Hd) = \sigma(va_i + \sum_{j=1}^n wt_{i,j} Hd_j) \quad (14)$$

Similarly, Activation probability of j<sup>th</sup> hidden unit can be formulated when a set of visible vector  $vs(vs_1, \dots, vs_i, \dots, vs_n)$  is represented as:

$$\text{Prob}(Hd_j = 1 | Vs) = \sigma(hb_j + \sum_{i=1}^m wt_{i,j} Vs_i) \quad (15)$$

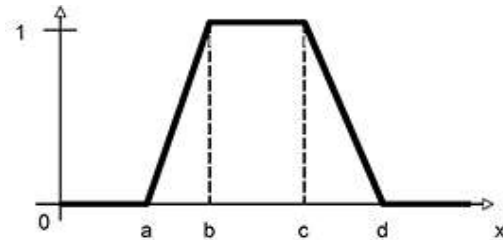


Fig. 2. Representation of Membership Function.

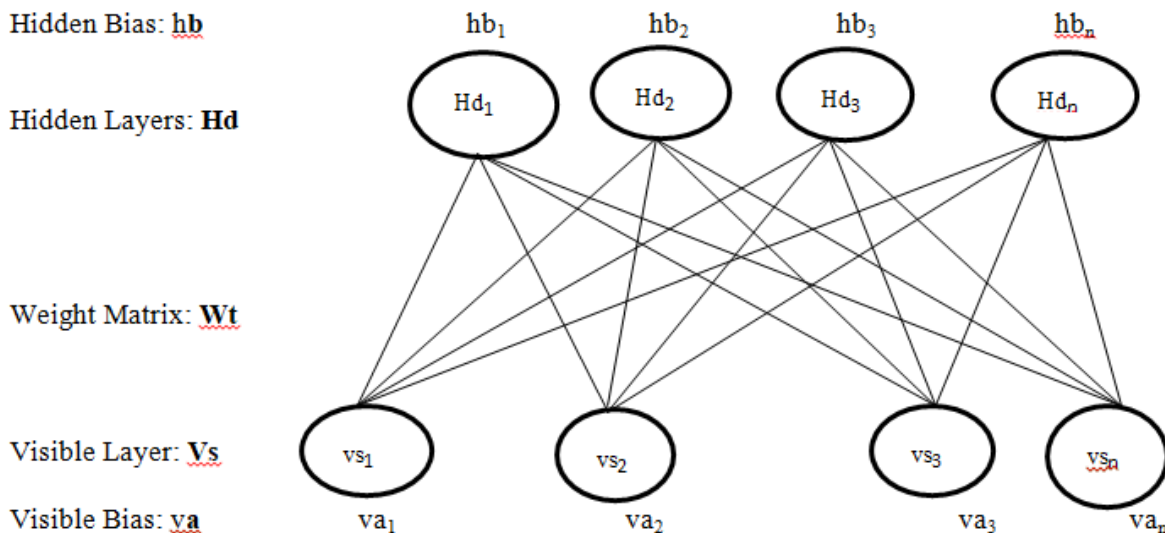


Fig. 3. Deep Belief Network RBM.

where,  $\sigma(\cdot)$  is a sigmoid function,  $w_{i,j}$  is the link weight among the  $i^{\text{th}}$  visible unit and the  $j^{\text{th}}$  hidden unit and the  $hb_j$  is the hidden bias of the  $j^{\text{th}}$  hidden unit.

To maximize the joint probability of the bunch of training inputs then it is signified as in “equation (16)”.

$$\arg \max_{wt} \prod_{vs \in Vs} Prob(Vs) \quad (16)$$

where, Vs is set of all training inputs of dataset.

#### D. Sea Turtle Foraging Algorithm

This algorithm is inspired by the sea turtle’s food searching behaviour. The sea turtle senses a kind of odor smell known as dimethyl sulfide came from their sources of food and they move towards the food source which gives out the strongest odor [20]. Ocean current also helps the turtle for their movement. The artificial foraging process of sea turtle is discussed in the subsequent steps

*Algorithm: Sea Turtle Foraging Algorithm.*

Steps involved:

1. Initialize N population of turtles.
2. Initialize the position of each turtle in a random manner

$$Pos_i(0)=[pos_1^i, pos_2^i, pos_3^i \dots pos_D^i] \quad (17)$$

where,  $i = 1$  to N number of turtles and D refers to D dimensional searching space which is contentious

3. Generate the initial velocities of turtles in a random manner  $Vel_i(0) = [vel_1^i, vel_2^i, vel_3^i \dots vel_D^i]$  and the velocity is controlled with in the predefined boundaries as follows:

$$Vel_{max} = \alpha[TUB-TLB] \quad (18)$$

$$Vel_{min} = -Vel_{max} \quad (19)$$

where, TUB and TUL are the upper bound and the lower bound of the D dimensional search space of the turtle.  $\alpha$  is the constant variable with the value ranges from 0 to 1.

4. Initial position of M food sources generated randomly

$$Fd_j(0) = [fd_1^j, fd_2^j, fd_3^j \dots, fd_D^j] \quad (20)$$

where,  $j = 1$  to M food sources which has D dimensional space for searching it.

5. Initial position of each source of food is given as the input into the objective function and estimate it to get the fitness value of that source of food.

6. Assign the position of each turtle as objective function and estimate it to get the fitness value of the concern turtle, turtle with the highest fitness value is recorded as I

$$I = \arg \max_i (fit_{pos_i(t)}) \quad (21)$$

where,  $f_{pos_i(t)}$  is the turtle  $i$ ’s fitness value at time t.

7. Velocity of each turtle is updated as shown:

$$vel_i(t) = vel_i(t-1) + \left( \frac{fit_{pos_i(t)} - fit_{pos_i(t-1)}}{fit_{pos_i(t-1)}} \right) (Pos_i(t) - Pos_i(t-1)) \quad (22)$$

where  $Pos_i(t)$  refers to the turtle position at time t and  $Pos_i(t-1)$  signifies the turtle position at time t-1.

8. Compute the ocean currents velocities by

$$VOC_i(t) = [voc_1^i, voc_2^i, voc_3^i \dots, voc_D^i] \quad (23)$$

$$VOC_i(t) = \alpha (Pos_i(t) - Pos_i(t-1))$$

9. Sum the turtle velocity to that of ocean current velocity to get the united velocity

$$UVOC = vel_i(t) + VOC_i(t) \quad (24)$$

10. If the turtle’s fitness value is less than that of the food source, then its contribution of food source (CFS) is represented as

$$CFS_j = \frac{fit_{FS_j}}{\sum_{q=1}^M fit_{FS_q}} \quad (25)$$

where,  $fit_{FS_j}$  refers to the food source  $j^{\text{th}}$  fitness value.

11. Determine the distance among the turtle and the food source using the formula

$$Dist_{ij} = ||Pos_i - FS_j|| \quad (26)$$

12. Compute the level of odor of the food source j seeming by the turtle i

$$C_{ij}(t) = (CFS_j * \exp\left(-\frac{Dist_{ij}^2}{2\sigma^2(t)}\right)) \quad (27)$$

$$\sigma(t) = \sigma_0 \exp\left(-\frac{t}{T}\right) \quad (28)$$

where,  $\sigma(t)$  indicates the level of fading of odor with short-lived time,  $\sigma_0$  is a persistent set to be equivalent to 1, and T is referring to iterations denoting the longest time the odor completely disappears.

13. Discover the turtle  $i$ ’s best food source, which has the highest value of  $C_{ij}$  compared to all other food sources available.

$$J = \arg \max_j (C_{ij}) \quad (29)$$

14. Update each turtles position as follows:

$$Pos_i(t+1) = Pos_i(t) + UVOC_i(t) + C_{ij}(t)(FS_j - Pos_i) \quad (30)$$

15. Stop the process if the maximum no of iteration is reached or else go to the step 6.

E. Proposed Algorithm: Fuzzy Deep Belief Network enriched with Sea Turtle Foraging Algorithm (EFDBN-STFA)

```

Input: Credit Card Dataset: CCDS
N : number of instances in the CCDS
Procedure
Stage 1: Feature Subset Selection
    • Apply preprocessing on CCDS
    • For I = 1 to N
        ◦ Calculate the Intuitionistic fuzzy Mutual Information by applying the equations (1) to (9)
        ◦ Generate Reduced Feature subset of CCDs to RCCDS
Stage 2: Convert the input to Fuzzy domain representation by applying equation (10)
Stage 3: Classifying the credit card transaction using EFDBN-STFA
Training Data: Tr-RCCDS; Test dataset Ts-RCCDS
Number of Layers NLY, Number of Epochs NE;
Number of hidden layers in FDBN is HDL1,...,HDLNL
Number. of units in each hidden layer HD1...HDN;
Weight Wt = { wt1, . . . , wtn}; biases Va, Hb
For G = 1 to NLY- 1
    For L = 1 to NE
        For Q = 1 to (TR + TS)
            ◦ Apply the equation (16) for supervised learning done by FDBN
            ◦ Compute the parameter values using the equations (14) and (15)
            ◦ Improve the EFDBNI hidden nodes HDL,N
            ◦ Call Sea turtle Algorithm for Weight optimization in EFDBN
            ◦ Classify the transaction based on the trained EFDBN using the equation (16)
End
End
    
```

IV. CONCEPTUAL FRAMEWORK

The proposed EFDBN-STFA credit card fraudulent transaction detection is deployed using python code. The performance analysis is done on two different credit card datasets. The other detection models used for comparison are Aggrandized Random Forest (RF), Aggrandized Kernel based Support Vector Machine (SVM) and Artificial Neural Network (ANN). The evaluation metrics used for examining the performance of the detection models are done using accuracy, precision and recall measures.

Accuracy: It is defined as the ratio of transaction which are correctly predicted to the total credit card transactions. This measure instantly specifies how well a model is trained.

$$Acc = \frac{Tot.no.of\ correctly\ detected\ geunie\ and\ fraudulent\ transactions}{total\ no\ of\ transaction\ in\ credit\ card\ dataset}$$

Precision: This metric is defined as the ratio of correctly predicted fraudulent transactions to the total transactions predicted as fraud.

$$PrCs = \frac{Tot.no.of\ correctly\ detected\ fraudulent\ transactions}{total\ no\ of\ transaction\ predicted\ as\ fraud}$$

Recall: This is signified as ratio of correctly predicted fraudulent transaction to the actual number of fraudulent transactions in the credit card dataset.

$$Rcl = \frac{Tot.no.of\ correctly\ detected\ fraudulent\ transactions}{total\ no\ of\ acutal\ fraudulent\ transactions\ in\ credit\ card\ dataset}$$

V. RESULTS AND DISCUSSIONS

By applying the normalized intuitionistic fuzzy mutual information feature methods, the five important features of each dataset is selected which is then used in further steps of the proposed system extended fuzzy deep belief network enriched with sea turtle foraging algorithm (EFDBN-STFA). This work is an extension of my previous work aggrandized random forest(RF) and aggrandized kernel based SVM(AKSVM+FPSO) which detects the credit card fraud transactions, thus leading to prediction based on the behavior patterns of the user with important features.

The “Fig. 4(a)” and “Fig. 4(b)” portraits the importance of features which involves in maximizing the relevancy among features with class and minimizes redundancy within features. By applying intuitionistic fuzzy mutual information (IFMI), the figures displays top five best features of two different credit card datasets. IFMI well treats the problem of inconsistencies in determining potential features of this credit card fraudulent detection. The features with highest score are considered for prediction process, using reduced feature subset which influences the process of credit card fraudulent detection more accurately.

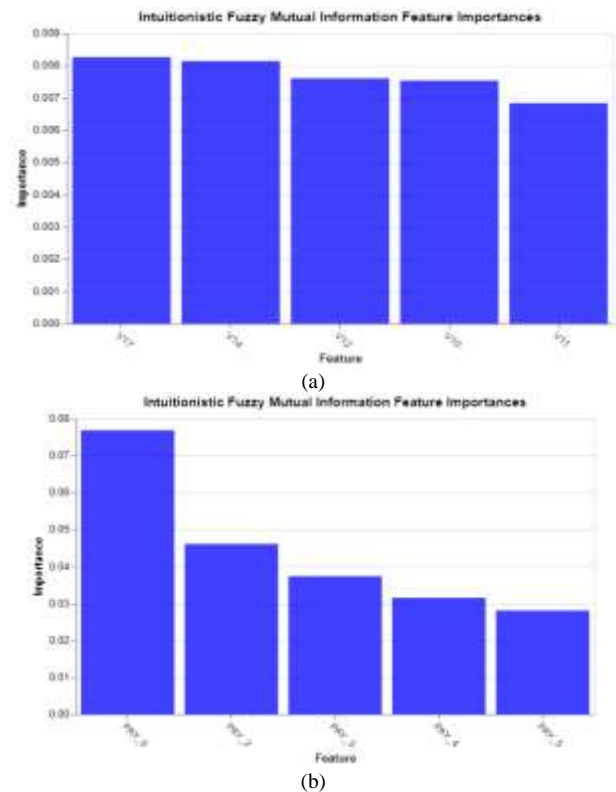


Fig. 4. (a) Feature Selection of Credit Card Dataset1, (b) Feature Selection of Credit Card Dataset2.

From the above “Fig. 5” it is observed that the performance of the proposed EFDBN-STFA produced highest accuracy rate of 96.2% for dataset1 and 95.3% for dataset 2. The other three classification models produce less accuracy because they failed to handle uncertainty in classifying credit card transactions as either genuine or fraudulent. With the knowledge of intuitionistic fuzzy mutual information is used to determine the relevancy among features with class, the redundancy is greatly reduced and the proposed EFDBN-STFA uses only the attributes with highest information score. Thus, it achieves to produce better accuracy while comparing other three models.

The “Fig. 6” illustrates the performance of the four different credit card fraud detection models based on the precision measure on two different credit card datasets. With the ability of the fine tuning their weights and biases assigned to the hidden nodes in fuzzy deep belief network, it surpasses the performance of other conventional detection models. With the enriched knowledge of sea turtle foraging algorithm this proposed model instead of assigning the weights in the random manner, they optimized the weight assignment of the deep belief network more precisely.

It is proved from the results shown in the “Fig. 7” which compares the performance of the four different credit card fraud detection models namely RF, SVM, ANN and EFDBN-STFA. The percentage of total relevant results correctly classified was done by the proposed model EFDBN-STFA. As the real-world credit card datasets cannot be handled in the crisp value to get appropriate interpretation, the fuzzy model which treats them in the form of linguistic terms using membership grade have greatly influence the process of fraud detection process. The IFMI induce the significant features which has to be used as the input is used by EFDBN for credit card fraud detection process. The merit of population-based metaheuristic searching is done by sea turtle foraging mechanism so that the weights assigned to the deep belief networks are fine-tuned based on the optimization and thus this enriched model achieves highest recall value compared to the other models.

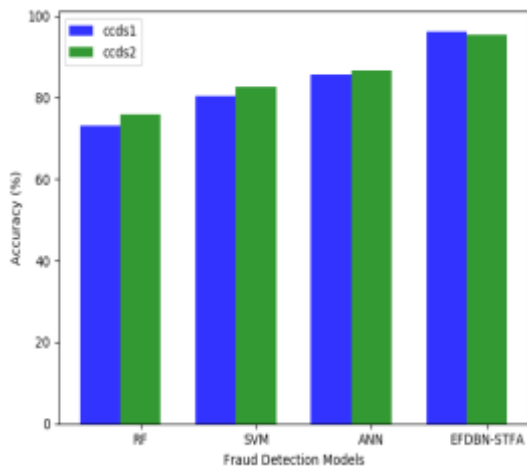


Fig. 5. Accuracy Comparison of Four different Credit Card Fraud Detection Models.

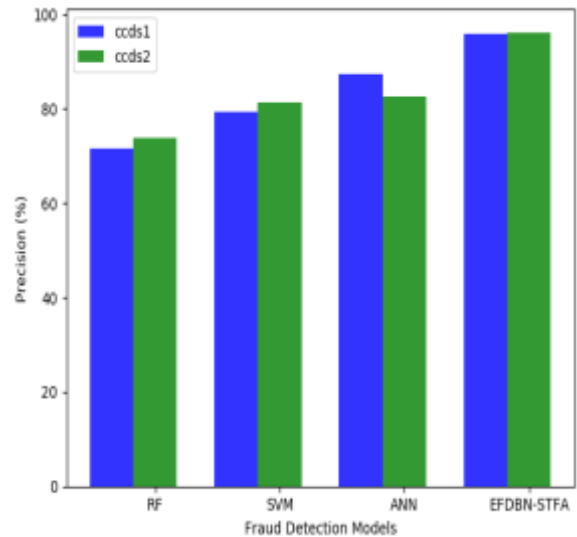


Fig. 6. Precision Measure Comparison of Four different Credit Card Fraud Detection Models.

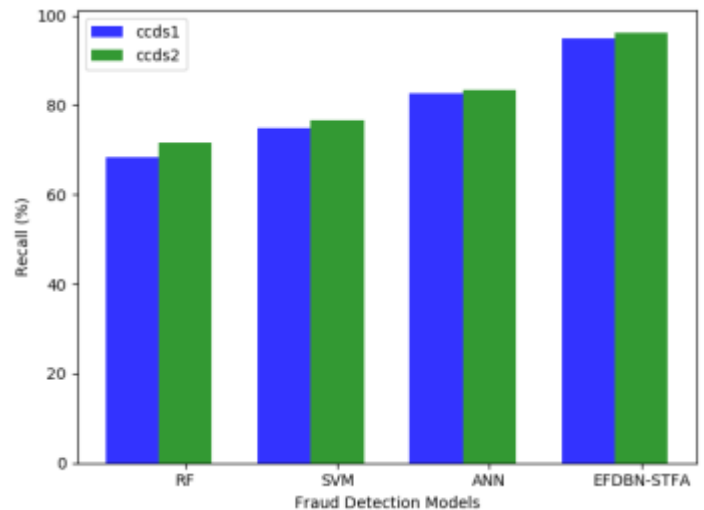


Fig. 7. Comparison of Recall Metric for Four different Credit Card Fraud Detection Models.

## VI. CONCLUSION

The ultimate motive of this proposed research work is to detect the credit card fraudulent transaction as it is continuous in nature. This work developed an enriched fraud detection model which handles the presence of vagueness and complexity in determine the pattern of transaction by focusing on three different dimensions. As a primary factor, significant feature of the datasets is selected to handle voluminous credit card datasets using the intuitionistic fuzzy mutual information about the features. With the reduced feature set, the fraud detection process is greatly influenced by fuzzy deep belief network, which gains the deep knowledge of the datasets and the relationship among the features using the stacked restricted Boltzmann machine. Instead of assigning the weights in a chaotic manner, with the inspiration of sea turtle foraging based optimization the weights assigned to the hidden layers are fine-tuned and thus the expected and the observed results produced an accepted outcome. The performance of EFDBN-

STFA is done on two different credit card datasets which are represented entirely in a different domain value. The results proved the empowerment of EFDBN-STFA for credit card fraudulent detection in presence of uncertainty by consequence higher detection rate of fraudulent transaction compared to the existing models. The proposed model restricts the frauds to happen while transactions and promote prevention of fraud in the future.

#### ACKNOWLEDGMENT

We authors thank all the contributors of this journal for considering the article. I would like to thank my guide for giving his support and encouragement in my work. Also, thank the authors of the references giving me the privilege to cite their article and enhance my knowledge. With the responsibility as Ph.D. Scholar at Nehru Arts and Science College, I thank all the teachers and my friends in giving their valuable ideas and support. Wish this article will be beneficial for future scholars.

#### REFERENCES

- [1] Wang, S., "A Comprehensive Survey of Data Mining-Based Accounting-Fraud Detection Research." International Conference on Intelligent Computation Technology and Automation, **vol.1**, pp.50-53, 2010.
- [2] Ngai, E.W.T., Hu, Y., Wong, Y. H., Chen, Y., & Sun, X, "The application of data mining techniques in financial fraud detection: A classification framework and an academic review of literature", Decision Support System 2010.
- [3] Askari.S and Hussain.A, "Credit Card Fraud Detection Using Fuzzy ID3", IEEE, Computing, Communication and Automation (ICCCA), 446-452, 2017.
- [4] Vaishnavi Nath Dornadula, Geetha S, "Credit Card Fraud Detection using Machine Learning Algorithms", International Conference On Recent Trends In Advanced Computing 2019, ICRTAC 2019 Computer Science 165 (2019) 631-641.
- [5] Imane sadgali, Nawal Sael, Faouzia Benabbou, "Fraud detection in credit card transaction using neural Networks", Association For Computing Machinery. ACM SCA2019, October 2-4.
- [6] Wen-Fang YU and Na Wang, "Research on Credit Card Fraud Detection Model Based on Distance" International Joint Conference On Artificial Intelligence 2009.
- [7] S P Maniraj, Aditya Saini, Swarna Deep Sarkar Shadab Ahmed, " Credit Card Fraud Detection using Machine Learning and Data Science", International Journal of Engineering Research & Technology (IJERT), **Vol. 8** Issue 09, September, 2019.
- [8] Andrea Dal Pozzolo, Giacomo Boracchi, Olivier Caelen, Cesare Alippi, and Gianluca Bontempi, (2017), "Credit Card Fraud Detection: a Realistic Modeling and a Novel Learning Strategy." IEEE Transactions on Neural Networks and Learning Systems · September, PP(99):1-14, 2017.
- [9] Navneet Kr. Kashyap, B.K. Pandey H.L. Mandoria, "Analysis of Pattern Identification Using Graph Database for Fraud Detection". Oriental Journal of Computer Science & Technology, August, **Vol. 9**, No. (2): Pgs. 81-91, 2016.
- [10] Salvatore J. Stolfo, David W. Fan, Wenke Lee and Andreas L., "Prodromidis Credit Card Fraud Detection Using Meta-Learning: Issues and Initial Results", In AAAI-97 Workshop on Fraud Detection and Risk Management, 1997.
- [11] V. Dheepa, R. Dhanapal, Behavior Based Credit Card Fraud Detection Using Support Vector Machines, ICTACT Journal on Soft Computing, Volume: 02, Issue: 04, 2012.
- [12] Chuang-Cheng Chiu and Chich-Yuan Tsai, "A Web Services-Based Collaborative Scheme for Credit Card Fraud Detection", Proceedings of the IEEE International Conference on e-Technology, e-Commerce and e-Service, pp. 177-181, 2004.
- [13] Tao Guo and Gui-Yang Li, "Neural Data Mining for Credit Card Fraud Detection", International conference on Machine Learning and Cybernetics, **Vol. 7**, pp. 3630-3634, 2008.
- [14] Suvasini Panigrahi, Amlan Kundu, ShamikrSural and A.K.Majumadar, "Credit Card Fraud Detection: A Fusion Approach Using Dempster Shafer Theory And Bayesian Learning", Information Fusion, **Vol. 10**, No. 4, pp. 354-363, 2009.
- [15] Kaggle Dataset1: <https://www.kaggle.com/mlg-ulb/creditcardfraud>
- [16] Default of credit cards client dataset: I-Cheng Yen, Department of Information Management, Chung Hua University, Taiwan.
- [17] Vergara, J. R. and Estevez, P. A. "A review of feature selection methods based on mutual information." Neural Computing and Applications, 24(1):175-186, 2014.
- [18] L. Visintin, R. H. S. Reiser and B. R. C. Bedregal, "Interval-Valued Intuitionistic Fuzzy Implications," 2011 Workshop-School on Theoretical Computer Science, Pelotas, RS, 2011, pp. 46-52, doi: 10.1109/WEIT.2011.22.
- [19] Feng Shuang, C. L. Philip Chen, "Fuzzy Restricted Boltzmann Machine and Deep Belief Network: A Comparison on Image Reconstruction", IEEE International Conference on Systems, Man, and Cybernetics (SMC), pp 1828-1833, 2017.
- [20] Daranat Tansui, Arit Thammano, "Sea Turtle Foraging Algorithm for Continuous Optimization Problems.", Proceedings of 2016 6th International Workshop on Computer Science and Engineering, (WCSE 2016), 17-19 , pp. 678 -681.

# Meta-Analysis of Artificial Intelligence Works in Ubiquitous Learning Environments and Technologies

Caitlin Sam<sup>1</sup>, Nalindren Naicker<sup>2</sup>, Mogiveny Rajkoomar<sup>3</sup>

Department of Information Systems  
Durban University of Technology  
Durban, KwaZulu-Natal

**Abstract**—Ubiquitous learning (u-learning) refers to anytime and anywhere learning. U-learning has progressed to be considered a conventional teaching and learning approach in schools and is adopted to continue with the school curriculum when learners cannot attend schools for face-to-face lessons. Computer Science, namely the field of Artificial Intelligence (AI) presents tools and techniques to support the growth of u-learning and provide recommendations and insights to academic practitioners and AI researchers. **Aim:** The aim of this study was to conduct a meta-analysis of Artificial Intelligence works in ubiquitous learning environments and technologies to present state from the plethora of research. **Method:** The mining of related articles was devised according to the technique of Preferred Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA). The complement of included research articles was sourced from the broadly used databases, namely, Science Direct, Springer Link, Semantic Scholar, Academia, and IEEE. **Results:** A total of 16 scientific research publications were shortlisted for this study from 330 articles identified through database searching. Using random-effects model, the estimated pooled estimate of artificial intelligence works in ubiquitous learning environments and technologies reported was 10% (95% CI: 3%, 22%;  $I^2 = 99.46\%$ ,  $P = 0.00$ ) which indicates the presence of considerable heterogeneity. **Conclusion:** It can be concluded based on the experimental results from the sub group analysis that machine learning studies [18% (95% CI: 11%, 25%),  $I^2 = 99.83\%$ ] was considerably more heterogeneous ( $I^2 = 99.83\%$ ) than intelligent decision support systems, intelligent systems and educational data mining. However, this does not mean that intelligent decision support systems, intelligent systems and educational data mining is not efficient.

**Keywords**—Educational data mining; intelligent systems; artificial intelligence; PRISMA; machine learning; ubiquitous learning

## I. INTRODUCTION

Ubiquitous learning (u-learning) is a shift from the e-learning paradigm which describes an environment that permits the use of ubiquitous computing devices to access teaching and learning contents by means of wireless networks at any time and in any location. U-learning is characterized by accessibility where the information is readily available whenever learners need its utilization. In u-learning the information remains on the platform and is always available to learners. There is immediacy where the information can be acquired instantly by the learners. The u-learning environment is interactive which allows learners to interact with teachers, peers, and experts effectively and efficiently via different

media. Furthermore, u-learning is context-aware where the environment can adjust to learners' real situations to necessitate adequate information for them [1-3]. U-learning has progressed in the recent unprecedented times of COVID-19 and climate change and its adoption is considered germane as conventional teaching and learning [4].

It is necessary to investigate scientific works in u-learning to understand the impact of scientific approaches in this evolving technology for teaching and learning that education systems around the globe have become increasingly reliant on. A systematic review and meta-analysis can put into perspective what studies have been conducted and expose gaps that exist for future studies. A systematic review provides an objective high-level overview on a research topic until present state. A meta-analysis is a statistical means of integrating diverse studies in a research area [5]. Limitations in existing work from the literature show that the impact of Artificial Intelligence and its applications in u-learning frameworks is not sufficiently addressed as a collective.

The systematic review and meta-analysis proposed by the authors have two main goals: 1) To present the evolution of artificial intelligence (AI) works in u-learning to map out the terrain for future scientific studies. 2) Highlight and analyze the specific application areas of predominate AI Algorithms in u-learning. This meta-analysis paper will provide more information to ensure that practitioners on the ground make well informed educational decisions on u-learning implementation and that scientific researchers gain valuable insights to make better research decisions for future trends.

## II. RELATED WORKS

In the study by Meyliana, Hidayanto and Budiardjo in [6], students' social media preference was analysed to increase student engagement with the university. Data was collected from 1021 students from fifty-eight Indonesian Universities using questionnaires. Entropy was used to process data and assign criteria weights for social media preference. Using the TOPSIS (Technique for Order Preference by Similarity to the Ideal Solution) method, it was established that the implementation of social media was more dependent on information quality as opposed to service quality. However, while comprehensiveness and usefulness of information was highly essential to students, they also valued system availability, efficiency, and fulfilment as it directly impacted their expectation and active learning process. TOPSIS was also used to rate the social media platforms which resulted in

LINE being the best social media platform and Podcast being the least likely platform to enhance student engagement with universities [6].

Since the advent of the primary Social Networking Site (SNS) as a novel way of communicating with other people, a lot of research has tried to identify theoretically and empirically the history of, the impact on, and the characteristics of the relationship between users and the SNS. There exists a behavioural studies research gap on the reasons why users join and participate in SNSs [7]. The study by Rad, Dahlan, Iahad, Nilashi and Ibrahim in [7] explored the influential factors causing users to adopt an SNS. A multi-criteria decision-making (MCDM) tool, fuzzy AHP (Analytical Hierarchy Process) was used to evaluate the level of importance that literature derived factors such as: performance expectancy, social influence, effort expectancy, trust, facilitating conditions, privacy, perceived enjoyment, self-efficacy, and attitude toward technology had on the adoption of an SNS. Data was collected from 291 University students in the field of SNS via questionnaires and the findings of the study were that trust, performance expectancy and security were critical influential facts in SNS adoption [7].

It is apparent that technology has changed the landscape of the learning environment and that the way in which school learners learn is enhanced by different modes of education. Classroom technology incorporates interactive learning technology such as e-book technology. In Malaysia, the acceptance of novel technology like e-book technology by school children was considered important [8]. The study by Elyazgi, Nilashi, Ibrahim, Rayhan and Elyazgi in [8], identified the interface factors of CCI (Children Computer Interaction) and the determinants of usability guiding e-book behavioural acceptance by 417 school learners. The combination of the TAM (Technology Acceptance Model) and the e-book technology-related literature review, the research hypotheses were established from the interrelationship of a detailed set of constructs. The research hypotheses built the measurement framework, which was quantified by a structured questionnaire comprising a five-point Likert scale. Using the questionnaire and TOPSIS the importance of interface factors was deuced. The analysis of data indicated positive results about perceived ease of use, perceived usefulness, learner behaviour, usability, and interface. The combination of CCI and TAM factors afforded results that showed that school learners accepted the use of e-books. The highest ranking was awarded to perceived ease of use whilst the lowest ranking was behaviour intention. The former was attributed to the functions and features of e-books which seemed to be easy to use. However, it was concerning that the e-book technology usability scale was lower than the interface scale, which inferred that school learners' e-book technology acceptance will improve if it is viewed as championing an elevated level of interactivity [8].

The paper by Omorogbe and Igbinosun in [9] examined the attributes that parents considered for school choice enrolment of their children or wards. Twelve attributes from four categories of school alternatives were studied in this work. In Benin City, a survey was indiscriminately conducted in the three local government areas. The AHP was

implemented to evaluate the attributes and intuitionist fuzzy TOPSIS was applied to rank the alternatives. The correctness and consistency of results were affirmed by the two metric functions that were used both producing the same result. The ranking of the schools was as follows in descending order: The Missionary schools (A4), private schools for middle class (A2), the premier private schools for the elite (A3), and the Public (government) schools (A1). It was concluded that by implementing a scientific approach to a humanistic system, appropriate and accurate results can be produced.

The study by Pires and Cota in [10], focused on developing a self-regulating adaptive intelligent system to enhance special education needs, extending the paradigm to special education needs, using non-linear methods to correct the learning path to cater for individual special needs, and extending the research paradigm to u-learning. The study was founded on inclusive education and multiple intelligences such that special needs learners were able to fill cognitive gaps in their learning process. The proposed architecture of the system was based on an intelligent structure supplemented by a Genetic Algorithm (GA) and a Chi-square statistical function. The study involved a sample population of 13 600 kindergarten learners of mixed abilities who were subject to an intelligent system comprising a Java developed GA module, and a XML & SCORM based module consisting of an LMS (learning management system) such as Moodle and a Knowledge Block repository [10].

The study by Angeli, Howard, Ma, Yang, and Kirschner [11] explained and addressed several key questions on the utilisation of data mining in educational technology classroom research. Previous studies conducted in Australia and Europe which used the data mining techniques fuzzy representations and association rules mining were presented as examples in this study. These studies investigated 115 university student and 12 978 school learners' behaviours, experiences, and learning within computer-assisted classroom activities. In the study employing fuzzy representations, questionnaire data was inductively explored. This study aptly depicted how data mining could be used by educational technologists to monitor and guide the integration of school-based technology. The inferences of the study were reviewed based on the need to create educational data mining tools that could present information, results, comments, explanations, and recommendations in profound ways to novice users in data mining such as teachers. The study using association rules mining involved comprehension clarity on how learners with cognitive differences interact with a simulation in problem solving. The study used Statistica as a data mining tool and illustrated how data mining could be employed to enhance evaluation practices of educational software in the educational technology field. Finally, matters associated with data privacy were addressed in both studies [11].

A high accuracy of students' performance prediction is useful in identifying the low performing students at the start of the learning process. This objective is achieved by machine learning where techniques are employed to uncover patterns or models of data which is valuable in decision-making [12]. The study by Belachew and Gobena in [12] applied machine learning concepts to the dataset obtained from the college of

computing and informatics of Wolkite University registries office. The study collected 1071 student's transcript data that consisted of their grades in all courses and their final GPA (Grade Point Average). Machine learning methods, Naïve Bayesian, Support Vector Machine (SMO) and Neural Networks were applied after pre-processing the data. Lastly, a model for each method was developed, and the performance and results of each model was evaluated. The aim of the study was to create a model using machine learning to derive conclusions on students' academic performance [12].

Information and communications technology (ICT) have become very important in all spheres of human life. They are utilized in various fields as information systems (IS) using numerous telecommunication media to afford end-users the capability to control digital data [13]. Additionally, the development of relatively new technology, propagated distance learning through e-Learning platforms in the last two decades. Recommendation systems have become progressively used in IS, more so in e-learning platforms. These systems operate to recommend and propose content of these platforms to end-users corresponding to their needs to allot the most information for learning [13].

The paper by El Mabrouk, Gaou and Rtili in [13], presented a data mining based intelligent hybrid recommendation system using Neural networks, Bayesian networks and Decision trees with a population size of 700 university students. This system comprised four parts. The first part for data collection and centre of interest creation was done via two modes, namely: explicit data collection which was based on end-users and what they entered on their profiles, and automatic and implicit data collection by offering a survey to users to gather information about their interest. The second part involved managing the information previously collected and developing the learning model, to categorise users who posted the content and catalogue content to forward the results to the recommendation module. The third part drew parallels between content and learners and made recommendations for learners. The fourth part created a recommendation by learner log file which was used in the upcoming recommendation. The results of the study proposition were satisfactory, and the system was optimised with regards to response, processing time and accuracy when compared to traditional recommendation [13].

To effectively execute critical pedagogical interventions for students' satisfactory and on-time graduation, students' future performance based on their academic records must be predicted accurately. Student performance in completing degrees' prediction has limited studies conducted which results in several new challenges such as: the disparity of student's differences regarding selected courses and backgrounds; courses are not uniformly informative in making accurate predictions; and the incorporation of students' evolving progress needs into the prediction [14]. The study by Xu, Moon and van der Schaar [14] developed a new machine learning method for student performance prediction in degree courses that could address the aforementioned key challenges. There were two major features of the proposed method: the first feature was a bilayered structure consisting of various base predictors and a flow of ensemble predictors to make

predictions on students' progressing performance; the second feature was a data-driven approach built on probabilistic matrix factorization and latent factor models to unearth course relevance, which is imperative for creating efficient base predictors. With extensive simulations on a dataset of 1169 university students collected over three years, the study revealed that the proposed method achieved greater performance compared to benchmark approaches [14].

The study by Khoshi, Gooshki and Mahmoudi in [15] used fuzzy TOPSIS and AHP to prioritise the effective qualifications of medical lecturers from the opinion of medical students from the allied Tehran University of Medical Sciences in 2013 to 2014. Two hundred medical students were chosen based on random sampling method and were surveyed in accordance with Cochran's formula. Research based questionnaires were used as data collection tools that were divided professional, technical, and individual parts. Experts approved content validity. By calculating the Cronbach's alpha ( $\alpha = 0,85$ ), reliability was confirmed by measuring the internal cohesion degree [15].

A commonly employed approach for student modelling is Bayesian Knowledge Tracing (BKT). It is a commonly used approach for student modelling. A versatile model that can be used for various tasks is Long Short-Term Memory (LSTM) [16]. The study by Mao, Lin and Chi in [16] compared BKT, a derivative IBKT (Intervention-BKT) and LSTM based on two student modelling tasks; learning gains and post-test scores prediction. An automatic skill discovery method (SK) was incorporated into all three models which augmented the exercise-skill assignments with a nonparametric prior. A total of six models was explored: IBKT, IBKT+SK, BKT, BKT+SK, LSTM, and LSTM+SK. One training dataset was collected from Cordillera which is a language physics intelligent tutoring system. The other training dataset was collected from Pyrenees which is a standard probability intelligent tutoring system. The results of the study revealed that BKT and BKT+SK outdid the other models on post-test scores prediction. Additionally, LSTM and LSTM+SK attained the highest area under the ROC curve (AUC), F1-measure, and accuracy on learning gains prediction. Moreover, the study demonstrated that the BKT+SK combination could dependably predict post-test scores only using the first 50% of the total training sequences. For early prediction of learning gain, making use of the first 70% of the total training sequences, LSTM delivered an equivalent prediction using all of the training sequences. The findings of the study revealed a learning environment that could predict learning gains and students' performance early and could afford an adaptive pedagogical strategy appropriately [16].

The study by Mohamed and Lamia in [17], involved the use of computers, tablets, smartphones, and other smart devices as auxiliary tools of contemporary teaching and learning methods. The flipped classroom approach was used as an element of IOT (Internet of Things) to encourage problem solving in a mathematical logic course both in and out of the classroom via an ITS (Intelligent Tutoring System). The sample population was 50 university students and the researchers employed on-way ANOVA to determine the effectiveness of the flipped classrooms. This study showed



that self-efficacy, perceived usefulness, perceived support, and compatibility for improving social ties are imperative precursors for continued use of flipped classrooms [17].

The study by Sirait, Fitriani, Hidayanto, Purwandari and Kosandi in [18] determined eleven criteria of social media platform selection to increase student participation in government activities and eight alternatives that best aligned to the criteria. Sirait *et al.* (2018) also established an order for assessing social media preferences and highlighted the social media that were popular to assist e-participation. Criteria generated was based on the theory of hedonic and utilitarian motivation. Data was collected from University students in Indonesia through a questionnaire and was processed via fuzzy AHP to ascertain the weight of social media to enhance e-participation in government activities. The TOPSIS method was employed to establish social media preferences. The results revealed that system quality of high importance along with hedonic gratification and information quality. The best SM application for student e-participation in priority order was: Line, Instagram, Path, YouTube, Facebook, Twitter, Blog and Wiki's [18].

Adaptive e-learning platforms offer personalized learning process mainly depending on learning styles. The conventional way of finding learning styles relies on requesting learners to evaluate their own behaviours and attitudes through questionnaires and surveys. This method renders several weaknesses such as the lack of learners' self-awareness of their own inclinations. Additionally, most learners feel bored when asked to complete a questionnaire and the conventional way presumes that learning styles cannot alter over time [19]. The paper by El Aissaoui, Oughdir and El Alloui in [19] proposed a generic approach for automatically identifying learning styles corresponding to a specific learning styles model. The study comprised two major steps: firstly, using web mining techniques to extract learning sequences from learners' log files on the Moodle platform; and secondly, using clustering algorithms to classify the extracted sequences via a specific learning style model. This experimental study was performed using a classification Learning Style Model called Felder-Silverman Model (FSM) and a clustering algorithm called Fuzzy C-Means and used a real-world dataset of 1235 university students. The goodness of the study's approach was compared with the K-means algorithm, MCQ method and the FCM algorithm. The results revealed that this two-step approach proved to be promising and outperformed traditional approaches [19].

Grouping learners specifically in high school is an essential process to classify and divide them into classes based on their interests and abilities as it helps them to thrive. Most schools apply academic grades to group learners but there are other approaches that exist. As this is an annual task with new learners, both the teachers and learners feel overwhelmed with grouping [20]. A solution to this repetitive task may be the implementation of a decision support system that can

automate the grouping process. An example of an unsupervised learning algorithm is a SOM (self-organising map) which uses an artificial neural network structure to afford a reduced dimensional representation of the given input. SOM is also a clustering technique [20]. The study by Purbasari, Puspaningrum and Putra [20] used SOM to academically group 275 school learners based on their national examination results and rapport books into three distinct clusters, namely, Social Sciences, Life Sciences and Linguistics.

It is a valuable endeavour for any educational institution to predict students' academic performance as the predictions can assist educators in supporting students who are at risk of failure. The process of educational data mining (EDM) is a machine learning and data mining technique which above other tasks can predict students' performance [21]. The study by El Aissaoui, El Alami El Madani, Oughdir, Dakkak and El Alloui in [21] proposed a methodology to create a student performance prediction model using multiple linear regression (MLR) which is a supervised machine learning technique. The three major steps of the methodology were: 1) pre-processing and analysing the variables/attributes of students using a group of statistical analysis methods; 2) using different methods for selecting the most crucial variables; 3) constructing diverse MLR models centred on the variables selected and using the k-fold cross-validation technique for comparing the 395 students' academic performance. The results obtained revealed that the model derived from the selected variables of the MARS (Multivariate Adaptive Regression Splines) method, outperformed the other constructed models [21].

AI methods, namely, Intelligent Decision Support Systems, Intelligent Systems, Educational Data Mining (EDM) and Machine Learning are captured in Table I below. Table I is categorised according to the following fields: Number (#), Authors, Year, Artificial Intelligence Methods and Techniques, Problem Focus, Number of Algorithms, Number of Datasets, Number of Variables and Size. The table captures studies in Artificial Intelligence in u-learning from 2015 to 2020. The problem focus column suggests that there were lots of studies conducted in higher education and fewer studies in school based education.

Table I also shows the techniques employed by each Artificial Intelligence method. The study numbers can be used to reference the study in the table. Intelligent Decision Support Systems techniques used were: Entropy (#1), TOPSIS (#1, #3, #13), AHP (#2, #4, #10, #13) and Fuzzy TOPSIS (#10). Intelligent Systems techniques used were: Genetic Algorithms (#5), ANOVA (#12) and Self Organising Maps (#15). Education Data Mining (EDM) techniques used were: Data Mining (#6). Lastly, Machine Learning techniques used were: Neural Networks (#7, #8), Naïve Bayes (#7, #8, #11, #14), Support Vector Machines (#7), Ensemble methods (#9), Deep Learning (#11), Clustering (#14) and K-fold cross validation (#16).

TABLE I. ARTIFICIAL INTELLIGENCE WORKS APPLIED IN U-LEARNING (2015-2020)

| #   | Authors  | Year | Artificial Intelligence Methods and Techniques  | Problem Focus   | #Algorithms | #Dataset/s | # Variable/s | Size   |
|-----|--|------|---|---|-------------|------------|--------------|--|
| 1.  | Meyliana, M., A.N. Hidayanto, and E.K. Budiardjo                   | 2015 | <b>Intelligent Decision Support Systems</b><br>Entropy and TOPSIS Method<br>10 alternatives   | Social Media preference for student engagement improvement in universities.   | 2           | 1          | 7            | 1021 University students                           |
| 2.  | Rad, M.S., Dahlan, H.M, Iahad, N.A, Nilashi, M. and Ibrahim, O.    | 2015 | <b>Intelligent Decision Support Systems</b><br>Fuzzy AHP  | Using a Multi-Criteria Decision-Making Approach for Assessing the Factors Affecting Social Network Sites Intention to Use.                    | 1           | 1          | 10           | 291 University Students                            |
| 3.  | Elyazgi, M.G., Nilashi, M., Ibrahim, O., Rayhan, A. and Elyazgi, S | 2016 | <b>Intelligent Decision Support Systems</b><br>Technique for Order Preference by Similarity to Ideal Solution (TOPSIS)<br>6 criteria/34 indexes                         | Evaluating the Factors Influencing E-book Technology Acceptance among School Children Using TOPSIS Technique.                                 | 1           | 1          | 34           | 417 School Learners                                |
| 4.  | Omorogbe, D. E. A. and Igbinosun, L. I.                            | 2016 | <b>Intelligent Decision Support Systems</b><br>integrated AHP-intuitionistic fuzzy TOPSIS<br>4 alternatives   | Parents Preference for Students' Choice of Urban Schools in Benin City, Nigeria: Integrated AHP Intuitionistic Fuzzy TOPSIS.                  | 1           | 1          | 12           | 144 Parents  |
| 5.  | Pires, J.M. and Cota, M.P.   | 2016 | <b>Intelligent Systems</b><br>Genetic Algorithm<br>The Chi-Square function  | "Intelligent" Adaptive Learning Objects applied to Special Education needs<br>Extending the eLearning paradigm to the u-Learning environment. | 2           | 1          | 4            | 13 600 Kindergarten Learners                       |
| 6.  | Angeli, C., Howard, S.K. Ma, J., Yang, J. and Kirschner, P.A.      | 2017 | <b>Educational Data Mining (EDM)</b><br>Association Rules<br>Fuzzy Representations  | Data mining in educational technology classroom research: Can it make a contribution?   | 2           | 2          | 5<br>6       | 115 University students<br>12978 School Learners   |
| 7.  | Belachew, E.B. and Gobena, F.A.                                    | 2017 | <b>Machine Learning</b><br>Neural Networks, Naive Bayesian and Support Vector Machine   | Student Performance Prediction Model using Machine Learning Approach: The Case of Wolkite University.   | 3           | 1          | 6            | 1071 University students                           |
| 8.  | El Mabrouk, M., Gaou, S. and Rtili, M.K.                           | 2017 | <b>Machine Learning</b><br>Neural networks, Bayesian networks and Decision trees.<br>1000 evaluations   | Towards an Intelligent Hybrid Recommendation System for E-Learning Platforms Using Data Mining.   | 1           | 1          | 1000         | 700 University Students                            |
| 9.  | Xu, J., Moon, K.H. and van der Schaar, M.                          | 2017 | <b>Machine Learning</b><br>A prediction layer and an ensemble prediction layer  | Predicting students' future performance based on their ongoing academic records.  | 2           | 1          | 5            | 1169 University students                           |
| 10. | Khoshi, A., Gooshki, H.S. and Mahmoudi, N.                         | 2018 | <b>Intelligent Decision Support Systems</b><br>Combined Fuzzy AHP and Fuzzy TOPSIS  | Hybrid Fuzzy AHP and Fuzzy TOPSIS to prioritise the qualifications of teachers. 6 criteria, 17 indexes.                                       | 2           | 1          | 6            | 200 University Students                            |
| 11. | Mao, Y. Lin, C. and Chi, M.  | 2018 | <b>Machine Learning</b><br>Long Short-Term Memory (LSTM)<br>Bayesian Knowledge Tracing (BKT)<br>Intervention- (IBKT) in combination with a Skills discovery method (SK) | Bayesian Knowledge Tracing (BKT) for predicting students' performance- Learning gains prediction and Test score prediction.                   | 6           | 2          | 33<br>11     | 169 University students<br>475 University students |

|     |  |      |  |   |   |   |    |                          |
|-----|--|------|--|---|---|---|----|--------------------------|
| 12. | Mohamed, H. and Lamia, M.  | 2018 | <b>Intelligent Systems ANOVA</b>   | Flipped classroom as an element of Internet of Things (IoT) into learning process.    | 1 | 1 | 9  | 50 University Students   |
| 13. | Sirait, A.D.S., Fitriani, W.R., Hidayanto, A.N., Purwandari, B. and Kosandi, M.    | 2018 | <b>Intelligent Decision Support Systems</b><br>Fuzzy AHP<br>TOPSIS   | Social media to support the development of e-government. 8 alternatives.              | 2 | 1 | 11 | 401 University students  |
| 14. | El Aissaoui, O., Oughdir, L. and El Alliou, Y.                                     | 2019 | <b>Machine Learning</b><br>Unsupervised algorithm (K-means) for grouping supervised algorithm (Naive Bayes) for classification                         | Adaptive e-learning system for detecting learning styles automatically.               | 2 | 1 | 17 | 1235 University Students |
| 15. | Purbasari, I.Y., Pusaningrum, E.Y. and Putra, A.B.S.                               | 2019 | <b>Intelligent Systems</b><br>Self-Organizing Map (SOM) algorithm  | An intelligent system which can automatically perform grouping on a list of students. | 1 | 1 | 9  | 275 School Learners      |
| 16. | El Aissaoui, O., El Alami El Madani, Y., Oughdir, L., Dakkak, A. and El Alliou, Y. | 2020 | <b>Educational Data Mining (EDM)</b><br>Data mining and machine learning techniques multiple linear regression (MLR) k-fold cross-validation technique | A Multiple Linear Regression-Based Approach to Predict Student Performance.           | 7 | 1 | 32 | 395 School Learners      |

### III. MATERIALS AND METHODS

#### A. PRISMA

Preferred Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA) is an evidence based least group of items for writing a report in meta-analyses and systematic reviews. PRISMA concentrates on giving an account of reviews and is employed as the foundation for reporting systematic reviews, specifically intervention evaluation [22]. Meta-analyses and systematic reviews are methodologically rigorous studies that establish the reference standard for developing evidence in various clinical, business fields and scientific studies [23]. The number of studies makes it essential to ensure reporting transparency and quality are strictly followed. It is important to note study quality and reporting quality are not congruent, a feebly reported study is of inadequate value as it is difficult to make a transparent and complete judgment of its value in the absence of all the required information [24].

A study that has been conducted well may be reported poorly and perhaps eclipsed by its reporting [24]. It is a proven fact that digression from the reporting guidelines in meta-analyses and systematic reviews can result in bias [25]. Such work has consequently led to the development of the PRISMA statement which is a twenty-seven-item checklist ensuring transparency in the reporting of a review [25]. Given the upsurge of systematic reviews in e-learning, m-learning, and u-learning it is crucial to ensure reviews are effectively reported. The current study is appositely reported using the PRISMA methodology to ensure education researchers, scientific method researchers, policy makers and instructors can make transparent and complete judgements to steer key scientific or education decisions.

In this meta-analysis and systematic review, the following databases were searched for and afforded the most applicable published articles: ScienceDirect, Semantic Scholar, SpringerLink, IEEE, Google Scholar, and Academia. Other

databases affording fewer applicable articles were Educase, Hindawi, Science Publishing Group, Taylor & Francis, MDPI, Modestum, Scopus, Eric, SAGE Journals, Emerald, and World Scientific. These databases were restricted to English papers published between 2015 to 2020 using the following combination of terms: Intelligent Decision Support Systems, Intelligent Systems, Machine Learning, Data Mining, Educational Data Mining, Ubiquitous Learning, u-learning, e-learning and m-learning. The search terms were separated or combined using Boolean operators such as ‘OR’, ‘AND’ and ‘NOT’. The studies identified by the search strategies were downloaded. A total of 330 published articles between the years 2015 and 2020 were identified as depicted in Fig. 1. The abstract and introduction sections were read in order to determine the eligibility the papers dealing with scientific methods used in u-learning research. As a corresponding process, reference lists of the appropriate studies were manually checked for any citations omitted by the electronic database searching. The scope was reduced to 16 scientific works which was analyzed in-depth in the meta-analysis.

#### B. Inclusion Criteria

As the focus of the study was on scientific methods used in u-learning studies, the criteria for inclusion of the articles were: Intelligent Systems (IS), Intelligent Decision Support Systems (IDSS), Machine Learning (ML); and Educational Data Mining (EDM) which were applied to u-learning studies. Since u-learning comprises e-learning and m-learning, these factors were also added to the inclusion criteria [26]. Additionally, the inclusion criteria were specific to scientific methods used in u-learning in school-based and university education.

#### C. Intelligent Systems (IS)

IS are technically innovative and autonomous systems that sense and react to the physical and social world to achieve human users’ goals. A study of how computers can comprehend and translate video sequences and static images into visual information emerged between the 1950’s to 1960’s

and has since evolved into a potent technology that is pivotal to most sectors such as education [10]. The fundamental factors that have attributed to this evolution are the exponential augmentation of the algorithms, memory capacity, and processor speed of technology. Research in intelligent systems faces numerous challenges, many of which relate to representing a dynamic physical world computationally [20].

#### D. Intelligent Decision Support Systems (IDSS)

IDSS is a decision support system that extensively uses AI techniques. AI techniques have been employed over a longstanding history as IS and knowledge-based systems in management IS [9]. Preferably, an IDSS should perform as a human consultant in helping decision makers to collect and examine evidence, detect and make a diagnosis on problems, and propose and evaluate solutions. The purpose of AI techniques rooted in an IDSS is to ensure that the previously mentioned tasks are completed by a computer under well-known decision parameters, while closely mimicking expert human capabilities [9]. In research, AI directed towards allowing systems to react to innovation and indecision in more adaptable ways is beginning to be used in IDSS. For instance, intelligent agents that conduct difficult cognitive tasks without human involvement have been utilized in a variety of decision support applications [15]. Competences of these intelligent agents consist of machine learning, knowledge sharing, automated inference, and machine learning. An array of AI techniques like fuzzy logic, rough sets, and case-based reasoning have also been employed to facilitate better performance of decision support systems in uncertain conditions [15].

#### E. Machine Learning (ML)

ML involves computer algorithm studies that develop instinctively through experience [14]. It is a subgroup of AI. ML algorithms construct a mathematical model centred on training data (sample data), to make decisions or predictions without being overtly programmed to do so [16]. ML algorithms are employed in a vast array of applications, for instance with computer vision and email filtering, where it is complex and not feasible to create traditional algorithms to execute the desired tasks [19]. ML is closely associated to computational statistics, which involves making predictions with computers. The mathematical optimization studies afford application, theory, and technique domains to the machine learning field. Data mining is a connected field of study [19].

#### F. Educational Data Mining (EDM)

EDM is a developing discipline, involved with emerging methods for investigating the rare and progressively large-scale data that are derived from educational settings. It also involves applying those methods to better comprehend learners, and the settings which they learn [11]. Irrespective of whether educational data is mined from; learners' use of

administrative data from universities and schools, computer supported collaborative learning, or interactive learning environments; it frequently has numerous levels of evocative hierarchy, which need to be decided by the data properties, rather than in advance. Other impacting factors in the study of EDM are context, time, and sequence [21].

#### G. Exclusion Criteria

Articles that were excluded were: written in dialects other than English; published before January 2015; study designs such as letters to editors, reviews, commentaries, editorials, book chapters, books, expert opinions, books, theses, and brief reports; and scientific applications to enhance u-learning. Also, articles that neglected to account for the inclusion criteria were excluded. All non-scientific publications on u-learning were excluded. Fig. 1 shows the flow diagram of the study based on the PRISMA methodology.

#### H. Quality Assessment

Information that aligned with the inclusion criteria were rooted out from the chosen studies. The study with 100% correlation to the inclusion criteria was applicable for the meta-analysis and systematic review. The first author extracted the following information from the studies which met the inclusion criteria: the author's name, year of publication, scientific model, size, and algorithms. The first and second authors assessed the qualities of each article included by using a critical appraisal tool for use in the systematic review for a study examining artificial intelligence works in ubiquitous learning environments and technologies.

#### I. Statistical Analysis

Data were extracted in a Microsoft excel spreadsheet, and analysis was carried out using statistical software. Heterogeneity among reported prevalence was assessed by computing p values of Higgins's  $I^2$  statistics;  $I^2$  was considered as significant at a p-value < 0.05. The DerSimonian and Laird's random-effects meta-analysis model [28] was used to determine the pooled effect size since the actual effect is not the same in all studies. We deal heterogeneity with subgroup analysis, meta-regression, and sensitivity analysis. Subgroup analysis was done based on study scientific model approaches. Besides, an effort to understand the sources of heterogeneity, univariate meta-regression analysis was conducted for sample size, and publication year. A forest plot [29] was used to describe pooled prevalence with 95% confidence intervals. The size of each box indicated the weight of the study, while each crossed line refers to a 95% confidence interval with the mean effect at the centre. The possibility of publication bias was assessed visually with funnel plots, and the objectivity test of Egger's test with a p-value less than 0.05 was considered evidence of publication bias.

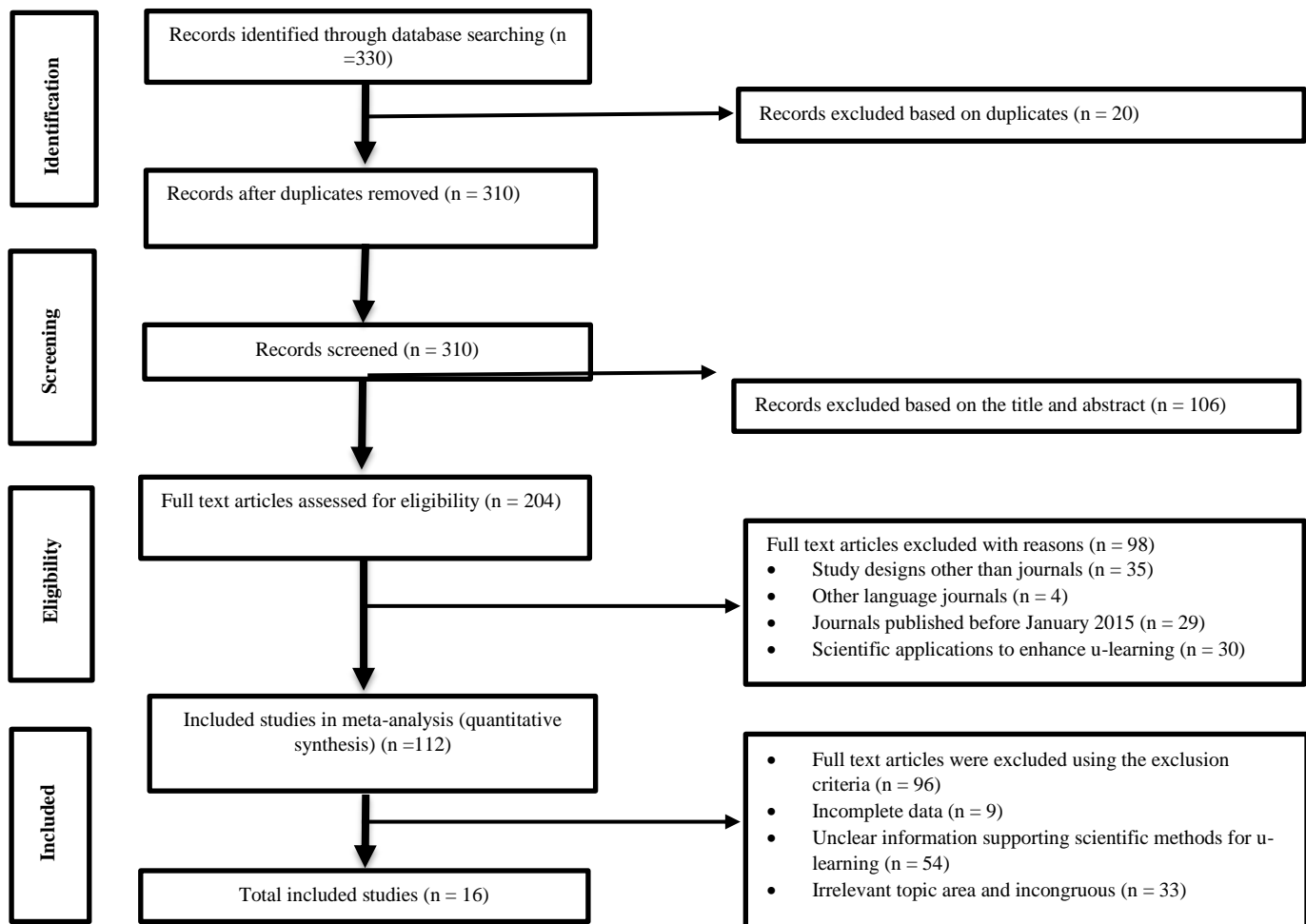


Fig. 1. Flow Diagram of the Database Searches (Preferred Reporting Items for Systematic Reviews and Meta-Analysis) Source: Moher, D., Liberati, A., Tetzlaff, J., Altman, D.G. and Group in [27].

#### IV. RESULTS

The pooled estimate showed that there is statistically significant heterogeneity between the studies. Hence, this necessitated the use of random-effects model. Thus, using random-effects model [28], the estimated pooled estimate of artificial intelligence works in ubiquitous learning environments and technologies reported by the 16 studies was 10% (95% CI: 3%, 22%;  $I^2 = 99.46%$ ,  $P = 0.00$ ) which indicates the presence of heterogeneity. The pooled estimate of artificial intelligence works in ubiquitous learning environments and technologies was presented using a forest plot (Fig. 2). From the forest plot in Fig. 2, the black dot at the centre of the grey box indicates the estimated prevalence point of each study, and the line indicates the 95% confidence interval of the estimates. The grey box indicates the weight of each study, contributing to the overall pooled prevalence estimate. The blue diamond represents the 95% confidence interval of the pooled rotavirus prevalence estimate.

##### A. Subgroup Analysis based on the Scientific Model Approaches

Subgroup analyses (Fig. 3) were carried out stratified scientific model approaches. The subgroup analysis by scientific model approaches was conducted to assess the potential heterogeneity between studies. Of the 16 studies, the highest pooled estimate was found in studies conducted with machine learning [18% (95% CI: 11%, 25%),  $I^2 = 99.83%$ ] followed by studies conducted with intelligent systems [13% (95% CI: 1%, 26%),  $I^2 = 0.0%$ ] while both intelligent decision support systems and educational data mining have the lowest percentage given as [7% (95% CI: 3%, 11%),  $I^2 = 95.21%$ ] and [7% (95% CI: 5%, 9%),  $I^2 = 0.0%$ ] respectively (Fig. 3). Meanwhile, the high heterogeneity could be as a result of significant variation in the number of students among various studies.

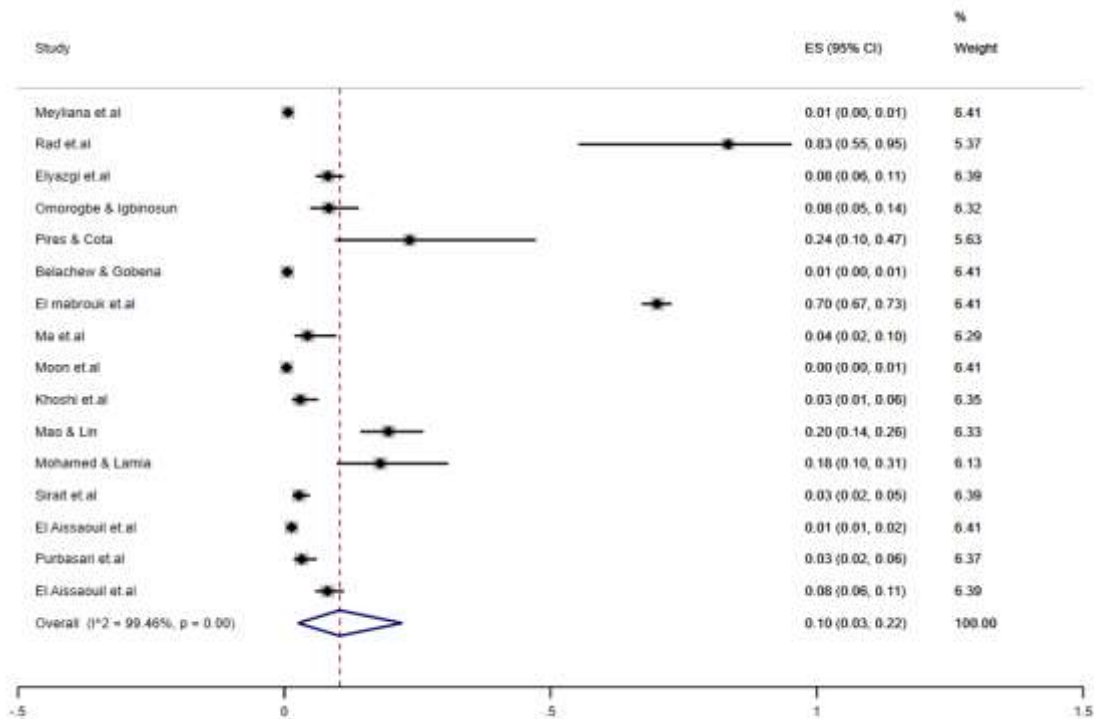


Fig. 2. Forest Plot Showing the Overall Pooled Estimate of Artificial Intelligence Works in ubiquitous Learning Environments and Technologies.

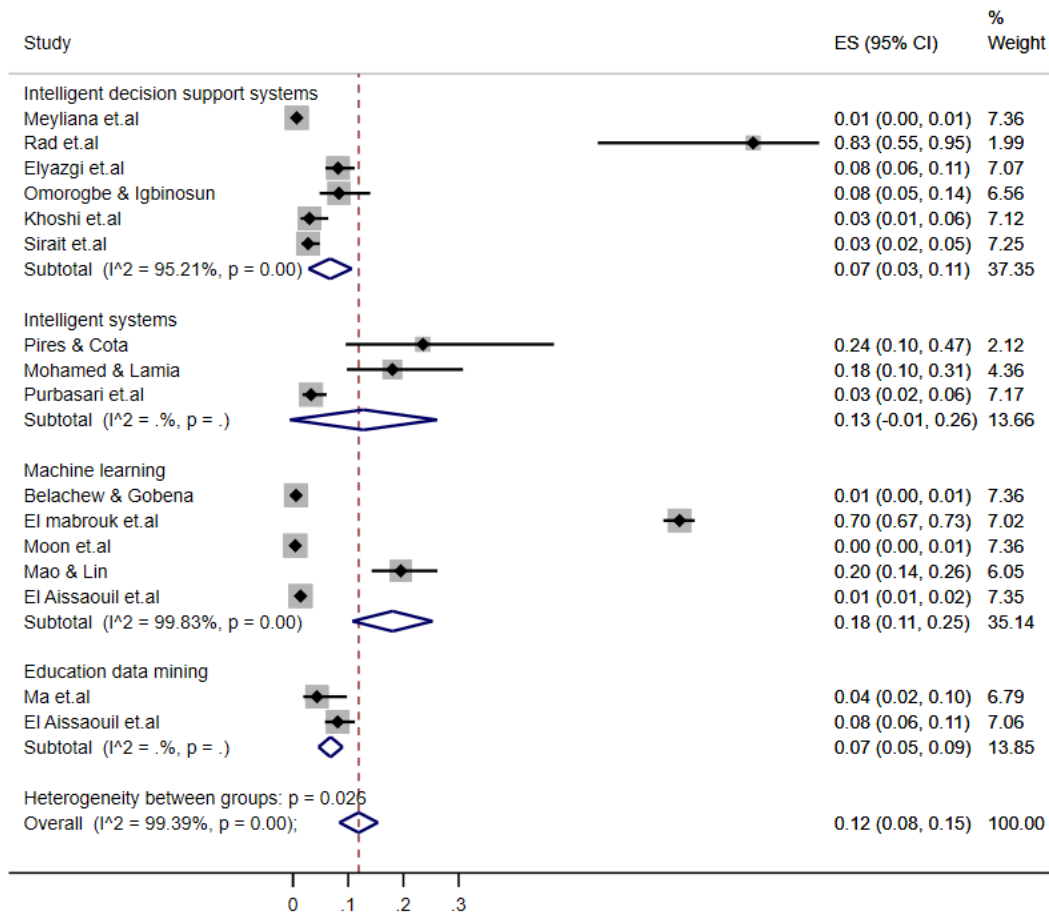


Fig. 3. Subgroup Analysis of Artificial Intelligence Works in ubiquitous Learning Environments and Technologies Stratified by Scientific Model Approaches.

## V. CONCLUSION

This study analysed different artificial intelligence works in ubiquitous learning environments and technologies based on the four major scientific approaches – intelligent decision support systems, intelligent systems, machine learning and educational data mining. The meta-analysis experimental result divulged that studies which used the machine learning approach gave the best result. Machine Learning studies have garnered significant interest and is variable in intervention effects being evaluated. Machine learning studies show heterogeneity in terms of sample size and variation in the choice of algorithms. Nevertheless, it does not mean that intelligent decision support systems, intelligent systems, and educational data mining are not efficient. The analysis show that while intelligent decision support systems is highly heterogeneous intelligent systems and educational data mining show no heterogeneity. It is observed from the studies that intelligent decision support systems although vary in intervention effects they have smaller datasets and employ fewer standard algorithms than studies in machine learning. Scope exists for future studies in all areas of artificial intelligence to proffer solutions in varied u-learning environments and technologies.

## ACKNOWLEDGMENT

Kind acknowledgement goes to the Durban University of Technology for making the resources available for this research project.

## REFERENCES

- [1] Ogata, H. and Yuno, Y., Context-Aware Support for Computer-Supported Ubiquitous Learning., in Proceedings of the 2nd IEEE International Workshop on Wireless and Mobile Technologies in Education. 2004. p. 27 - 34.
- [2] Chiu, P.S., Kuo, Y., Huang, Y. and Chen, T., A Meaningful Learning based u-Learning Evaluation Model, in Eighth IEEE International Conference on Advanced Learning Technologies. 2008. p. 77-81.
- [3] Yahya, S., Ahmad, E.A. and Jalil, K.A., The definition and characteristics of ubiquitous learning: A discussion. International Journal of Education and Development using Information and Communication Technology, 2010. 6(1): p. 117-127.
- [4] Miller, M.D., Going Online in a Hurry: What to Do and Where to Start. The chronical of higher education, Moving online now., 2020.
- [5] Tabuenca, B., Wu, L. and Tovar, E., The PRISMA: A Visual Feedback Display for Learning Scenarios. in Proceedings of 13th International Conference on Ubiquitous Computing and Ambient Intelligence UCAmI 2019, 2019. Toledo, Spain.
- [6] Meyliana, M., Hidayanto, A.N. and Budiardjo, E.K., Evaluation of social media channel preference for student engagement improvement in universities using entropy and TOPSIS method. Journal of Industrial Engineering and Management, 2015. 8(5): p. 1676-1697.
- [7] Rad, M.S., Dahlan, H.M, Iahad, N.A, Nilashi, M. and Ibrahim, O., Using a Multi-Criteria Decision-Making Approach for Assessing the Factors Affecting Social Network Sites Intention to Use. Journal of Soft Computing and Decision Support Systems 2015. 2(3): p. 20-28.
- [8] Elyazgi, M.G., Nilashi, M., Ibrahim, O., Rayhan, A. and Elyazgi, S., Evaluating the Factors Influencing E-book Technology Acceptance among School Children Using TOPSIS Technique. Journal of Soft Computing and Decision Support Systems., 2016. 3(2): p. 11-25.
- [9] Omorogbe, D.E.A.a.I., L.I. , Parents Preference for Students' Choice of Urban Schools in Benin City, Nigeria: Integrated AHP Intuitionistic Fuzzy TOPSIS. An International Multi-disciplinary Journal, Ethiopia., 2016. 10(2): p. 254-265.
- [10] Pires, J.M.a.C., M.P. , "Intelligent" Adaptive Learning Objects applied to Special Education needs: Extending the eLearning paradigm to the uLearning environment, in 11th Iberian Conference on Information Systems and Technologies (CIISTI). 2016: Las Palmas, Spain. p. 1-6.
- [11] Angeli, C., Howard, S.K. Ma, J., Yang, J. and Kirschner, P.A. , Data mining in educational technology classroom research: Can it make a contribution? . Computer and Education, 2017. 113: p. 226-242.
- [12] Belachew, E.B.a.G., F.A. , Student Performance Prediction Model using Machine Learning Approach: The Case of Wolkite University. International Journal of Advanced Research in Computer Science and Software Engineering, 2017. 7(2): p. 46-50.
- [13] El Mabrouk, M., Gaou, S. and Ritli, M.K. , Towards an Intelligent Hybrid Recommendation System for E-Learning Platforms Using Data Mining. International Journal of Emerging technologies in Learning, 2017. 12(6): p. 52-76.
- [14] Xu, J., Moon, K.H. and van der Schaar, M., A Machine Learning Approach for Tracking and Predicting Student Performance in Degree Programs. IEEE Journal of Selected Topics in Signal Processing, 2017. 11(5): p. 742-753.
- [15] Khoshi, A., Gooshki, H.S. and Mahmoudi, N. , The data on the effective qualifications of teachers in medical sciences: An application of combined fuzzy AHP and fuzzy TOPSIS methods. . Data in Brief, 2018. 21(2018): p. 2689-2693.
- [16] Mao, Y., Lin, C. and Chi, M. , Deep Learning vs. Bayesian Knowledge Tracing: Student Models for Interventions. Journal of Educational Data Mining. , 2018. 10(2): p. 29-54.
- [17] Mohamed, H. and Lamia, M. , Implementing flipped classroom that used an intelligent tutoring system into learning process Computers and Education, 2018. 124(2018): p. 62-76.
- [18] Sirait, A.D.S., Fitriani, W.R., Hidayanto, A.N., Purwandari, B. and Kosandi, M. , Evaluation of social media preference as e-participation channel for students using fuzzy AHP and TOPSIS. In 4th International Conference on Computing, Engineering, and Design, ICCED 2018. 2019: Bangkok: Institute of Electrical and Electronics Engineers Inc. p. 158-163.
- [19] El Aissaoui, O., Oughdir, L. and El Alloui, Y. , A fuzzy classification approach for learning style prediction based on web mining technique in e-learning environments. . Education and Information Technologies, 2019. 24(2019): p. 1943-1959.
- [20] Purbasari, I.Y., Puspaningrum, E.Y. and Putra, A.B.S. , Using Self-Organizing Map (SOM) for Clustering and Visualization of New Students based on Grades. Journal of Physics: Conference Series, 2020. 1569(2020): p. 1-6.
- [21] El Aissaoui, O., El Alami El Madani, Y., Oughdir, L., Dakkak, A. and El Alloui, Y. , A Multiple Linear Regression-Based Approach to Predict Student Performance. Advanced Intelligent Systems for Sustainable Development, 2020. 1102(2020): p. 9-23.
- [22] Chatzigeorgiou, I.M. and and Andreou, G.T., A systematic review on feedback research for residential energy behavior change through mobile and web interfaces. Renewable and Sustainable Energy Reviews, 2020. 135(2021): p. 1-16.
- [23] Moher, D., Shamseer, L., Clarke, M., Ghersi, D., Liberati, A., Petticrew, M., Shekelle, P., Stewart, L.A. and PRISMA-P Group. , Preferred reporting items for systematic review and meta-analysis protocols (PRISMA-P) 2015 statement. Systematic Reviews, 2015. 4(1): p. 1-9.
- [24] Gundogan, B., Fowler, A. and Agha, R. , Assessing the compliance of systematic review articles published in leading dermatology journals with the PRISMA statement guidelines: A systematic review protocol. International Journal of Surgery Protocols, 2018. 10-12(2018): p. 1-4.

- [25] Manriquez, J., Andino-Navarrete, R., Cataldo-Cerda, K. and Harz-Fresno, I. , Bibliometric characteristics of systematic reviews in dermatology: A cross-sectional study through Web of Science and Scopus. *Dermatologica Sinica*, 2015. 33(2015): p. 154-156.
- [26] Lee, Y.-H., Hsieh, Y-C and Hsu, C-N. , Adding Innovation Diffusion Theory to the Technology Acceptance Model: Supporting Employees' Intentions to use E-Learning Systems. *Journal of Educational Technology and Society*., 2011. 14(4): p. 124-137.
- [27] Moher, D., Liberati, A., Tetzlaff, J., Altman, D.G. and Group P. , Preferred reporting items for systematic reviews and meta-analyses: the PRISMA statement. *The British Medical Journal*., 2009. 39(9): p. 873–880.
- [28] DerSimoniana, R. and Laird, N. , Meta-analysis in clinical trials revisited. . *Contemporary Clinical Trials* ., 2015. 45: p. 139-145.
- [29] Verhagen, A. and Ferreira, M.L. , Forest Plots. *Journal of Physiotherapy* ., 2014. 60(3): p. 170-173.



# Hate Speech Detection in Twitter using Transformer Methods

Raymond T Mutanga<sup>1</sup>, Nalindren Naicker<sup>2</sup>, Oludayo O Olugbara<sup>3</sup>  
ICT and Society Research Group, Department of Information Systems  
Durban University of Technology, South Africa  
Durban, 4000

**Abstract**—Social media networks such as Twitter are increasingly utilized to propagate hate speech while facilitating mass communication. Recent studies have highlighted a strong correlation between hate speech propagation and hate crimes such as xenophobic attacks. Due to the size of social media and the consequences of hate speech in society, it is essential to develop automated methods for hate speech detection in different social media platforms. Several studies have investigated the application of different machine learning algorithms for hate speech detection. However, the performance of these algorithms is generally hampered by inefficient sequence transduction. The Vanilla recurrent neural networks and recurrent neural networks with attention have been established as state-of-the-art methods for the assignments of sequence modeling and sequence transduction. Unfortunately, these methods suffer from intrinsic problems such as long-term dependency and lack of parallelization. In this study, we investigate a transformer-based method and tested it on a publicly available multiclass hate speech corpus containing 24783 labeled tweets. DistilBERT transformer method was compared against attention-based recurrent neural networks and other transformer baselines for hate speech detection in Twitter documents. The study results show that DistilBERT transformer outperformed the baseline algorithms while allowing parallelization.

**Keywords**—Attention transformer; deep learning; neural network; recurrent network; sequence transduction

## I. INTRODUCTION

Social media platforms such as Twitter are publicly accessible digital resources for online communication and collaboration. Despite its popularity and convenience, Twitter is increasingly being used to spread hate speech. The level of anonymity granted by Twitter makes it conducive for the dissemination of hateful speech about people. Furthermore, a proportional relationship between hate speech propagation and the occurrence of hate-related crimes is highlighted in other studies [1, 2]. Given the high volume and nature of messages posted on Twitter, it is imperative to develop ways to curb the dissemination of hateful messages.

Currently, social media companies such as Twitter and Facebook employ human annotators to manually delete messages deemed to be hateful [3]. Moreover, users of these platforms are encouraged to flag and report contents they perceive to be inimical to the public. Nevertheless, these methods are labor-intensive and subject to human judgment [4]. The grave consequences of hate speech propagation and inherent limitations of human annotators have necessitated the

development of automated hate speech detection methods that use machine learning algorithms. Machine learning algorithms can be classified into two broad categories, which are classical machine learning and deep learning. Both methods have been exploited and tested for hate speech detection in earlier studies.

Classical algorithms depend on feature engineering, a process which is complex and time-consuming. The complexity of the feature engineering process negatively impacts the capture of semantic and syntactic text representations [5]. Deep learning algorithms perform end-to-end training by allowing highly predictive representations to be effectively coded. Deep neural networks such as recurrent neural network (RNN) can preserve sequence information over time, thereby integrating contextual information better in classification tasks [6]. However, their inherently sequential nature prohibits parallelization, thereby increasing processing time. Moreover, RNN suffers from the limitation of long-term dependency, making it less effective as the hiatus between where information appears and the point where the information is required increases. This is particularly important in context-dependent applications such as hate speech detection.

Researchers have created techniques based on recurrent neural networks in conjunction with the attention mechanism to solve some of these problems. Such attention-based recurrent neural networks allow for the modeling of dependencies regardless of distance between the input or output sequences [7, 8]. However, the inclusion of recurrent neural network prohibits parallel processing and negatively impacts processing time. Due to such limitations, some recent works have focused on improving attention mechanisms. Research in this direction has given birth to transformers that perform sequence transduction entirely based on attention [9]. This allows for capturing relevant information that might be contained in every word within a sentence while allowing parallel processing. Consider the following statements for an example. “Foreigners must fall. They are taking our jobs. Some of them are stealing from us”. Current approaches which are mostly based on traditional deep learning algorithms fail to capture that the word “Some” in the third sentence refers to foreigners because of their reliance on past hidden states to capture dependencies with previous words. Transformer algorithms are designed to capture such long-term dependencies using positional embedding to remember word order in sequences [9]. In addition, parallelization enables faster training when compared to traditional deep learning approaches that are based on sequential processing [9].

Consequently, this research seeks to enhance hate speech detection by capturing long-term dependencies using transformer methods while allowing parallel processing. The remainder of this paper is organized as follows. Section II reviews the related works. Section III describes the materials and methods of the study. Experiments and results of experiments are presented in Section IV. Finally, Section V presents the conclusions and future works.

## II. RELATED WORK

Supervised machine learning techniques are the predominant approach used for automated hate speech detection [10]. Hate speech detection can be modeled in machine learning as a dichotomous class or multiclass classification problems that can be addressed adequately using either classical learning algorithms or deep learning algorithms. Classical learning algorithms rely on manually engineered features while deep learning algorithms automatically learn features from the input data, instead of adopting handcrafted features [5]. The unstructured nature of human language presents many intrinsic challenges to automated text classification methods. One key challenge faced by existing methods of hate speech detection is the failure to capture long-term dependencies. This leads to loss of contextual information, which is vital for semantic interpretation. Deep learning algorithms, particularly the recurrent neural network (RNN) algorithms, have been the de-facto methods in handling sequence data such as text [11, 12]. However, they have been limited in the length of sequences they can capture [13]. Transformers are a promising way for capturing long-term dependencies in textual data. However, the technical barriers that need to be surmounted to adapt transformers to automated hate speech detection are not trivial.

RNN algorithms such as long short-term memory (LSTM) [14] and gated recurrent units [15] were developed specifically to address the problem of long-term dependencies which other machine learning algorithms suffer from. The Vanilla RNN works by assigning more weights to prior data points of a sequence, making it suitable for classification of textual data in a way that facilitates improved semantic analysis [16]. Nevertheless, RNN is prone to problems of exploding gradient and vanishing gradient during backpropagation training [17].

The LSTM has a chain-like structure of the Vanilla RNN, but further incorporates multiple gates to control the quantity of data that are allowed into every node state. LSTM is especially helpful in minimizing the vanishing gradient problem [18]. In addition, the LSTM preserves long-term dependencies efficaciously compared to the Vanilla RNN [16], thereby allowing the algorithm to capture more context. Despite these benefits, the LSTM cannot capture long term dependencies to arbitrary lengths. Specifically, the performance of the LSTM drops as sequence length increases beyond thirty words [7].

Further research that is aimed at addressing the problem of long-term dependencies has given birth to the attention mechanism [7]. Attention allows modeling of dependencies irrespective of the distance between input and output sequences [8]. Attention mechanisms work on the assumption that every word in each sentence is relevant. This allows for the capture

of context that may be necessary when classifying subjective text such as hate speech. Attention-based models have been investigated with success in text-related tasks [19-21]. However, they are used in conjunction with RNNs [9] and therefore are unable to process word sequences in parallel. For a large corpus of text, this may significantly affect the processing time.

Recent adaptations of attention approach have shifted methods progressively from RNNs to self-attention and transformers [22]. Transformer has rapidly become the dominant architecture for natural language processing (NLP), outperforming RNNs in natural language generation and natural language understanding [23]. The transformer architecture scales well with training data and model size while facilitating efficient parallelization and capturing long-range sequence features. In addition, transformers allow transfer learning by fine-tuning large pre-trained language models for downstream NLP tasks with a relatively small number of training examples, resulting in an improved performance regardless of dataset size [24]. This is particularly important when dealing with highly imbalanced datasets with few instances of hate speech.

There are several types of transformer methods that have been investigated with success in NLP research. Bidirectional encoder representations from text (BERT) [22] has surpassed previous performance benchmarks in common NLP tasks [25]. BERT uses vast unlabeled data for creating models whose parameters can be tuned as desired for smaller supervised data to improve performance. The success of BERT has led to the development of several algorithms based on BERT architecture. These algorithms include RoBERTa [26], DistilBERT [27] and XLNET [28]. RoBERTa is an enhancement of BERT, which is trained on a bigger dataset to improve performance while DistilBERT learns a streamlined version of BERT. XLNET is a generalized autoregressive pre-training method that aims to reconstruct the original data from corrupted input.

## III. MATERIALS AND METHODS

In this section, we present materials and methods used in this study, including acquisition and structure of experimental dataset and setup.

### A. Experimental Dataset

Multiclass hate speech and offensive (HSO) language dataset was used in this study for model validation. The dataset was developed, first used by authors in [29], and it was distributed through CrowdFlower. This dataset contains 24783 Twitter text messages that have been labeled into one of the following three classes: 'neutral', 'Offensive' and 'Hate' where 77.4% of the messages are labeled as 'neutral', 16.8% as 'Offensive' and 5.8% as 'Hate'. In this paper, the hate speech detection has been solved as a three classes classification problem.

### B. Experimental Setup

The proposed methods of this study were implemented using Python programming language. Keras library was used to implement the attention-based LSTM method. The proposed

method was implemented using Hugging face transformers class embedded in Python. Experiments were conducted on a computer running Windows 10 operating system with the configuration of Intel(R) Core (TM) i5-8250U CPU @ 1.60GHz (8 CPUs), 1.8GHz, 8 GB RAM and 500 Gigabytes hard disk drive.

### C. Preprocessing

Due to the colloquial nature of Twitter messages, Twitter data are highly unstructured and contain a lot of noise that can affect method accuracy. Consequently, it was deemed necessary to preprocess all Tweets to remove less predictive text features. Preprocessing is widely known to improve performance of classification methods [30] while reducing processing time. The labeled dataset was initially processed to normalize Twitter text as follows:

- Removal of the following noise characters i.e. :| : , ; & ; ! ? \.
- Normalization of hashtags into words, so that for instance, “#muslimsmustfall” becomes “Muslims must ‘fall’”. Hashtags are used to prefix tweets but do not give any valuable information. We have developed a python method to normalize hashtags.
- Lowercasing and stemming to reduce word inflexions.
- Removal of tokens that appear in the dataset less than 5 times that is elimination of low frequency terms at a threshold of 5.

### D. Proposed Method

Transformers mirror the standard NLP machine learning method pipeline that includes the processing of data, application of a method, and making predictions. This approach was selected because of its inherent self-attention mechanism that allows it to capture long-term dependencies while allowing parallel processing of input features. The capture of long-term dependencies allows the transformer to perform anaphora resolution, which is one of the major challenges in text processing [31]. However, most transformer models require substantial computational resources, thereby limiting their applicability in resource-constrained environments [32]. For example, they cannot run on portable devices. Furthermore, these models are slower at inference times, making them unfeasible for real-time situations. To address these problems, we propose DistilBERT a streamlined version of BERT that uses only half the number of parameters of BERT [27] but retains the performance of BERT in many text processing tasks [33] while making the inference 60% faster than BERT [34]. DistilBERT was created by removing token type embeddings and pooler from the default architecture of BERT [35]. DistilBERT further reduced the number of layers by 50%, thereby significantly reducing the footprint of the model.

In this study, we have used the distilBERT base uncased model with 66 million parameters pretrained on the Toronto Book Corpus and English Wikipedia [27]. The data used in the experiment was first preprocessed using the steps discussed in Section III of this paper. The preprocessed data were then split where a random 80% of instances were allocated for parameter

fine-tuning and training, while 20% random instances were allocated for evaluating the performance of the final model. Fig. 1 shows the architecture of the proposed DistilBERT method of this study. The hyperparameter fine tuning component of the architecture is essential and will be explained in the subsequent section.

### E. Model Parameters

The hyperparameter tuning is a crucial step when customizing pre-trained models to specific tasks. As shown in Table I, we optimized our method for hate speech detection by altering sequence length, batch size, early stopping patience and number of epochs. The maximum sequence length was set at 280 in line with the Twitter limit of allowable characters. The optimal number of epochs in our experiments was set to 4. We have configured the early stopping patience technique to prevent our method from overfitting. We set the early stopping patience value at 4, therefore, training is terminated if the evaluation loss fails to improve for four successive evaluations. The evaluation batch size defines the number of examples that are processed concurrently during the training session. In our experiments, we set the evaluation batch size to 256 because it was the largest batch size that our processor could handle effectively. The proposed method has been evaluated using five state of the art metrics outlined in Section IV of this paper.

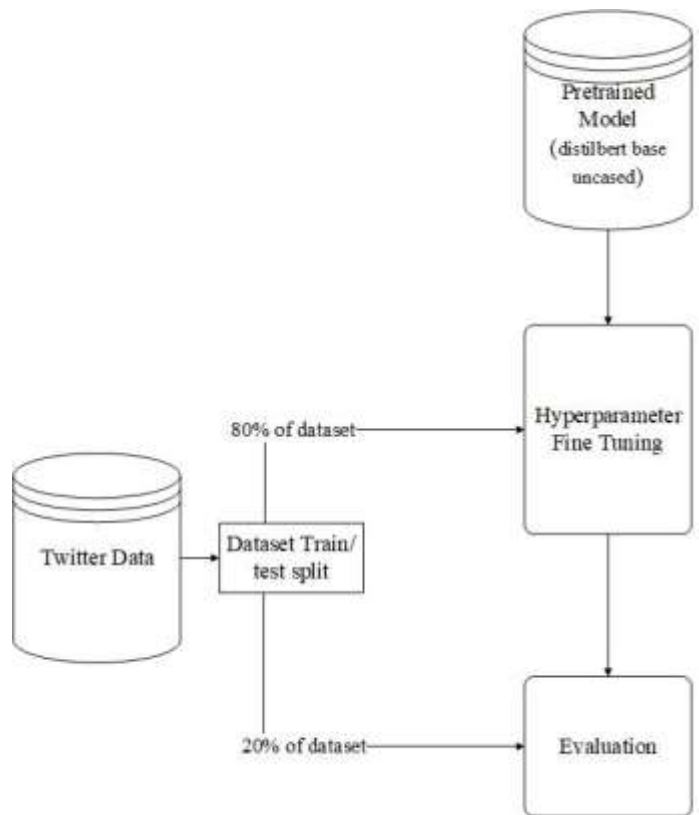


Fig. 1. Architecture of the Distilbert Hate Speech Detection Method.

TABLE I. PARAMETERS OF THE PROPOSED APPROACH

| Parameter | seq_length | epochs | Early_stopping | Batch_size |
|-----------|------------|--------|----------------|------------|
| Value     | 280        | 4      | 4              | 256        |

#### IV. RESULTS AND DISCUSSION

Results of the proposed DistilBERT method was compared against results computed by BERT, XLNet, RoBERTa and attention-based LSTM. We split the dataset in the ratio of 80:20 for model training and testing, respectively. The algorithms were analyzed in terms of six standard functional metrics of accuracy, precision, recall and F-measure, Mathews correlation coefficient (MCC) and evaluation loss. The results are presented based on the ability of the models to detect hate tweets.

##### A. Analysis of Accuracy

The experimental results of the proposed DistilBERT method, along with five baseline algorithms are presented in Table II and Fig. 2. It can be observed that the proposed DistilBERT method recorded the highest average accuracy of 92%. It is worth mentioning that the differences in accuracy scores for all transformer-based methods were negligible. This may be attributed to the fact that they all use standard extensively tested pre-trained models. Expectedly all transformer-based algorithms performed better than the LSTM with Attention. The least performing transformer method had an accuracy of 89%, which is superior to LSTM with attention which had 66% accuracy. This trend is because of the ability of the transformers to capture long-term dependencies better than LSTM with Attention.

##### B. Analysis of Precision

It can be observed from Table III and Fig. 3 that the DistilBERT (base-uncased) and XLNet algorithms jointly recorded the highest precision score of 75% whilst LSTM with attention recorded the least precision score of 65.9%. Although the LSTM with attention recorded the least result, it should be noted that this score is higher than scores recorded by methods using classical machine learning [29]. This result confirms the literature position that attention improves performance in NLP tasks [19].

TABLE II. ACCURACY SCORES OF FOUR BENCHMARK HATE SPEECH DETECTION ALGORITHMS AND THE PROPOSED DISTILBERT MODEL ON THE HSO DATASET

| Algorithm           | Method Name                 | Accuracy |
|---------------------|-----------------------------|----------|
| BERT                | bert-base-uncased           | 0.90     |
| RoBERTa             | robert-base                 | 0.91     |
| RoBERTa             | robert-base-openai-detector | 0.90     |
| XLNet               | xlm-mlm-en-2048             | 0.91     |
| LSTM with Attention |                             | 0.66     |
| DistilBERT          | distilbert-base-uncased     | 0.92     |

TABLE III. PRECISION SCORES OF FOUR BENCHMARK HATE SPEECH DETECTION ALGORITHMS AND THE PROPOSED DISTILBERT MODEL ON THE HSO DATASET

| Algorithm           | Method Name                 | Precision |
|---------------------|-----------------------------|-----------|
| BERT                | bert-base-uncased           | 0.74      |
| RoBERTa             | robert-base                 | 0.74      |
| RoBERTa             | robert-base-openai-detector | 0.72      |
| XLNet               | xlm-mlm-en-2048             | 0.75      |
| LSTM with Attention |                             | 0.66      |
| DistilBERT          | distilbert-base-uncased     | 0.75      |

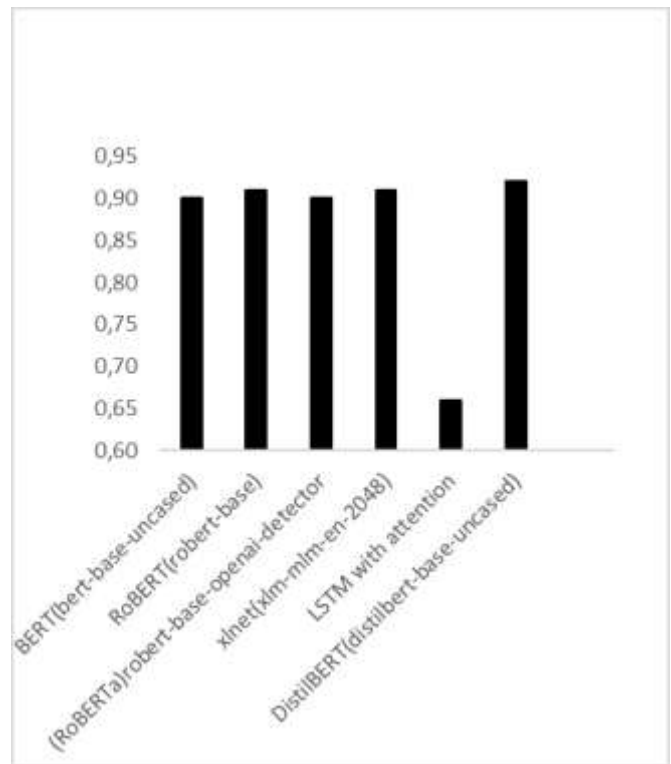


Fig. 2. Illustration of Accuracy Scores using Four Benchmarking Hate Speech Detection Algorithms and the Proposed DistilBERT Model.

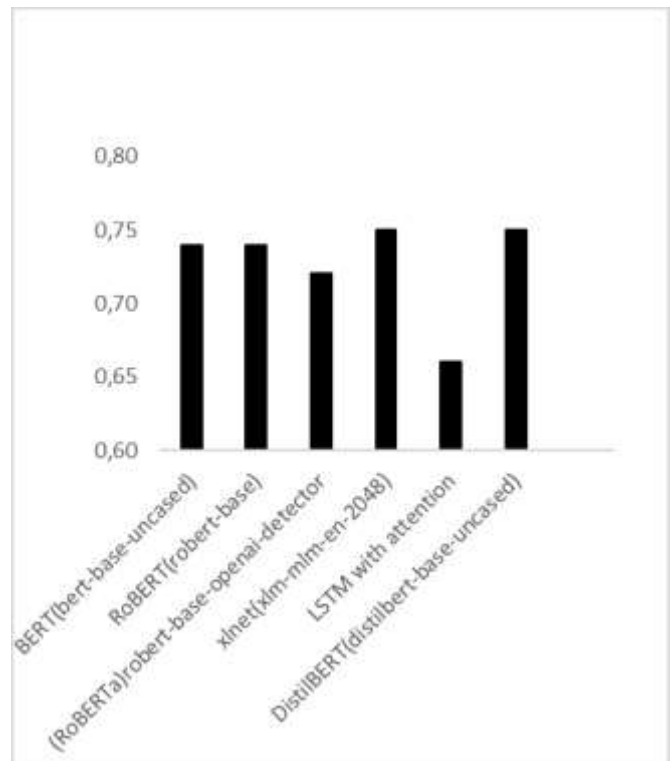


Fig. 3. Illustration of Precision Scores using Four Benchmarking Hate Speech Detection Algorithms and the Proposed DistilBERT Model.

C. Analysis of Recall

Results from Table IV and Fig. 4 show that DistilBERT and XLNet recorded the average recall score of 75% to demonstrate its superior over other algorithms explored in this study. LSTM with attention had the least recall score of 66%. Although the LSTM with attention performed inferior in our experiments, it should be noted that it performed superior to an earlier study on the same dataset for the task of hate speech detection [29].

D. Analysis of MCC Scores

Table I lists the MCC scores calculated for the overall test tweets selected from the experimental dataset. It can be observed that our proposed method recorded the highest MCC score of 75%. Fig. 5 shows that the difference in MCC scores for all algorithms explored in this study is negligible. The worst performing algorithm was RoBERTa (robert-base-openai-detector) which recorded a MCC score of 71% while the best performing algorithm was DistilBERT (distilbert-base-uncased) which recorded a MCC score of 75%.

TABLE IV. RECALL SCORES OF FOUR BENCHMARK HATE SPEECH DETECTION ALGORITHMS AND THE PROPOSED DISTILBERT MODEL ON THE HSO DATASET

| Algorithm           | Method Name                 | Recall |
|---------------------|-----------------------------|--------|
| BERT                | bert-base-uncased           | 0.72   |
| RoBERTa             | robert-base                 | 0.65   |
| RoBERTa             | robert-base-openai-detector | 0.63   |
| XLNet               | xlm-mlm-en-2048             | 0.69   |
| LSTM with Attention |                             | 0.66   |
| DistilBERT          | distilbert-base-uncased     | 0.75   |

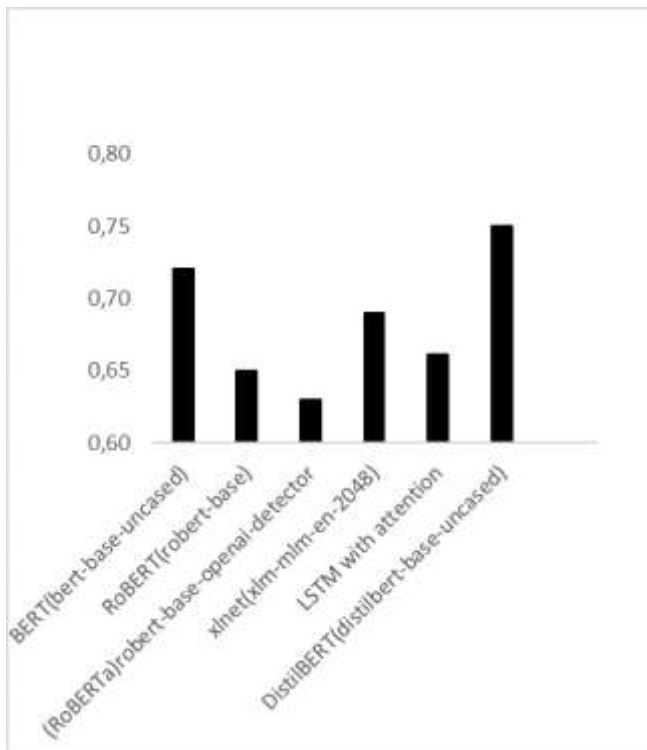


Fig. 4. Illustration of Recall Scores using four Benchmarking Hate Speech Detection Algorithms and the Proposed DistilBERT Model.

TABLE V. MCC SCORES OF FOUR BENCHMARK HATE SPEECH DETECTION ALGORITHMS AND THE PROPOSED DISTILBERT MODEL ON THE HSO DATASET

| Algorithm           | Method Name                 | MCC  |
|---------------------|-----------------------------|------|
| BERT                | bert-base-uncased           | 0.73 |
| RoBERTa             | robert-base                 | 0.73 |
| RoBERTa             | robert-base-openai-detector | 0.71 |
| XLNet               | xlm-mlm-en-2048             | 0.74 |
| LSTM with Attention |                             | 0.72 |
| DistilBERT          | distilbert-base-uncased     | 0.75 |

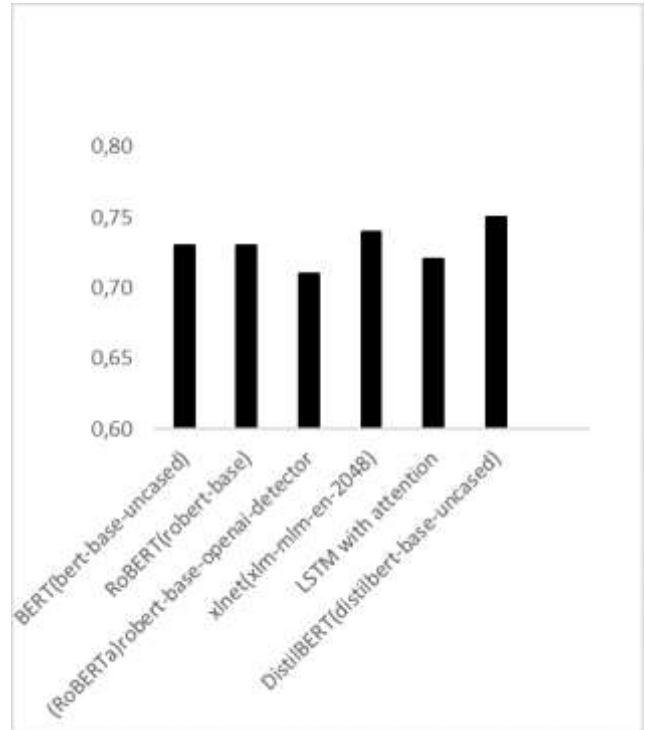


Fig. 5. Illustration of MCC Scores using Four Benchmarking Hate Speech Detection Algorithms and the Proposed DistilBERT Model.

E. Analysis of Evaluation Loss

Table VI shows evaluation loss recordings for the experiments carried out in this study. Fig. 6 clearly shows that our proposed method recorded the best (lowest) evaluation loss of 28% while the LSTM with attention recorded the worst evaluation loss of 36%. This shows that our proposed method maximized predictive capability while minimizing the misclassification error rate more than any of the baseline algorithms.

TABLE VI. EVALUATION LOSS SCORES OF FOUR BENCHMARK HATE SPEECH DETECTION ALGORITHMS AND THE PROPOSED DISTILBERT MODEL ON THE HSO DATASET

| Algorithm           | Method Name                 | Eval loss |
|---------------------|-----------------------------|-----------|
| BERT                | bert-base-uncased           | 0.32      |
| RoBERTa             | robert-base                 | 0.32      |
| RoBERTa             | robert-base-openai-detector | 0.33      |
| XLNet               | xlm-mlm-en-2048             | 0.31      |
| LSTM with Attention |                             | 0.36      |
| DistilBERT          | distilbert-base-uncased     | 0.28      |

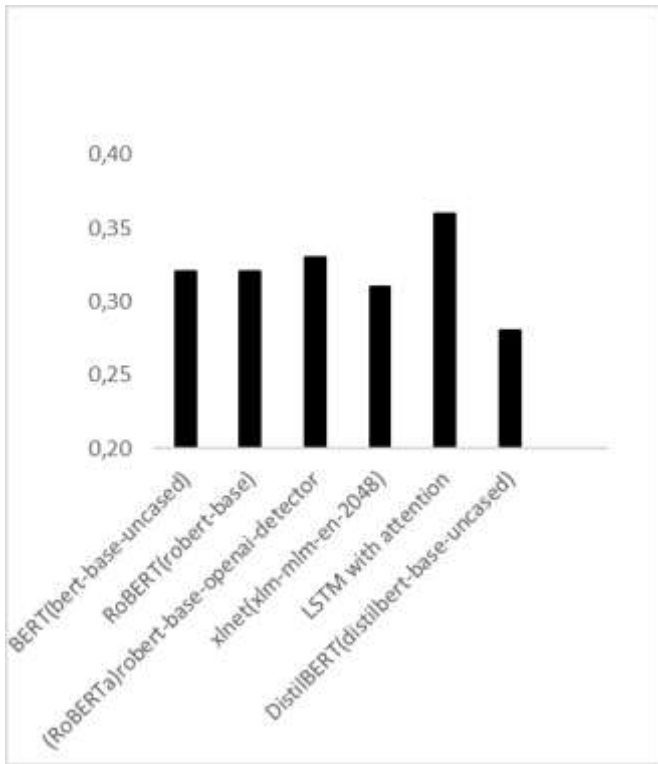


Fig. 6. Illustration of Evaluation Loss Scores using Four Benchmarking Hate Speech Detection Algorithms and the Proposed DistilBERT Model.

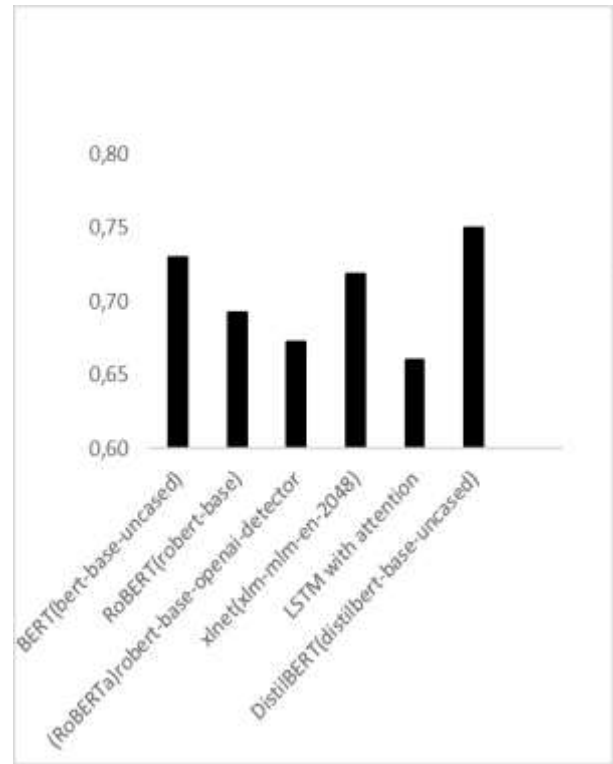


Fig. 7. F-Measure Scores Per Algorithm Illustration of F-Measure Scores using Four Benchmarking Hate Speech Detection Algorithms and the Proposed DistilBERT Model.

#### F. Analysis of F-Measure

Table VII and Fig. 7 show the F-measure scores of the algorithms explored in this study. It can be observed that the DistilBERT (distilbert-base-uncased) recorded the best F-measure score of 75% while LSTM with attention recorded the lowest F-measure score of 66%. Although DistilBERT has fewer layers and parameters, it outperformed all other transformer algorithms explored in this study. The superior performance of DistilBERT may be attributed to the chosen hyperparameters during experimentation. The same hyperparameters were used to train all the models. It can be argued that the used parameters are not necessarily the optimal combination of hyperparameters for each model explored in this study. Careful selection of the best hyperparameters may improve performance of models such as BERT and RoBERTa.

TABLE VII. F-MEASURE SCORES OF FOUR BENCHMARK HATE SPEECH DETECTION ALGORITHMS AND THE PROPOSED DISTILBERT MODEL ON THE HSO DATASET

| Algorithm           | Method Name                 | F-Measure |
|---------------------|-----------------------------|-----------|
| BERT                | bert-base-uncased           | 0.73      |
| RoBERTa             | robert-base                 | 0.69      |
| RoBERTa             | robert-base-openai-detector | 0.67      |
| XLNet               | xlm-mlm-en-2048             | 0.72      |
| LSTM with Attention |                             | 0.66      |
| DistilBERT          | distilbert-base-uncased     | 0.75      |

Comparative results based on five different metrics from this work show that the transformer models consistently outperform the LSTM with attention. The superior performance of transformer demonstrates that limitations of LSTM, which are inefficient sequence transduction and lengthy processing time have been adequately addressed by the transformer method in hate speech detection.

#### V. CONCLUSION AND FUTURE WORK

Given the societal implications of hate speech, it is crucial that systems that can accurately distinguish between hate speech, offensive language and neutral speech are developed. Despite concerted efforts from social media companies, governments, and academia, hate speech detection remains a challenging problem in the society of today. In this paper, we have explored several transformer-based methods for hate speech detection. We have evaluated the effectiveness of our method using six state of the art metrics. The results showed that the DistilBERT, a distilled version of BERT, outperforms all transformer-based baseline methods and the attention-based LSTM explored in this study. We, therefore, conclude that the proposed method can be used to learn effective information for the classification of hate speech in resource-constrained environments because it is computationally inexpensive. In addition, transformers facilitate transfer learning, allowing them to be used where training data is limited. It is common for

hate speech on social media to be expressed in more than one language. For example, most people in Africa codeswitch their native languages with French, Portuguese, or English language. In future work, we plan to explore multilingual pre-trained models for the task of hate speech detection. The data used in this study were limited to textual Twitter texts only, whereas hate speech on Twitter may be expressed through different data formats such as images and videos. For example, a user may post a video inciting hate speech on Twitter and still go undetected. This limitation calls for the development of multimodal datasets that include other formats of data. Future study will develop methods that integrate both textual and image data for hate speech detection.

#### ACKNOWLEDGMENTS

We would like to thank the Durban University of Technology for making funding opportunities and materials for experiments available for this research project. Kind acknowledgement to Dr. T. Adeliyi for his contributions.

#### REFERENCES

- [1] H. Hosseinmardi, S. A. Mattson, R. I. Rafiq, R. Han, Q. Lv, and S. Mishra, "Detection of cyberbullying incidents on the instagram social network," arXiv preprint arXiv:1503.03909, 2015.
- [2] Z. Zhang, D. Robinson, and J. Tepper, "Detecting hate speech on twitter using a convolution-gru based deep neural network," in European semantic web conference, 2018: Springer, pp. 745-760.
- [3] Z. Waseem and D. Hovy, "Hateful symbols or hateful people? predictive features for hate speech detection on twitter," in Proceedings of the NAACL student research workshop, 2016, pp. 88-93.
- [4] G. K. Pitsilis, H. Ramampiaro, and H. Langseth, "Effective hate-speech detection in Twitter data using recurrent neural networks," Applied Intelligence, vol. 48, no. 12, pp. 4730-4742, 2018.
- [5] T. Young, D. Hazarika, S. Poria, and E. Cambria, "Recent trends in deep learning based natural language processing," iee Computational intelligence magazine, vol. 13, no. 3, pp. 55-75, 2018.
- [6] S. Sohagir, D. Wang, A. Pomeranets, and T. M. Khoshgoftaar, "Big Data: Deep Learning for financial sentiment analysis," Journal of Big Data, vol. 5, no. 1, p. 3, 2018.
- [7] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," arXiv preprint arXiv:1409.0473, 2014.
- [8] Y. Kim, C. Denton, L. Hoang, and A. M. Rush, "Structured attention networks," arXiv preprint arXiv:1702.00887, 2017.
- [9] A. Vaswani et al., "Attention is all you need," in Advances in neural information processing systems, 2017, pp. 5998-6008.
- [10] A. Schmidt and M. Wiegand, "A survey on hate speech detection using natural language processing," in Proceedings of the Fifth International Workshop on Natural Language Processing for Social Media, 2017, pp. 1-10.
- [11] X. Wang, W. Jiang, and Z. Luo, "Combination of convolutional and recurrent neural network for sentiment analysis of short texts," in Proceedings of COLING 2016, the 26th international conference on computational linguistics: Technical papers, 2016, pp. 2428-2437.
- [12] M. S. Islam, S. S. S. Mousumi, S. Abujar, and S. A. Hossain, "Sequence-to-sequence Bangla sentence generation with LSTM recurrent neural networks," Procedia Computer Science, vol. 152, pp. 51-58, 2019.
- [13] Z. Mhammedi, A. Hellicar, A. Rahman, and J. Bailey, "Efficient orthogonal parametrisation of recurrent neural networks using householder reflections," in International Conference on Machine Learning, 2017, pp. 2401-2409.
- [14] S. Hochreiter and J. Schmidhuber, "Long short-term memory," Neural computation, vol. 9, no. 8, pp. 1735-1780, 1997.
- [15] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, "Empirical evaluation of gated recurrent neural networks on sequence modeling," arXiv preprint arXiv:1412.3555, 2014.
- [16] K. Kowsari, K. Jafari Meimandi, M. Heidarysafa, S. Mendu, L. Barnes, and D. Brown, "Text classification algorithms: A survey," Information, vol. 10, no. 4, p. 150, 2019.
- [17] 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.
- [18] R. Pascanu, T. Mikolov, and Y. Bengio, "On the difficulty of training recurrent neural networks," in International conference on machine learning, 2013, pp. 1310-1318.
- [19] A. M. Rush, S. Chopra, and J. Weston, "A neural attention model for abstractive sentence summarization," arXiv preprint arXiv:1509.00685, 2015.
- [20] Y. Wang, M. Huang, X. Zhu, and L. Zhao, "Attention-based LSTM for aspect-level sentiment classification," in Proceedings of the 2016 conference on empirical methods in natural language processing, 2016, pp. 606-615.
- [21] B. Yang, L. Wang, D. Wong, L. S. Chao, and Z. Tu, "Convolutional self-attention networks," arXiv preprint arXiv:1904.03107, 2019.
- [22] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," arXiv preprint arXiv:1810.04805, 2018.
- [23] T. Wolf et al., "Transformers: State-of-the-art natural language processing," arXiv preprint arXiv:1910.03771, 2019.
- [24] W. Antoun, F. Baly, and H. Hajj, "AraBERT: Transformer-based model for Arabic language understanding," arXiv preprint arXiv:2003.00104, 2020.
- [25] T. Rajapakse, "To Distil or Not To Distil: BERT, RoBERTa, and XLNet." <https://towardsdatascience.com/to-distil-or-not-to-distil-bert-roberta-and-xlnet-c777ad92f8> (accessed 28 July, 2020).
- [26] Y. Liu et al., "Roberta: A robustly optimized bert pretraining approach," arXiv preprint arXiv:1907.11692, 2019.
- [27] 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.
- [28] Z. Yang, Z. Dai, Y. Yang, J. Carbonell, R. R. Salakhutdinov, and Q. V. Le, "Xlnet: Generalized autoregressive pretraining for language understanding," in Advances in neural information processing systems, 2019, pp. 5753-5763.
- [29] T. Davidson, D. Warmesley, M. Macy, and I. Weber, "Automated hate speech detection and the problem of offensive language," in Eleventh international aaai conference on web and social media, 2017.
- [30] A. K. Uysal and S. Gunal, "The impact of preprocessing on text classification," Information Processing & Management, vol. 50, no. 1, pp. 104-112, 2014.
- [31] E. Cambria, S. Poria, A. Gelbukh, and M. Thelwall, "Sentiment analysis is a big suitcase," IEEE Intelligent Systems, vol. 32, no. 6, pp. 74-80, 2017.
- [32] H. Sajjad, F. Dalvi, N. Durrani, and P. Nakov, "Poor Man's BERT: Smaller and Faster Transformer Models," arXiv preprint arXiv:2004.03844, 2020.
- [33] B. Cheang, B. Wei, D. Kogan, H. Qiu, and M. Ahmed, "Language Representation Models for Fine-Grained Sentiment Classification," arXiv preprint arXiv:2005.13619, 2020.
- [34] S. Minaee, N. Kalchbrenner, E. Cambria, N. Nikzad, M. Chenaghlu, and J. Gao, "Deep learning based text classification: A comprehensive review," arXiv preprint arXiv:2004.03705, 2020.
- [35] B. Büyükoç, A. Hürriyetöglü, and A. Özgür, "Analyzing ELMo and DistilBERT on Socio-political News Classification," in Proceedings of the Workshop on Automated Extraction of Socio-political Events from News 2020, 2020, pp. 9-18.

# VerbNet based Citation Sentiment Class Assignment using Machine Learning

Zainab Amjad<sup>1</sup>, Imran Ihsan<sup>2</sup>  
Department of Creative Technologies  
Air University, Islamabad, Pakistan

**Abstract**—Citations are used to establish a link between articles. This intent has changed over the years, and citations are now being used as a criterion for evaluating the research work or the author and has become one of the most important criteria for granting rewards or incentives. As a result, many unethical activities related to the use of citations have emerged. That is why content-based citation sentiment analysis techniques are developed on the hypothesis that all citations are not equal. There are several pieces of research to find the sentiment of a citation, however, only a handful of techniques that have used citation sentences for this purpose. In this research, we have proposed a verb-oriented citation sentiment classification for researchers by semantically analyzing verbs within a citation text using VerbNet Ontology, natural language processing & four different machine learning algorithms. Our proposed methodology emphasizes the verb as a fundamental element of opinion. By developing and assessing the proposed methodology and according to benchmark results, the methodology can perform well while dealing with a variety of datasets. The technique has shown promising results using Support Vector Classifier.

**Keywords**—Citation content analysis; sentiment analysis; semantic analysis; ontology; natural language processing

## I. INTRODUCTION

Sentiment Analysis is a method to categorize and recognize feelings, thoughts, ideas, or sentiments conveyed in a text, to determine the writer's intentions. Sentiment analysis depends on sentiment polarity and sentiment score [1]. Sentiment polarity [2] is the emotion expressed in a text, it can be positive, negative, or neutral, while sentiment score is based on one of the three models; Bag-of-words (BOW) model [3], part-of-speech (POS) model [4], and semantic relationships. In the Bag-of-words model, a text is described as the bag of its words, irrespective of grammar and word organization. POS tagging model identifies words in each language as one of many groups to define the role of a word. Categories of part-of-speech in the English language include nouns, adjectives, verbs, adverbs, etc. [5]. The last model is the semantic relationship, it is an association between the meanings of words.

Citation is a reference to a published source or even an unpublished one [6]. "Citation Sentiment Analysis" deals with the relationship between the citing paper and the cited paper to measure the quality of published work. Researchers usually need to analyze numerous scientific papers to find relevant articles to their work of research. Due to the significantly growing number of scientific papers, this task of analysis is

time-consuming and complicated. To resolve this issue there exists many researchers [7]–[9] who deal with the sentiment analysis of citation sentences to improve bibliometric measures. Such applications can help scholars in the period of research to identify the problems with the present approaches, unaddressed issues, and the present research gaps [10].

There are two existing approaches for Citation Sentiment Analysis: Qualitative and Quantitative [7]. Quantitative approaches consider that all citations are equally important while qualitative approaches believe that all citations are not equally important [9]. The quantitative approach uses citation count to rank a research paper [8] while the qualitative approach analyzes the nature of citation [10].

However, qualitative analysis of a citation is deeper than the simple sentiment analysis of a citation sentence. There is a need to explore the reason for a citation [9]. Charles [11] is an author of the book titled "The Informed Writer", wrote in his book "It is you who decides; what materials you need, discovers the connections between different pieces of information, evaluates the information". Thus, the author of a research paper creates a cognitive relationship between the citing paper and the cited paper while citing. Another research suggests that authors use verbs to assert their sentiment while citing another research [12], [13]. Therefore, verbs are the most important grammatical terms used in a research paper to express a stance towards another research and to provide a rhetorical context. The choice of a verb in a citing sentence plays an important role. Using Part-of-Speech tagging, it is now possible to tag verbs in a citing sentence using Natural Language Processing techniques. Combining the sentiment polarity and verbs in a citation sentence can help to understand the true nature of the author's intent.

This research aims to replace traditional citation sentiment analysis techniques by taking an ontological approach by using VerbNet Ontology and Mapping Graph [9] between verbs used within a citation to formulate opinions and its evaluation model that can identify the role of verbs in citation sentiment analysis. Section 2 describes the literature review and Section 3 has our proposed methodology. In Sections 4 and 5 experiments and results are delineated. Section 6 concludes the paper.

## II. LITERATURE REVIEW

ACL Anthology Network dataset is a collection of 8736 citations from 310 research papers [10]. This sentiment corpus is a manually created dataset that can be used for automatic



classification citation sentences. In the experiments, using supervised classifiers an F-Score of 0.797 was achieved using 10-fold cross-validation. Later, a context-enhanced citation sentiment detection was performed on the same dataset [14]. In this experiment, the dominant sentiment in the citation is considered as the context that represents more than one sentiment in a citation. The effect of context windows of different lengths on the performance of a sentiment analysis system was also studied [15].

Niket Tandon and Ashish Jain [16] proposed a new technique to generate a structured summary of research papers. The proposed methodology classified citation context into one or more of five classes using a Language Model (LM) approach. Random k-Label sets with Naïve Bayes algorithm was used as the baseline to achieve 68.5% average precision. Xiaojun Wan & Fang Liu [17] used the Regression method to automatically evaluate the strength value each citation, and the strength value was used to measure the significance and influence of paper and the author. For this purpose, the Support Vector Regression method [18] was used. Bilal Hayat [19] proposed a novel automated method for the classification of citation sentiments as positive and negative. Sentiment lexicon was used to classify the citation by picking a window size of five sentences and for sentiment analysis, the Naïve Bayes classifier was used. The technique was assessed on a manually annotated dataset that consists of 150 research papers and the results depicted 80% accuracy. Cheol Kim and George R. Thoma [20] presented an automated technique to classify the sentiments articulated in Comment-on sentences using the Support Vector Machine (SVM) with a Radial Basis Kernel Function (RBF) and a Bag-of-Words input features constructed on n-grams word statistics. Jun Xu [21] presented the citation sentiment analysis of the citations in clinical research papers. For this purpose, the discussion section from 285 clinical trial papers was selected and extracted the n-grams, sentiment lexicons, and structure features. The citations were classified using Machine Learning methods and performance was evaluated using the 10-fold cross-validation method to achieve 0.8 Micro F-score and 0.719 Macro F-score.

Marco Valenzuela [22] proposed a supervised classification method that states the task of classifying meaningful citations with either two classes (important vs. non-important citation) or four classes (incidental: related work, incidental: comparison, important: using the work, important: extending the work.) Their approach used both direct citations and indirect citations. They achieved a precision of 65% for a recall of 90%. Faiza Qayyum and Muhammad Tanvir Afzal [7] presented a binary citation classification approach, using metadata-based parameters and cue-terms. Their work is close to the approach proposed by Valenzuela [22] which is the combination of metadata and content-based features, also used two types of parameters: Metadata based parameters (Titles, Authors name, Keywords, Categories, and References) and content-based parameters (Abstract and Cue-phrases). The experiments are performed on two annotated data sets, which were evaluated by using SVM, KLR, and Random Forest classifiers. The proposed model achieved 0.68 precision.

In 2018, Zehra Taskin [23] conducted a content-based citation analysis for Turkish research and they concluded that using computational linguistics for the evaluation of citation contexts provides better results. They divided the citation text into for main classes, meaning, purpose, shape, array. This research was significant for the evaluation of citation text by context. Imran Ihsan [9] proposed a Citation's Context and Reasons Ontology (CCRO) that helped to identify citations' relations using dominant verbs from citation sentences. The proposed ontology created 8 classes all extracted from Positive, Negative, and Neutral sentiments. The extracted verb was mapped to the relevant classes in CCRO based on the sentiment of the verb in a citation text. The results illustrate that the proposed ontology is reliable and complete.

VerbNet [24] is an ontology-based on Stanford Linguist Beth Levins's English Verb Classes [25]. The ontology is a lexical resource that includes both semantic and syntactic information about its contents that houses over 230 verb classes. CCRO [9] has created a knowledge-based known as "Mapping Graph" among the verbs with predicative complements in the English Language, the verbs extracted from the selected corpus using NLP and CCRO classes. Combining VerbNet Ontology and Mapping Graph proposed in CCRO, this research uses Natural Language Processing techniques to extract and map verbs within a citation sentence for semantic-based citation sentiment analysis using various machine learning algorithms.

### III. PROPOSED METHODOLOGY

Fig. 1 displays the main process blocks of our proposed methodology. The methodology has four major blocks: datasets, preprocessing, feature selection, and machine learning algorithms. Details of individual blocks are.

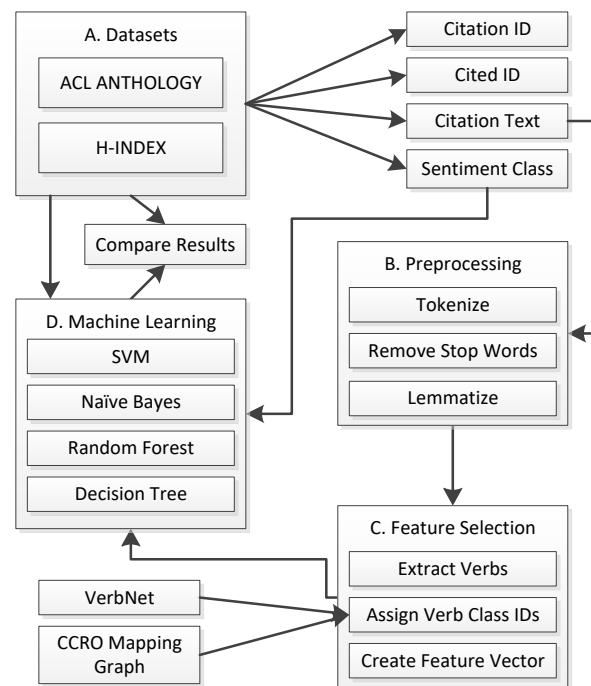


Fig. 1. Methodology.

### A. Datasets

Two datasets are employed. One is the publicly available ACL Anthology Dataset while the second is the manually curated H-Index dataset. ACL Anthology Dataset comprises of all the papers published by ACL and Computational Linguistics journal. Athar [26] manually constructed a dataset comprising of 8738 citation sentences, labeled with Citing Paper ID, Cited Paper ID, Citation Sentences, and their sentiment polarity (Positive, Negative, and Neutral). The second dataset [27] is a specific version of the ANN dataset [13] comprising of 701 citation sentences with their sentiment polarity. The distribution of all three classes in both datasets is shown in Fig. 2. Kindly note, the two datasets are employed for comparative study purposes only.



Fig. 2. Sentiment Class Distribution in both Datasets.

### B. Preprocessing

After selecting datasets next step is pre-processing on citation texts. The process comprises four steps. The first step is punctuation removal. Punctuation includes full stop, comma, and brackets, etc. used in writing to separate sentences and to clarify meaning. The second step is splitting up a sequence of citation text strings into pieces such as words, symbols called tokens. The third step is stop-words removal where commonly appearing words like 'is', 'a', 'an', 'it', 'which', etc. are considered as stop words and removed. The presence of stop words induces extra noise in different NLP problems that can negatively affect the results. In the last step, all the words in citing sentences are changed in their root terms. It does not simply chop off variations but uses a lexical knowledge like WordNet to gain an accurate form of words.

### C. Feature Selection

In feature selection, the first step is to extract verbs from tokenized citation sentences and can be achieved using part-of-speech tagging (POS). POS is also known as grammatical tagging. This technique marks the words from a text to a specific part of speech. In our experiments, only verbs are tagged. After tagging the next step is to assign a class ID using VerbNet. The VerbNet maps the verbs to their corresponding class. It is a lexical resource that includes both semantic and syntactic information about its contents. For this mapping, a Mapping Graph is used. Using the knowledge base and the extracted verbs in each sentiment, the Mapping Graph has been formulated [9] that provides a high level of abstraction

on CCRO classes. Based on the citation context, one such property can be attributed to multiple classes. Therefore, the combination becomes a graph rather than a tree where one individual verb can belong to multiple classes based on the citation's sentiment, making the class semantically coherent.

### D. Machine Learning

Based on our literature review, most of the researchers have used Support Vector Classifier (SVC), Naïve Bayes, and Random Forest Machine Learning Algorithms for the evaluation. Therefore, four algorithms were used in our proposed methodology. We have utilized the Support Vector Machine (SVM) with RBF kernel and degree 2, Naïve Bayes, Decision Tree, and Random Forest with a total no. of 10 trees and 0 maximum depth. As we have a class imbalance problem, which can lead to biasness of outcomes by always predicting the incidental class accurately. To solve this problem, we have used the SMOTE filter [28] in python. This solved the class imbalance problem by equalizing the number of classes. We have macro averaged the results of precision, recall, and F1- score.

## IV. EXPERIMENTS

The experiments performed are divided into three levels. The first level of the experiment describes data preprocessing. The second level uses VerbNet Ontology to extract and map verbs from citation sentences on its class ID. The third level is to apply machine learning algorithms to classify a citation in three sentiment classes.

### A. Data Preprocessing

To preprocess both datasets, a Python application was developed to remove punctuations, tokenize and remove stop-words, and lemmatization. The application using Spacy and NLTK Libraries. The resultant is a set of tokens in their base format. The sample output is shown in Fig. 3.

### B. Extract Verbs

The second experiment was to extract verbs from the preprocessed citing sentences. This step was achieved using Part of Speech (POS) tagger using NLTK using algorithms from similar research [13]. The total number of unique verbs extracted from the AAL dataset was 555, and from the H-Index dataset were 337. The total occurrence of verbs in the AAL dataset was 18,789 and, in the H-Index dataset were 700. Later, these unique verbs were assigned IDs using VerbNet class IDs. Kindly note, a verb can be a part of multiple classes in VerbNet making it a graph rather than a tree. Table I shows some sample verbs and their assigned VerbNet Class ID.

### C. Machine Learning Models

We have utilized the Support Vector Machine (SVM) with RBF kernel and degree 2, Naïve Bayes, Decision Tree, and Random Forest with a total no. of 10 trees and 0 maximum depth. As both datasets have a Class Imbalance problem that can lead to biases of outcomes by always predicting the incidental class accurately, SMOTE filter [28] was used. This solved the class imbalance problem by equalizing the number of classes. For the evaluation of results, macro averaged results were tabulated for precision, recall, and F1- score.

```

0 [analyze, set, article, identify, major, operation, edit, extract, sentence, include, remove, extraneous, stress, et
ract, sentence, combine, reduce, sentence, sentence, syntactic, transformation,...]
1 [table, 3, sample, compression, compression, angle, rat, baseline, 970, 100, dt2step, 2206, 221, apade, 1009, 210,
humor, 2007, 383, table, 4, wear, rating, automatic, compression, rally, add, e...
2 [53, relate, work, discussion, trusted, model,
essentially, belong, category, work, mail, et, al. . . 1909, jing, mckinnon, 2006]
3 [1999, propose, summarization, system, base, draft, revision, jing, mckinnon, 2008, process, system, base, extraction,
cutandpaste, generation, abstracter, perform, cutandpaste, operation, jing, m...
4 [find, deletion, lead, part, occur, summary, unlike, case, jing, mckinnon, 2009]
5 [automatic, text, summarization, approach, offer, reasonably, wellperform, approximation, identify, important, sente
nce, lin, buy, 2002, schiffman, et, al, 2002, erkan, rahy, 2004, mihalcea, te...
6 [al, 1994, compression, sentence, automatic, translation, approach, knight, narco, 2008, hidden, markov, model, jing,
mckinnon, 2008, topic, signature, base, method, lin, buy, 2000, lacatuso, et, ...
7 [generally, accept, kind, postprocessing, perform, improve, final, result, shortening, fuse, revise, material, grefes
statte, 1958, mail, gate, blowdown, 1990, jing, mckinnon, 2009, harzily, et, al...

```

Fig. 3. Citation Texts after Preprocessing.

TABLE I. EXTRACTED VERBS VS VERBNET CLASS IDS

| No | Verbs    | Class ID's   |
|----|----------|--|
| 1  | Analyze  | ['assessment-34']  |
| 2  | Set      | ['braid-41.2.2', 'force-59-1', 'image_impression-25.1', 'preparing-26.3-2', 'put-9.1-2'] |
| 3  | Identify | ['characterize-29.2-1-1']  |
| 4  | Remove   | ['banish-10.2', 'remove-10.1']   |
| 5  | Combine  | ['mix-22.1-1-1']   |
| 6  | Reduce   | ['limit-76']   |
| 7  | Extract  | ['remove-10.1']  |
| 9  | Add      | ['mix-22.1-2']   |
| 10 | Relate   | ['say-37.7-1']   |
| 12 | Perform  | ['performance-26.7-1']   |
| 13 | Think    | ['consider-29.9-2', 'wish-62']   |
| 14 | Find     | ['declare-29.4-1-1-2', 'get-13.5.1']   |
| 15 | Lead     | ['accompany-51.7', 'force-59']   |
| 16 | Occur    | ['occurrence-48.3']  |
| 17 | Offer    | ['future_having-13.3', 'reflexive_appearance-48.1.2']                                    |
| 18 | Accept   | ['approve-77', 'characterize-29.2-1-1', 'obtain-13.5.2']                                 |
| 19 | Improve  | ['other_cos-45.4']   |
| 20 | Generate | ['engender-27']  |
| 21 | Neglect  | ['neglect-75-1-1']   |
| 22 | Produce  | ['create-26.4', 'performance-26.7-2']  |
| 23 | Observe  | ['conjecture-29.5-2', 'investigate-35.4', 'say-37.7-1', 'sight-30.2']                    |

D. Performance Evaluation

1) *Classification accuracy*: Classification Accuracy is the ratio of the number of correct predictions to the total number of input examples. Classification accuracy is calculated by using the formula shown in Eq. 1.

$$Accuracy = \frac{No\ of\ Correct\ Predictions}{Total\ No\ of\ Predictions\ Made} \quad (1)$$

2) *Precision (Positive Predictive Value)*: Precision is a metric that counts the number of correct positive predictions made by the algorithm. It was calculated using the formula shown in Eq. 2.

$$Precision = \frac{True\ Positive}{True\ Positive + False\ Positive} \quad (2)$$

3) *Recall*: Recall is the metric that counts the number of correct positive predictions made from all positive predictions. It was calculated using the formula in Eq. 3.

$$Recall = \frac{True\ Positive}{True\ Positive + False\ Negative} \quad (3)$$

4) *F-Score*: F-measure combines both recall and precision into a single measure that has both the properties. Alone, neither recall nor precision expresses the complete story. So, once recall and precision have been calculated, both scores were combined into the calculation of F-measure. It is calculated by using the formula in Eq. 4.

$$F - Score = \frac{2 * Precision * Recall}{Precision + Recall} \quad (4)$$

V. RESULTS

After all pre-processing is applied on both datasets, the datasets are passed to the model for training. To evaluate the performance of our algorithms, the datasets were divided into two sets (Training and Validation set). 70% of the labeled data was used for training and 30 % of the labeled data was set aside for validation. After the training phase, 30% of the data was used to find out the accuracy of the algorithm. This labeled data was passed to the trained model. The model assigned labels to the verbs. These labels were then compared with the actual labels of the data. This comparison showed that our model was able to label all the verbs with an accuracy of 90%.

A. Results on AAL Dataset

Fig. 4 shows performance evaluation on AAL Dataset for four different classifiers, whereas Fig. 5 shows the precision-recall curve. The results show that SVM and Random Forest have performed better than Decision Tree and Naïve Bayes.

B. Results on H-Index Dataset

Fig. 6 shows performance evaluation on H-Index Dataset for four different classifiers, whereas Fig. 7 shows the precision-recall curve. The results show that SVM and Decision Tree have performed better than the other two.

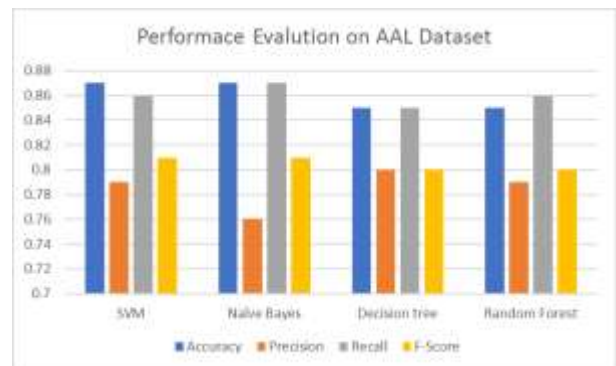


Fig. 4. Results on AAL Dataset.

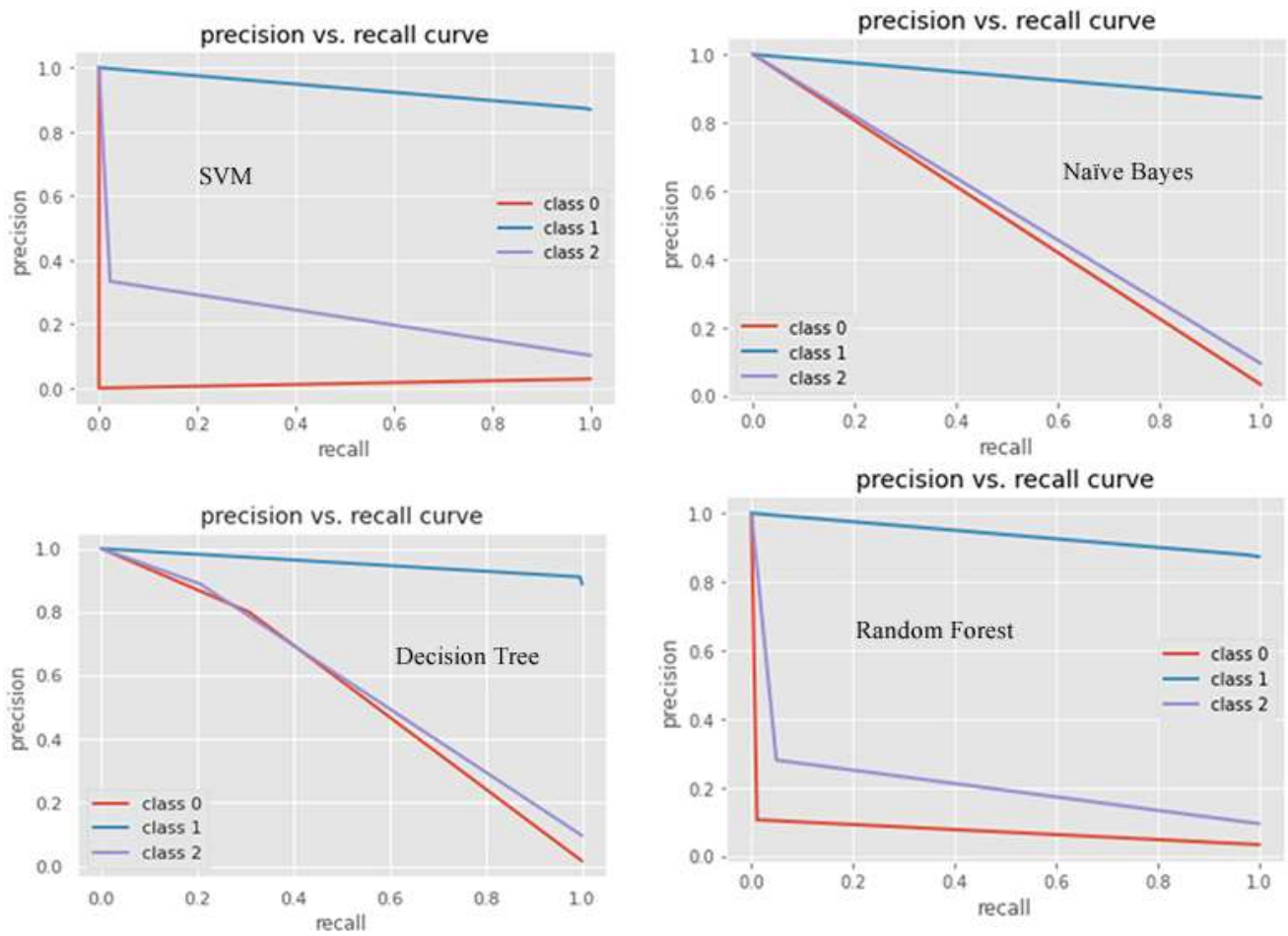


Fig. 5. Precision vs Recall Curve on 4 MLAs for AAL Dataset.

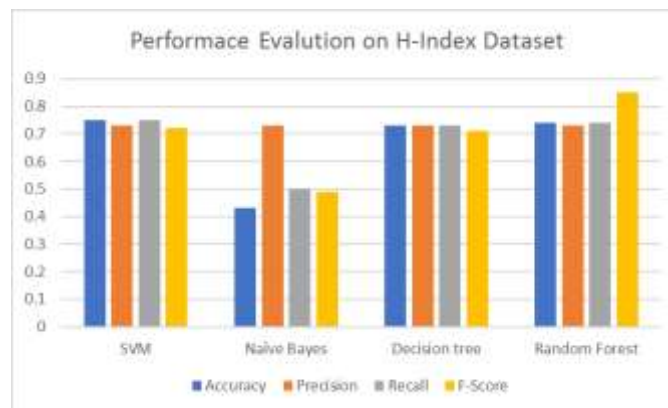


Fig. 6. Results on H-Index Dataset.

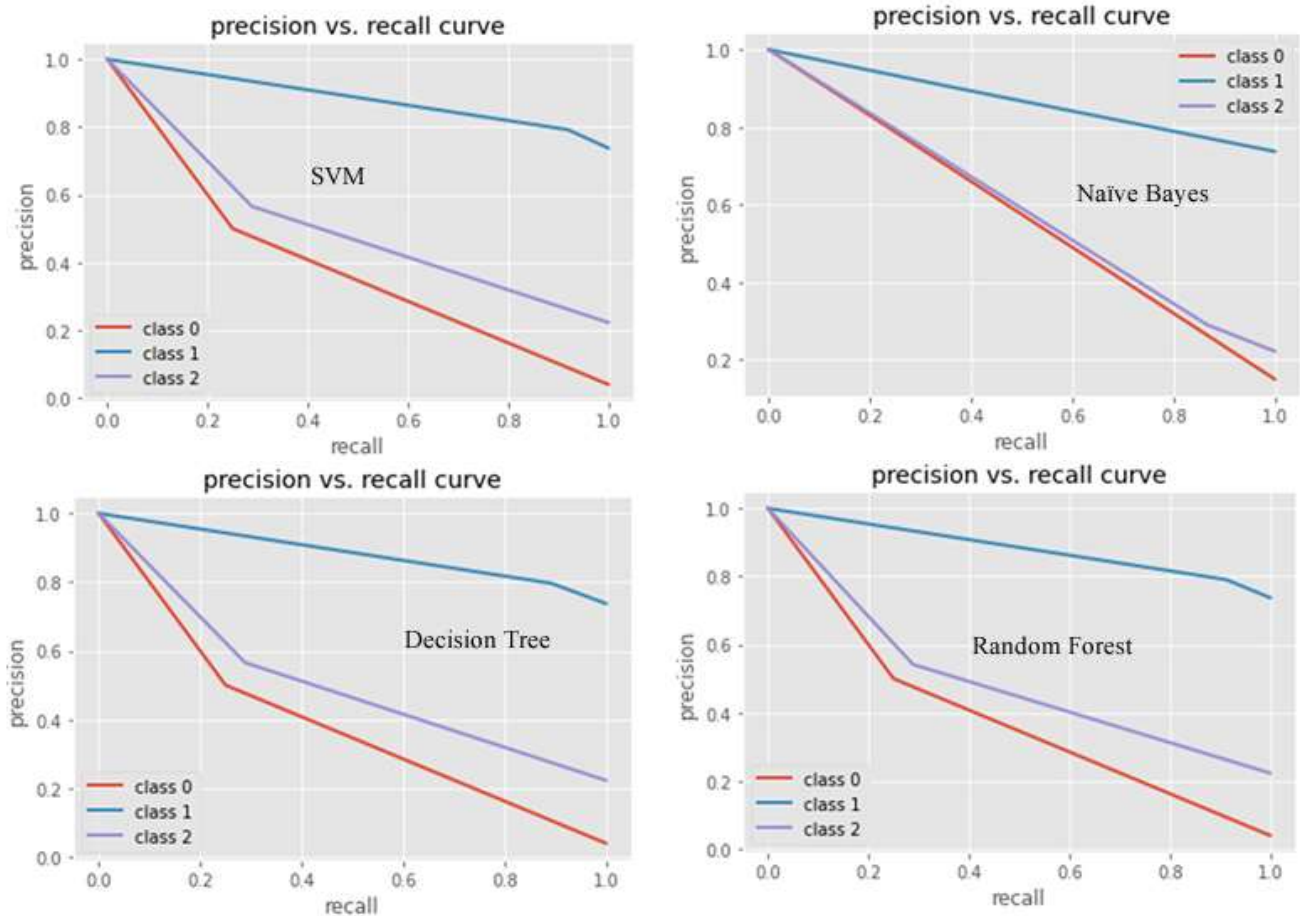


Fig. 7. Precision vs Recall Curve on 4 MLAs for H-Index Dataset.

### C. Combined Results

Combined results show that SVM has given better results than Naïve Bayes, whereas Random Forest has given the best results for both datasets implying that the extracted verbs as features have shown promising results using Support Vector Classifier and Random Forest as compare to Naïve Bayes.

## VI. CONCLUSION

Research is a continuous and recursive process. Every research paper and articles are built on some prior knowledge in the field. Research papers include citations to the external resources to discuss the work done by the previous researcher. With the rapid development in the research area, it becomes challenging for researchers to recognize quality research work. We have explored various existing approaches where classification methods mostly use nouns, adjectives, etc. as features. This paper proposes a new verb-based approach as an important term of opinion. We have extracted opinion structures that regard the verb as an essential component. We have used publicly available ACL Anthology Citation Dataset and our curated H-Index dataset for experiments. The experiments show 90% accuracy using Random Forests.

### REFERENCES

[1] M. El-Din, "Analyzing Scientific Papers Based on Sentiment Analysis," Cairo University, 2016.

[2] B. Liu, "Sentiment analysis and opinion mining," *Synth. Lect. Hum. Lang. Technol.*, 2012, doi: 10.2200/S00416ED1V01Y201204HLT016.

[3] Y. Zhang, R. Jin, and Z. H. Zhou, "Understanding bag-of-words model: A statistical framework," *Int. J. Mach. Learn. Cybern.*, 2010, doi: 10.1007/s13042-010-0001-0.

[4] C. N. dos Santos and R. L. Milidiú, "Part-of-speech tagging," in *SpringerBriefs in Computer Science*, 2012.

[5] F. Benamara, C. Cesarano, A. Picariello, D. Reforgiato, and V. S. Subrahmanian, "Sentiment analysis: Adjectives and adverbs are better than adjectives alone," 2007.

[6] S. B. Shum, "Evolving the Web for Scientific Knowledge: First Steps Towards an ÖHCI Knowledge WebÖ TodayÖs HCI digital library," *Interfaces (Providence)*, vol. 39, pp. 1–9, 1998.

[7] F. Qayyum and M. T. Afzal, "Identification of important citations by exploiting research articles' metadata and cue-terms from content," *Scientometrics*, 2019, doi: 10.1007/s11192-018-2961-x.

[8] D. W. O. Fleischmann, , An-Shou Cheng and Kenneth R., Ping Wang,Emi Ishita, "The Role of Innovation and Wealth in the Net Neutrality Debate," *Commun. Inf. Lit.*, 2009, doi: 10.1002/asi.

[9] I. Ihsan and M. A. Qadir, "CCRO: Citation's Context Reasons Ontology," *IEEE Access*, vol. 7, pp. 30423–30436, 2019, doi: 10.1109/ACCESS.2019.2903450.

[10] A. Athar, "Sentiment analysis of citations using sentence structure-based features," *ACL HLT 2011 - 49th Annu. Meet. Assoc. Comput. Linguist. Hum. Lang. Technol. Proc. Student Sess.*, no. June, pp. 81–87, 2011, [Online]. Available: <http://dl.acm.org/citation.cfm?id=2000976.2000991>.

[11] C. Bazerman, *The Informed Reader: Contemporary Issues in the Disciplines*. 2010.

- [12] G. Thompson and Y. Yiyun, "Evaluation in the reporting verbs used in academic papers," *Appl. Linguist.*, vol. 12, no. 4, pp. 365–382, 1991, doi: 10.1093/applin/12.4.365.
- [13] I. Ihsan, S. Imran, O. Ahmed, and M. A. Qadir, "Sentiment Based Study of Citations Reporting Verb Corpus Using Natural Language Processing," *Corporum J. Corpus Linguist.*, vol. 2, no. 1, pp. 25–35, 2019.
- [14] A. Athar and S. Teufel, "Context-enhanced citation sentiment detection," 2012.
- [15] V. Qazvinian and D. R. Radev, "Identifying non-explicit citing sentences for citation-based summarization," 2010.
- [16] N. Tandon and A. Jain, "Citation Context Sentiment Analysis for Structured Summarization of Research Papers," in the 35th German Conference on Artificial Intelligence (KI-2012), 2012, no. i, pp. 98–102.
- [17] C. Jochim and H. Schütze, "Towards a generic and flexible citation classifier based on a faceted classification scheme," 24th Int. Conf. Comput. Linguist. - Proc. COLING 2012 Tech. Pap., no. December 2012, pp. 1343–1358, 2012, [Online]. Available: <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.379.2126>.
- [18] H. F. Eid, A. Darwish, A. Ella Hassanien, and A. Abraham, "Principle components analysis and support vector machine based Intrusion Detection System," in Proceedings of the 2010 10th International Conference on Intelligent Systems Design and Applications, ISDA'10, 2010, pp. 363–367, doi: 10.1109/ISDA.2010.5687239.
- [19] B. H. Butt, M. Rafi, A. Jamal, R. S. Ur Rehman, S. M. Z. Alam, and M. B. Alam, "Classification of research citations (CRC)," in CEUR Workshop Proceedings, 2015, vol. 1384, no. January, pp. 18–27.
- [20] I. C. Kim and G. R. Thoma, "Automated classification of author's sentiments in citation using machine learning techniques: A preliminary study," 2015, doi: 10.1109/CIBCB.2015.7300319.
- [21] J. Xu, Y. Zhang, Y. Wu, J. Wang, X. Dong, and H. Xu, "Citation Sentiment Analysis in Clinical Trial Papers," *AMIA ... Annu. Symp. proceedings. AMIA Symp.*, vol. 2015, pp. 1334–1341, 2015.
- [22] M. Valenzuela, V. Ha, and O. Etzioni, "Identifying meaningful citations," *AAAI Work. - Tech. Rep.*, vol. WS-15-13, pp. 21–26, 2015, [Online]. Available: <http://ai2-website.s3.amazonaws.com/publications/ValenzuelaHaMeaningfulCitations.pdf>.
- [23] Z. Taşkın and U. Al, "A content-based citation analysis study based on text categorization," *Scientometrics*, vol. 114, no. 1, pp. 335–357, 2018, doi: 10.1007/s11192-017-2560-2.
- [24] A. M. Giuglea and A. Moschitti, "Semantic role labeling via FrameNet, VerbNet and PropBank," 2006, doi: 10.3115/1220175.1220292.
- [25] B. Levin, *Book reviews*, vol. 37, no. 3. University of Chicago Press, 2012.
- [26] D. R. Radev, P. Muthukrishnan, V. Qazvinian, and A. Abu-Jbara, "The ACL anthology network corpus," *Lang. Resour. Eval.*, vol. 47, no. 4, pp. 919–944, 2013, doi: 10.1007/s10579-012-9211-2.
- [27] I. Ihsan and M. A. Qadir, "CCRO: H-Index Dataset," 2020. <https://github.com/imranihsan/CCRO/blob/master/CCRO-CC2.csv>.
- [28] J. A. Sáez, J. Luengo, J. Stefanowski, and F. Herrera, "SMOTE-IPF: Addressing the noisy and borderline examples problem in imbalanced classification by a re-sampling method with filtering," *Inf. Sci. (Ny)*, 2015, doi: 10.1016/j.ins.2014.08.051.

# DBSR: A Depth-Based Secure Routing Protocol for Underwater Sensor Networks

Ayman Alharbi

Department of Computer Engineering  
Collage of Computer Science and Information Systems  
Umm Al-Qura University, Mecca  
Saudi Arabia

**Abstract**—Depth-Based protocol has gained considerable attention as an efficient routing scheme for Underwater Wireless Sensor Networks UWSNs. It requires only depth information to perform the routing process. Despite this feature, UWSNs which operate with the employment of DBR protocol are vulnerable to depth spoofing attack. In this paper, Depth Based Secure Routing protocol is proposed to overcome this vulnerability. DBSR modifies traditional DBR routing algorithm by securing the depth information which is embedded in the header part of DBR packet. In addition to that, each node verifies the sender's identity based on a digital signature scheme. We extensively evaluate the overhead and performance gain of DBSR for two signature schemes based on Elliptic Curve Cryptography method considering various network conditions. The simulation study is performed using NS3-based simulator. Our results show that DBSR can avoid depth-spoofing attack by achieving 95% and 85% delivery ratios under low and high network loads respectively. Contrary to popular belief, results show that careful utilization of cryptographic techniques is justifiable without significant overhead on the communication cost.

**Keywords**—UWSN; DBR; ECC

## I. INTRODUCTION

Water covers more than 70% of earth planet. The nature of underwater world includes valuable resources such as unique minerals, various food sources, and other undiscovered sites. Traditionally, a diver or marine underwater vehicle collect data from fixed sensors which were used in order to observe underwater information. However, due to the harsh and unsafe environment of underwater world, scientists continue developing more tools to be used remotely [1], [2]. Moreover, this approach is not suitable for real-time applications such as military surveillances. As a result, Underwater Wireless Sensor Networks (UWSNs) have emerged as a promising technology due to their unique features [3] [4] [5]. First, underwater sensor networks provide useful sensing capabilities that can be used for long-term and short-term monitoring. They have the capability to be operated days, weeks, months even years wirelessly, hence, enable of wide sensing fields such as: temperature, salinity, current movements, video, image, chemical sensing [6][7]. Second, high density feature allows extensive discovering and exploration of wide underwater areas. Third, real-time sensing and monitoring missions can be achieved using underwater

sensor networks [8]. Fourth, when unexpected failure occurred in any sensor in the network, rapid error detection and remote fixing are applicable features using underwater sensor networks [9]. Fifth, compared to traditional underwater equipment, underwater sensor networks offer the possibility of re-configuring sensors remotely and eliminate the need for physically accessing underwater sites. In summary, UWSNs help to transmit data through wide distances and harsh circumstances. Figure 1 shows an example of UWSN architecture. In this architecture, each underwater node is capable of gathering, relaying data through different transmission media (acoustic or optical) waves to the surface gateway. After collecting data from underwater nodes, gateways nodes transmit the observed data to the base station using radio waves.

There has been considerable effort to enhance the performance of UWSN for different objectives e.g. delay, power, mobility and other performance goals [10], [11]. Depth based routing protocol [12] was proposed in order to enhance routing functionality of UWSNs. In this protocol, the forwarding procedure depends mainly on the depth information of the source of each received packet. In other word, the main advantage of DBR is that it free doesn't depend on complex geographic computation and perform a free localization method. In traditional DBR protocol, the routing mechanism is based on depth information of each forwarder node. When the node receive a packet, it checks the forwarder/sender depth then check if it is candidate for forwarding the received packet. If the sender's depth is larger than its depth, it hold the packet for a certain amount of time called "holding time". If this condition is not met, the packet will immediately be dropped. However, this approach is vulnerable to serious security attack namely, depth-spoofing attack [13]. In this attack, sender's depth could be compromised and utilized for malicious activities. If an adversary succeeded in gaining information about current topology, it can easily deploy a malicious node at an excellent location where it can receive packets form different nodes in the network. In this case, the attacker will forward that packet with a fake depth projecting a better position. Accordingly, any node located at the attacker's transmission range will drop their packets after receiving attacker's spoofed message.

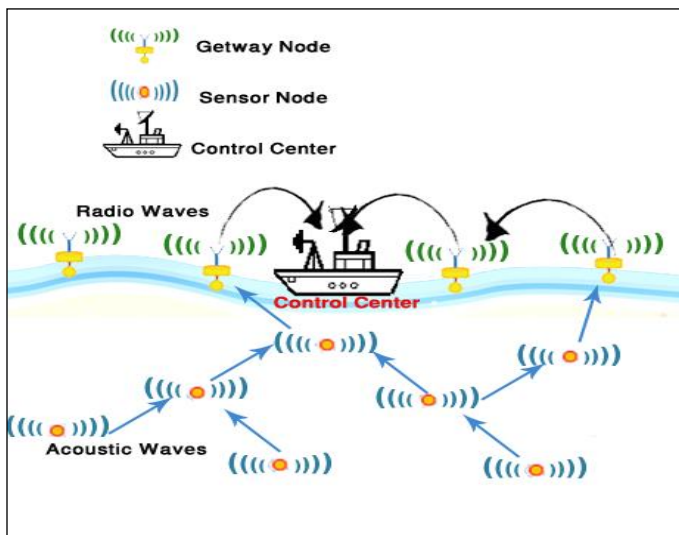


Fig. 1. UWSN Architecture.

In this paper, we propose Depth Based Secure Routing protocol to mitigate the abovementioned vulnerability. We seek to study the performance analysis of securing DBR protocol based on cryptography approach extensively. The performance analysis considers state-of-the-art encryption schemes designated energy-constrained devices.

The rest of the paper is organized as follows: Section II points out relevant background and details the attack model. In Section III, related work is reviewed. Section V presents the proposed solution. Section VI highlights the main findings based on simulation results. Finally, Section VII summarize author's conclusions.

## II. BACKGROUND

Due to the unique challenges of aquatic environment, proactive or reactive general routing protocols do not behave well underwater and considered very costly. Hence, Geographic routing protocols provides better performance and most suitable for UWSN's.

### A. Geographic Information-based Routing

Geographic routing protocols depend on location information of sensor nodes. The key factor of establishing routing paths between the source and destination is the location of each node. VBF [14] and DBR [12] are both among the most common geographic information-based routing protocols which are built for UWSNs. In this subsection, we review the main functionality of the two protocols. However, the main focus will be on DBR to highlight recent enhancements to this protocol, and its existing vulnerability named as the depth spoofing attack.

1) *Vector Based Forwarding (VBF) Routing:* In VBF, a path from source to destination will be limited to only few number of nodes which satisfy certain geographic conditions. Therefore, the vector allow will only high benefit nodes to be participated in the forward process. Other nodes will discard the packets to save energy. As it mentioned, participation in routing depends on which nodes fall in the path between

source and destination. The packet compose of three fields, the sender (A), the sink (B), and the forwarder. The algorithm will find the routing vector from the sender to the sink A to B, then forward packets along the path. Each node belong to the path can forward the packets based on calculating a specific factor. This factor is called the "desirableness factor" which determines the suitability of each candidate forwarding node. As a result, the node will discard if the calculated factor is large. Accordingly, if the results is 0, it will be optimal node for forwarding the packet.

2) *Depth Based Routing Protocol DBR:* The basic idea behind DBR is that, a node need only the recognize sender's depth to decide whether it is eligible to forward the received packet or drop it. Hence, only optimal forwarder nodes will be considered among the routing process. Therefore, one of the main advantages of DBR is that any node inside the network doesn't need to have any information about the current topology or further locations information of other nodes.

Among the preparation process of the transmission of a generated packet, an important step which the source node incorporates its depth in the header part. While all nodes which are located at the same transmission range of the sender will receive the packet, a packet will be dropped if the sender's depth  $d_s$  is lower than the receiving node's depth  $d_r$ . Consequently, if  $d_s > d_r$  receiving node will be candidate to forward the packet.

It is worth mentioning that candidate node will keep the received packet for a certain period of time called "holding time ( $T$ )". This factor is used for the advantage of calculating the closest node to the surface which will be the qualified forwarder for the packet. The holding time can be calculated as :

$$T = K(Y - \delta) \quad (1)$$

where  $Y$  is the communication range for the node and  $\delta$  is the difference between  $d_s$  and  $d_r$ .  $K$  is a constant which can be used to determine the maximum holding time.

a) *DBR enhancement:* Energy-Efficient Depth-Based Routing EEDBR protocol was proposed by [15]. The key difference between DBR and EEDBR is the decision of selecting the forward node. In EEDBR, the decision will be based on both the depth and the residual energy of the forwarder. The protocols requires that each node will broadcast its residual energy and depth to its neighbors frequently. Unfortunately, this drawback add more overhead to the routing procedure since more packets are required.

Light-weight depth-based routing LDBR was also proposed to enhance DBR for underwater wireless sensor network [16]. As in EEDBR, the residual energy is also taken into consideration when candidate nodes receive packets from the sender. However, the enhancement was made to the decision of determining the optimal forwarder. In addition of depth condition, the residual energy of sender and next hop packet will be one of the main factors to determine the optimal forwarder node.



An Improved Adaptive Mobility of Courier Nodes in Threshold-Optimized BDR Protocol IAMCTD was proposed [17]. In this proposed approach, the authors designed an improved DBR protocol to deal with real-time sensitive applications. In addition of depth, other network parameters also considered for routing decision such as network density. The protocol improved the network lifetime as well as transmission loss.

An Optimized Depth-Based Routing Protocol for Underwater Wireless Sensor Networks ODBR was also proposed [18]. The protocol address a shortcoming point which exists in DBR protocol. The nodes closest to the sink will lose energy more than other network nodes. Therefore, the proposed algorithm ensured an optimized method for energy balancing between all nodes. In ODBR, nodes with high traffic will be marked for a specific zone in order to reduce the energy consumption. Hence, ODBR improves lifetime, throughput and energy consumption of UWSNs.

### III. RELATED WORK

Authors in [19] presented a study for evaluating the cost of digital signature schemes. They highlighted the usability of applying digital signature schemes for the environment of UWSNs. However, the study only considered the cost of signing while the verification cost is assumed to be performed by the sink only. Moreover, the routing protocol is not specified. In [20], authors mitigated the effect of depth-spoofing attack by combining between authentication-based method and thresholding strategy. Unlike DBR protocol, node will compare its depth with lower/upper bounds threshold. Hence, if the receiver's depth falls within range of threshold, it forwards the packet. This window is calculated randomly by the sender depending on its neighbors information. However, the proposed approach suffered high transmission cost since each source node will depend on a randomized threshold window which may rise the cost significantly.

### IV. ATTACK MODEL

We adapt the attack model presented in [6] to validate our proposed solution. In this attack model, the attacker first will try to eavesdrop the transmission and listen to at least two nodes in the area. Figure 2 illustrates the attack strategy. The source S is in the same range of first forwarder node f1. The attacker is also able to hear transmissions generated by the source because it is located in the same transmission range. When the source transmit packet, f1 should compare its depth with the source's one. F1 will decide to hold the packet for a certain threshold since its depth is the lower than the source. The packet should be forwarded by f1 to f2 by the end of holding time, however, since the attacking node received the same packet, it forward it with spoofed depth assuming lower value than f1. Accordingly, f1 will drop the packet which has been delayed due to the holding time period. Consequently, all other nodes in the attacker transmission range will receive the same packet with smaller depth and then will drop their packet. As a result, the attacker is network.

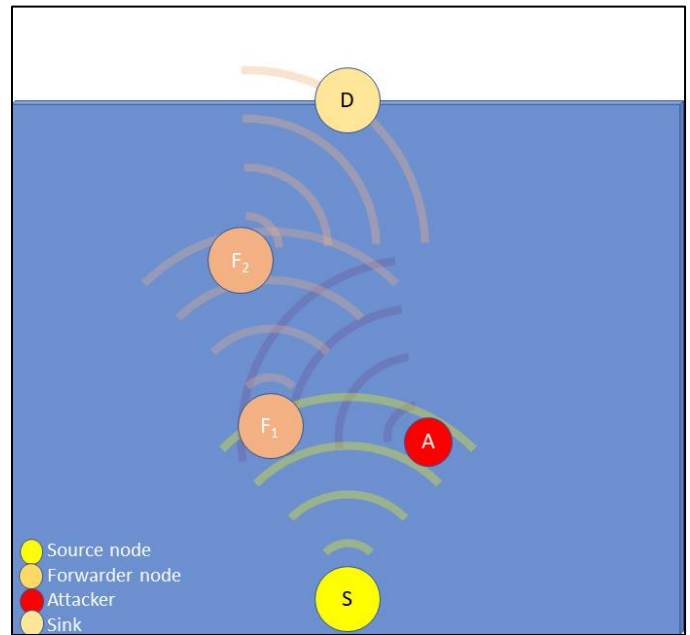


Fig. 2. Depth-Spoofing Attack.

### V. PROPOSED METHOD

#### A. Justifications

- Our proposed method is based on light cryptography techniques using elliptic curve algorithm which utilize fewer bits to secure transmission.
- The proposed method secure DBR routing protocol one of the most widely used in the are of underwater sensor networks.
- There is lack of extensive performance analysis of utilizing cryptography techniques to secure DBR protocol.

#### B. Assumptions

- The attacker and the legitimate nodes have the same transmission range.
- The attacker has better capability so that he/she will choose the best location in the network based on many factors such as (depth, weak links .. etc.).
- We assume homogeneous nodes i.e. all nodes have the same level of energy.
- We ignore the computation cost since it is negligible to communication cost for UWSNs [20].
- All public/private keys inside each node are secured by strong hardware-rooted encryption platform that makes it infeasible to compromise a node.
- We assumes that no additional nodes will be added to the network. Therefore, we leave the scalability issue for future work.

Able to prevent further transmissions and crash the DBSR methodology.

To mitigate the effect of depth spoofing attack we need to protect that sensitive information on which the best forwarder for the packet is selected. As depicted in Figure 3 DBR protocol encapsulates the sender's depth with each packet as a part of the header. Unquestionably, this tiny information is very sensitive and it is highly vulnerable to the depth-spoofing attack. Our proposed approach add the following additional security steps to the existing protocol:

- 1) Pair of security keys (private/public) will be assigned to each node before deployment. The private key will be used for signing whereas the public key used for verification.
- 2) Each sensor node will be configured with its own private key and a list of all public keys of other node.
- 3) Pair of security keys (private/public) will be assigned to each node before deployment. The private key will be used for signing whereas the public key used for verification.
- 4) Each sensor node will be configured with its own private key and a list of all public keys of other nodes.
- 5) The DBSR packet header, shown in shown in Figure 3 contains two additional fields: The recent forwarder ID and the forwarder signature.
- 6) The sender/forwarder calculates the signature of: Packet ID (source ID , sequence number), forwarder ID, and forwarder depth.
- 7) The sender/forwarder places its signature and then transmits the whole packet to the next hop.
- 8) A receiver first verifies the signature. If the verification fails, the packet is considered malicious and it will be ignored.

Figure 4 summarizes the steps of the DBSR protocol. In the deployment stage, each node will be configured with pair of public/private key. In addition, each node stores list of all public keys of other nodes. In the operation stage, each node will verify sender's depth before accepting the received packet.

### C. Proposed Signature Scheme

UWSNs have unique characteristics which restrict the available resources for cryptography operations. Due to the limitation of available energy at underwater wireless nodes, ECC-based schemes are considered more suitable for wireless devices [19]. Therefore, we consider two efficient ECC-based algorithms for DBSR design, ECDSA and BLS digital signatures methods. As can be depicted from algorithm 1, the signing process leads to additional bits, therefore, the overhead will be investigated per hop considering various network conditions.

## VI. PERFORMANCE ANALYSIS

### A. Experiment Settings

To evaluate the overhead of authentication on the network, we set up three scenarios as follow: First, we highlight the behavior of DBR under different network conditions and various message lengths. Second, we investigate the effect of depth-spoofing attack on DBR. Third, we study the overhead of DBSR considering two encryption schemes.

Throughout the experiment, to study the effect of network capacity, we run the simulation for different generation bit rate of source nodes, namely 10 b/s, 50 b/s and 100 b/s for light, medium and high loads, respectively. Also, for effect of message length, we interchange the value of message length between 100b and 1kb. Moreover, we investigate the effect of attacker capability by varying number of malicious nodes between 1 and 7 nodes. Finally, for the effect of different signature schemes, we interchange the overhead between 40 bits and 20 bits for ECDSA and BLS respectively. Other simulations settings can be summarized in Table I.

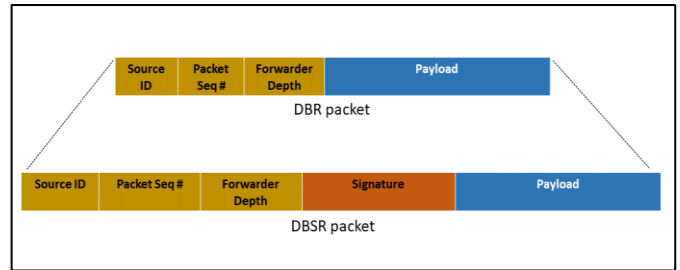


Fig. 3. DBR and DBSR Packet Format.

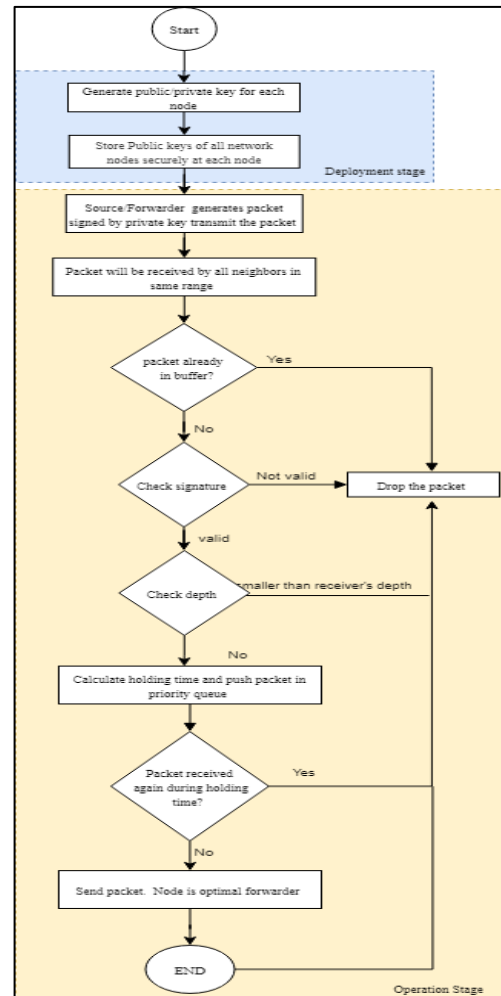


Fig. 4. DBSR Flowchart.

**Algorithm 1** DBSR sender signature generation based on ECDSA

1. **Compute** the header hash  $\mathbf{h} = \text{hash}(\text{header})$
2. **Compute** a random number  $\mathbf{k}$  where  $1 \leq \mathbf{k} \leq (p - 1)$
3. **Compute** random point  $(\mathbf{x1}, \mathbf{y1}) = \mathbf{k} \times \text{Base point}(x,y) \text{ mod } \mathbf{p}$   
 $\mathbf{r} = \mathbf{x1} \text{ mod } \mathbf{n}$
4. **Check**  $\mathbf{r} \neq 0$ , if yes then **repeat** steps 1 to 3
5. **Compute** signature  $\mathbf{S} = (\mathbf{k}^{-1} (\mathbf{h} + \text{sender's private key} \times \mathbf{r})) \text{ mod } \mathbf{n}$

**Algorithm 2** DBSR signature verification by receiver based on ECDSA

1. **Compute**  $\mathbf{w} = \mathbf{S}^{-1} \text{ mod } \mathbf{n}$
2. **Compute**  $\mathbf{u1} = (\mathbf{h} \times \mathbf{w}) \text{ mod } \mathbf{n}$   
 $\mathbf{u2} = (\mathbf{r} \times \mathbf{w}) \text{ mod } \mathbf{n}$
3. **Compute**  $\mathbf{C}(x2,y2) = \mathbf{u1} \times \text{Base point}(x,y) + \mathbf{u2} \times \text{sender's public key}$
4. **Check**  $\mathbf{x2} \text{ (mod } \mathbf{p}) \neq \mathbf{r}$ , if yes **reject and report to sink**

TABLE I. SIMULATION PARAMETRES

| Parameter              | Value        |
|------------------------|--------------|
| Network area           | 500m × 500 m |
| Network Density        | 3            |
| Number of source nodes | 7            |
| Communication Ranges   | 150 m        |
| Interference Ranges    | 300 m        |
| Interference Range     | 300 m        |
| Total number of nodes  | 30 nodes     |
| Channel Bit Rate       | 10000 b/s    |
| ECDSA signature size   | 40 bits      |
| BLS signature size     | 20 bits      |

**B. DBSR vs DBR**

To show the effectiveness and benefits of DBSR, we compare the delivery ratio and packet loss percentage of DBR protocol with DBSR protocol under low, medium and high traffic loads.

As can be observed from Figure 5 and Figure 6, DBR will be significantly affected by active attack especially when the networks operates under low traffic. The number of delivered packets decreases from 600 to 50 packets. On the other hand, DBSR improved the performance by achieving 585 successfully delivered packets.

**C. Effect of Attacker Capability**

Unquestionably, when the number of attacking nodes increases, the chance of attack effect increases. As can be seen in Figure 7 and Figure 8, the effect of the number of attackers on the delivery ratio is dominant. The more attacked nodes, the lower the delivery ratio. In the worst case, the attackers can reduce the delivery ratio by 91%. Similarly, the more

attacked nodes, the higher the packet loss ratio. In the worst case, the attackers can increase the packet loss ratio by 90%. As can be observed also, the effect of the attack is more severe in lightly loaded networks. In addition, it is worth mentioning that when the network operates under no attack, the main factor affecting delivery ratio is the network load. Delivery ratio in a highly loaded network can be as low as 36% as can be seen in Figure 8.

**D. Effect of Message Length**

As mentioned previously, we evaluate the effect of different message lengths by considering small a long values for each generated message. As can be depicted by Figure 8, the effect of message length is negligible especially in lightly loaded networks (<2% change in delivery ratio). This is encouraging because it indicates that adding authentication can be justified in lightly loaded network.

**E. Effect of Authentication**

As can be seen in Figure 9, the overhead of authentication is smaller for long messages. Same results show that overhead of BLS authentication is smaller than ECDSA. This is due to the fact that BLS generates shorter signatures than ECDSA. The difference between BLS and ECDSA is more significant for the case of short messages in a highly loaded networks. This is due to the fact the overhead of authentication almost doubles the network load which leads to excessive packet loss due to congestion.

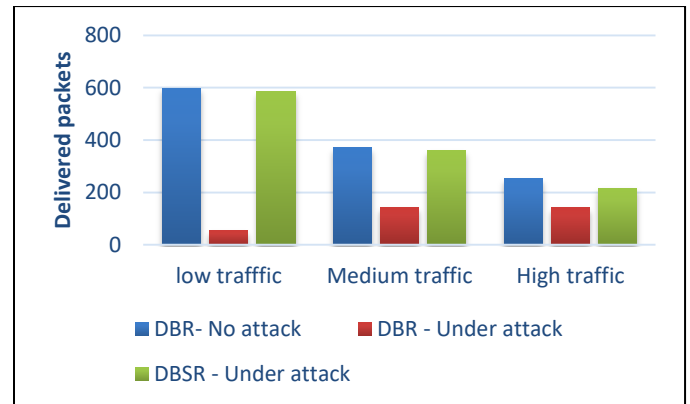


Fig. 5. DBR vs DBSR Delivered Packets.

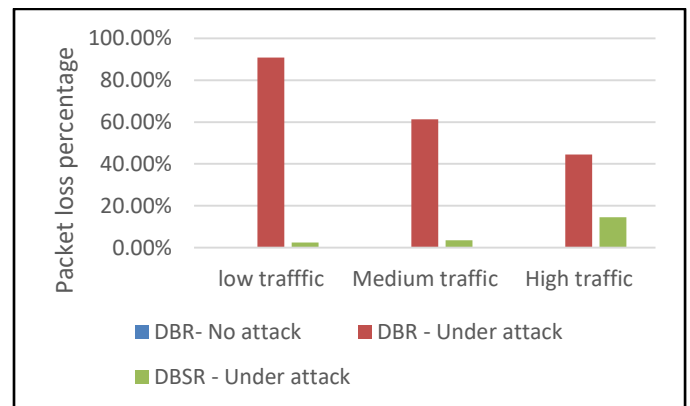


Fig. 6. DBR VS DBSR Packet Loss.

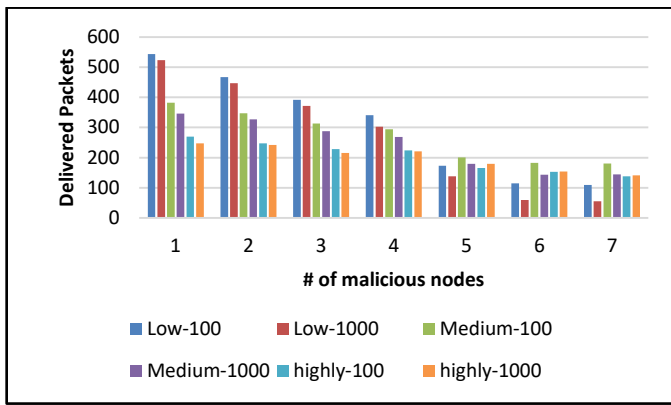


Fig. 7. Effect of Attacker Capability.

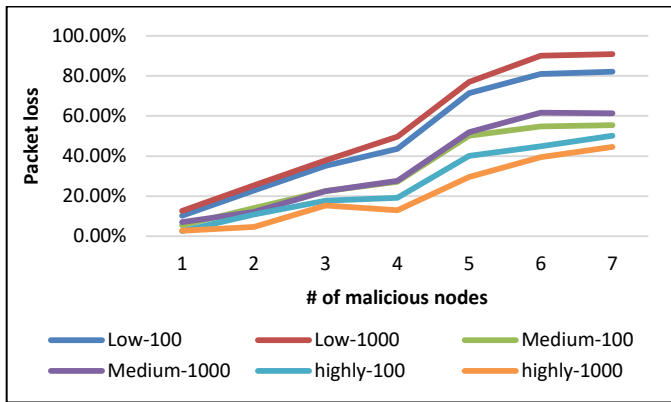


Fig. 8. Effect of Message Lengths.

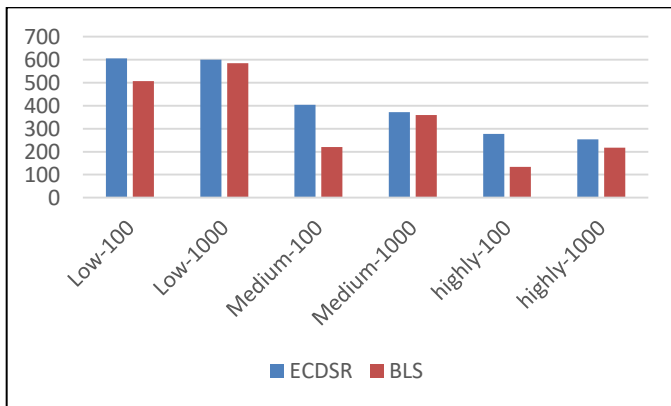


Fig. 9. ECDSA Vs BLS.

## VII. CONCLUSION AND FUTRE WORK

In this paper, a security improvement to DBR routing protocol based on ECC scheme was introduced. The proposed approach eliminates depth-spoofing attack by securing the depth information. The proposed method suggests the use of an encryption mechanism, namely ECC algorithm. The key contribution is the study of signature overhead considering various network parameters. Simulation results show that the proposed scheme archives high delivery ratio. Results also show that there is high dependency between certain network parameters such as: network load and packet nominal length on one hand and signature overhead on the other hand. Our

future work to extend the proposed approach to allow key exchange an interesting question that arose from this research is how to enhance this approach so a new node can be added to the network without need for mutual authentication between nodes.

## REFERENCES

- J. Watt, M. R. Phillips, C. E. A. Campbell, I. Wells, and S. Hole, "Wireless Sensor Networks for monitoring underwater sediment transport," *Science of the Total Environment*, vol. 667, Elsevier B.V., pp. 160–165, Jun. 01, 2019, doi: 10.1016/j.scitotenv.2019.02.369.
- G. Yang, L. Dai, G. Si, S. Wang, and S. Wang, "Challenges and Security Issues in Underwater Wireless Sensor Networks," in *Procedia Computer Science*, 2019, vol. 147, pp. 210–216, doi: 10.1016/j.procs.2019.01.225.
- I. F. Akyildiz, D. Pompili, and T. Melodia, "Underwater acoustic sensor networks: Research challenges," *Ad Hoc Networks*, vol. 3, no. 3, pp. 257–279, May 2005, doi: 10.1016/j.adhoc.2005.01.004.
- A. A. Sheikh, E. Felemban, M. Felemban, and S. B. Qaisar, "Challenges and opportunities for underwater sensor networks," in *Proceedings of the 2016 12th International Conference on Innovations in Information Technology*, IIT 2016, Mar. 2017, doi: 10.1109/INNOVATIONS.2016.7880021.
- H. Alhomyani, R. Ammar, H. Albarakati, and A. Alharbi, "Deployment strategies for underwater sensing and processing networks," in *Proceedings - IEEE Symposium on Computers and Communications*, Aug. 2016, vol. 2016-August, pp. 358–363, doi: 10.1109/ISCC.2016.7543766.
- R. W. L. Coutinho, A. Boukerche, L. F. M. Vieira, and A. A. F. Loureiro, "Underwater Sensor Networks for Smart Disaster Management," *IEEE Consumer Electronics Magazine*, vol. 9, no. 2, pp. 107–114, Mar. 2020, doi: 10.1109/MCE.2019.2953686.
- M. Juhari, K. Ibrahim, H. Tembine, and J. Ben-Othman, "Underwater Wireless Sensor Networks: A Survey on Enabling Technologies, Localization Protocols, and Internet of Underwater Things," *IEEE Access*, vol. 7, pp. 96879–96899, 2019, doi: 10.1109/ACCESS.2019.2928876.
- N. Javaid, U. Shakeel, A. Ahmad, N. Alrajeh, Z. A. Khan, and N. Guizani, "DRADS: depth and reliability aware delay sensitive cooperative routing for underwater wireless sensor networks," *Wireless Networks*, vol. 25, no. 2, pp. 777–789, Feb. 2019, doi: 10.1007/s11276-017-1591-1.
- G. Han, X. Long, C. Zhu, M. Guizani, and W. Zhang, "A High-Availability Data Collection Scheme based on Multi-AUVs for Underwater Sensor Networks," *IEEE Transactions on Mobile Computing*, vol. 19, no. 5, pp. 1010–1022, May 2020, doi: 10.1109/TMC.2019.2907854.
- N. Li, J.-F. Martínez, J. Meneses Chaus, and M. Eckert, "A Survey on Underwater Acoustic Sensor Network Routing Protocols," *Sensors*, vol. 16, no. 3, p. 414, Mar. 2016, doi: 10.3390/s16030414.
- S. Sahana, K. Singh, R. Kumar, and S. Das, "A review of underwater wireless sensor network routing protocols and challenges," in *Advances in Intelligent Systems and Computing*, 2018, vol. 638, pp. 505–512, doi: 10.1007/978-981-10-6005-2\_51.
- H. Yan, Z. J. Shi, and J. H. Cui, "DBR: Depth-based routing for underwater sensor networks," in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 2008, vol. 4982 LNCS, pp. 72–86, doi: 10.1007/978-3-540-79549-0\_7.
- M. Zuba, M. Fagan, J. H. Cui, and Z. Shi, "A vulnerability study of geographic routing in underwater acoustic networks," in *2013 IEEE Conference on Communications and Network Security*, CNS 2013, 2013, pp. 109–117, doi: 10.1109/CNS.2013.6682698.
- P. Xie, J. H. Cui, and L. Lao, "VBF: Vector-based forwarding protocol for underwater sensor networks," in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 2006, vol. 3976 LNCS, pp. 1216–1221, doi: 10.1007/11753810\_111.

- [15] A. Wahid, S. Lee, H. J. Jeong, and D. Kim, "EEDBR: Energy-efficient depth-based routing protocol for underwater wireless sensor networks," in *Communications in Computer and Information Science*, 2011, vol. 195 CCIS, pp. 223–234, doi: 10.1007/978-3-642-24267-0\_27.
- [16] S. Gul, S. H. Jokhio, and I. A. Jokhio, "Light-weight depth-based routing for underwater wireless sensor network," in 2018 International Conference on Advancements in Computational Sciences, ICACS 2018, Apr. 2018, vol. 2018-January, pp. 1–7, doi: 10.1109/ICACS.2018.8333483.
- [17] N. Javaid, M. R. Jafri, Z. A. Khan, U. Qasim, T. A. Alghamdi, and M. Ali, "IAMCTD: Improved adaptive mobility of courier nodes in threshold-optimized DBR protocol for underwater wireless sensor networks," *International Journal of Distributed Sensor Networks*, vol. 2014, Nov. 2014, doi: 10.1155/2014/213012.
- [18] T. Ahmed, M. Chaudhary, M. Kaleem, and S. Nazir, "Optimized depth-based routing protocol for underwater wireless sensor networks," in *ICOSST 2016 - 2016 International Conference on Open Source Systems and Technologies*, Proceedings, Jan. 2017, pp. 147–150, doi: 10.1109/ICOSST.2016.7838592.
- [19] E. Souza, H. C. Wong, I. Cunha, A. A. F. Loureiro, L. F. M. Vieira, and L. B. Oliveira, "End-to-end authentication in under-water sensor networks," in *Proceedings - International Symposium on Computers and Communications*, 2013, pp. 299–304, doi: 10.1109/ISCC.2013.6754963.
- [20] M. Zuba, M. Fagan, Z. Shi, and J. H. Cui, "A resilient pressure routing scheme for underwater acoustic networks," in 2014 IEEE Global Communications Conference, GLOBECOM 2014, Feb. 2014, pp. 637–642, doi: 10.1109/GLOCOM.2014.7036879.

# Product Recommendation in Offline Retail Industry by using Collaborative Filtering

Bayu Yudha Pratama<sup>1</sup>, Indra Budi<sup>2\*</sup>, Arlisa Yuliawati<sup>3</sup>

Faculty of Computer Science  
Universitas Indonesia, Depok  
West Java, Indonesia

**Abstract**—The variety of purchased products is important for retailers. When a customer buys a specific product in a large number, the customer might get benefit, such as more discounts. On contrary, this could harm the retailers since only some products are sold quickly. Due to this problem, big retailers try to entice customers to buy many variations of products. For an offline retailer, promoting specific products based on the markets' taste is quite challenging because of the unavailability of information regarding customers' preferences. This study utilized four years of purchase transaction data to implicitly find customers' ratings or feedback towards specific products they have purchased. This study employed two Collaborative Filtering methods in generating product recommendations for customers and find the best method. The result shows that the Memory-based approach (k-NN Algorithm) outperformed the Model-based (SVD Matrix Factorization). Another finding is that the more data training being used, the better the performance of the recommendation system will result. To cope with the data scalability issue, customer segmentation through k-Means Clustering was applied. The result implies that this is not necessary since it failed to boost up the models' accuracy. The result of the recommendation system is then applied in a suggested business process for a specific offline retailer shop.

**Keywords**—Recommendation system; offline retail store; memory-based collaborative filtering; customer segmentation

## I. INTRODUCTION

Recommendation system is a collection of tools and techniques to provide products or services suggestions for users [1]. The existence of this system allows companies to develop a marketing strategy, attract more customers, and increase sales. Therefore, many companies try to implement a recommendation system for their business interest. Recommendation system has been applied in a variety of industries. It can be found in the entertainment domain (music, movies, TV shows, books), news or tourism sites, e-commerce, e-library, and e-learning systems [2].

Even though the recommendation system has been extensively used in e-commerce domain as described in [2] and [3], research in [4] argued that it can also be implemented in the traditional retail stores. They suggested personalization as the next possible strategy for this kind of retailers. Personalization establishes a one-to-one relationship between

the retailer and the customer. By using a one-to-one relationship, a retailer can remember details and preferences for each customer. These preferences can be utilized to identify customer personal needs, wants, and demands. This personalization strategy can be realized by the implementation of recommendation system. E-commerce has already implemented the recommendation system with many benefits such as boosting up customer level of interaction, increasing sales, the diversity of items sold, customer satisfaction or loyalty, and also understanding customers' demand better [5]. Such benefits are expected to be achieved in traditional or offline retail stores.

A traditional or offline retail store differs from e-commerce in several aspects. The first is that traditional retail store still having a physical store for storage, display, and transaction. This is costly for them to keep rarely sold items in inventory [4]. On the other hand, the diversity of customers' demands always increasing. Retailers must be able to correctly identify customers' demand as well as offering a variety of products. This is a way to keep the goods in inventory to keep moving. It differs from the e-commerce setting on which they do not always have physical storage to keep their products.

The other distinction between offline and online commerce is that in an offline retail store, it is often difficult to access customers' purchase history and observe their purchase behavior. In some cases, the offline stores only keep the transactions without knowing who the buyers are. Of course, this is different from the online based stores, which the customers' identities and activities are recorded in the system. From this situation, there is a problem related to the unavailability of the user related data. Moreover, it is difficult for all customers to give feedback towards items they have bought. The impact is that the predictions are often poor when other users or customers are looking for recommendations for the rarely rated items [4]. This is different from e-commerce business where users usually explicitly asked to give a rating when their transaction is completed. The availability of such feedbacks or ratings is crucial for recommendation system studies since customers' feedback or ratings becomes one of available sources to derive any information that may be useful for other customers [6].

\*Corresponding Author

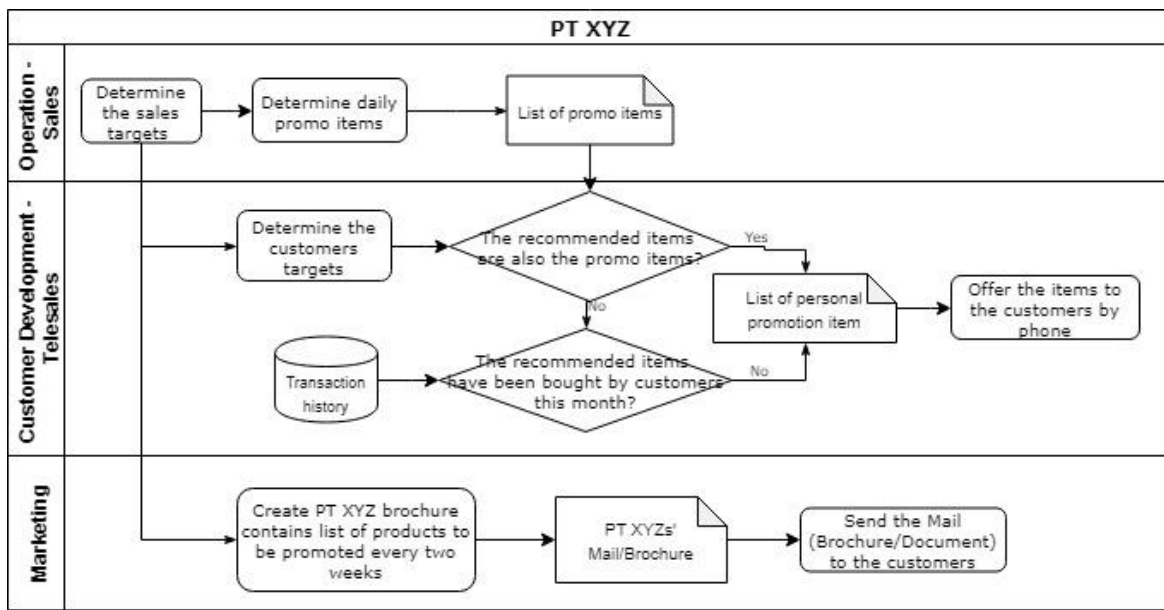


Fig. 1. The Business Process of PT XYZ in Promoting their Products.

This study is conducted based on the situation in a wholesales retail store in Indonesia, called PT XYZ. This retailer is categorized as a Broad-Deep-Mix retail store, which has large variety of products and heterogeneous customers [4]. This retailer has a personalized marketing scheme, by offering specific items either by phone or by sending product brochure to specific customers as presented in Fig. 1. Based on previous explanation, there are problems related to personalization in offline retail store. Inaccurate promotion target, either the items to be offered or the targeted customers can lead to the failure in selling those items. Hence, it will fail to fulfill their purpose to sell various products to customers, as well as fail to improve their profit. Therefore, the aim of this paper is to find a suitable approach to give better recommendations for customers on an offline retailer (specifically for PT XYZ) both empirically and practically, as well as considering the personalization approach that has already been used.

Following section will describe the related works and summarized proposed works about recommendation system in offline retailer. It continued to the next two sections with the materials and the research methodology. The result and discussion are presented in following section, and at the end of this paper the conclusion of this study will be presented.

## II. LITERATURE REVIEW

### A. The Problem of Limited Customers' Data

Specific for offline retailer, the main issue of this study is related to the absence of explicit product ratings or feedback by customers. As the main input for the process of the recommendation system, it is important to define such things. A research in [7] summarizes about some experiments to deal with this problem. Many of them use the association rule that ended up with the lack of personalization. Study by [7] itself utilized the smart fitting room, i.e. the IT artifact that gives product recommendation to the customers through a screen stored in the individual cabin. Started from the use of Association Rule Mining, their study shows that combine the

information from customers' interaction to the screen with the contextual information about products could improve the product recommendation in fashion stores. Another study about fashion retailer as written in [8] combines the online product click data and offline product sale data to reflect the preference of the customers. This experiment concludes that it is better to substitute than complement the products in the recommendation system. The percentage of purchase by using the former approach is higher than the latter. Nevertheless, there is less information related to the use of online and offline data combination. Recalling the situation of PT XYZ, those used in these previous studies did not owned by PT XYZ. But the idea is that these two retailers utilized the availability system that interacts with their customers, e.g. the smart fitting room and the online system to combine with the offline product sale. In PT XYZ, there is a membership system. Customers who register for membership in this retail store have their transaction history recorded in the stores' database. This study utilizes this data to generate the customers' feedback.

### B. Techniques in Recommendation System Study

There are several techniques to identify which items are recommended for specific users as summarized in [1]. They are distinguished based on the domain, knowledge, or the algorithm being used.

1) *Collaborative Filtering (CF)*: This approach gives a recommendation of items based on the similar preferences of other users in the past. The similarity between users/items can be inferred from their previous behavior such as rating or buying history.

2) *Content-based Filtering (CBF)*: This approach recommends items that are similar to what the user liked in the past. The similarity of items is obtained from their feature and description. The recommendation is compiled from the attribute information of items.

3) *Demographic recommender system*: The recommendation is given based on the demographic information, such as location, language, and age of user or customers. This approach implies that people with different demographic background should not receive the same recommendation.

4) *Knowledge-based recommender system*: By using this approach, system gives a recommendation based on specific knowledge over items from experts. It is then matched up to the items' benefit for users. The similarity is implied from the user needs and items' function. This approach identifies the similarity based on match "answer/solution" to users' "question/problem description".

5) *Community-based recommender system*: The recommendation is implied from the preference of users' circle friend. This approach is popular on social network-based system.

Based on the previously mentioned techniques, there are some drawbacks on some techniques to be applied to an offline retailer, specifically for the case of PT XYZ. CBF approach is not applicable for this case since traditional retailers usually do not store comprehensive description about their products rather than only consists of name, price, main categories, and sub-categories (e.g. dry-food/fresh-food/non-food). There is also limited information about customers' profile except for members of the stores (if any), so approach based on demographic is also not suitable for this case. Similarly, offline stores also do not maintain how their customers connected each other. By this condition it is difficult to get the recommendation based on customers' circle or community. Lastly, offline stores usually provide various kinds of items, so find experts for various kinds or categories of items is another problem for employing knowledge-based approach. From this analysis, it implies that CF is the most suitable approach for the case of PT XYZ.

### C. Data Scalability vs Customer Segmentation

Another issue in the recommendation system in general is the data scalability. It is caused by the huge amount of data that leads to the accuracy problem of the recommendation system [9], [10]. In some cases, this issue is related to the algorithm or approach being used to build the system, i.e. the use of Collaborative Filtering. Its performance on scalability is still poor given a huge user and item base [11]. The previous study in [9], [10], [12], as well as [13], try to include the customer segmentation process to cope with this challenge. This process is also used to identify profitable customers [9] [10], not only to make the data become smaller.

Commonly used approach to differentiate the customers into several segments is the simple-yet-powerful RFM model (Recency, Frequency, and Monetary). Recency is defined as the last time (in a month) a customer completed a transaction, Frequency describes how many total transactions for each customer, and Monetary calculates how much they buy in value [13]. Some examples of algorithms that can be used for segmenting customers are Artificial Neural Network (ANN) [10], *k*-Means clustering [9], Expectation-Maximization (EM) [13]. Both studies by [9] and [13] use this RFM model as the

segmentation method and shows a satisfiable result. The difference is the former use Association Rules and hybrid method, while the later use *k*-NN. The performance of former study was affected by using hybrid method while the later by customer segmentation. Nevertheless, the study by [9] is based on homogeneous store, i.e. only sell one kind of item. It is different from the case used in [13], as well as PT XYZ that are selling various items (heterogeneous retailer).

### D. Proposed Work

Based on the previous analysis in offline retailer about the problem of limited customers' data, various techniques to use, and the challenge in the data scalability, this study propose some steps to be applied in the case of PT XYZ as a wholesale offline retailer. The first, to cope with the unavailability of the customers' ratings data, this study takes advantage of the membership system applied in PT XYZ. Basically, customers' activities recorded in the system is elaborated, and the rating data is generated implicitly from the customers' purchase pattern based on a specific transformation metric. An assumption is made related to the result of this transformation: "The more frequently customers purchase an item, the higher the rating they implicitly give". The second, Collaborative Filtering will be used in this study. Since there are two approach in this technique, this study also tries to find the best approach. The last, to cope with the data scalability, this study adapt the use of RFM model to apply the customer segmentation. An experiment is employed to elaborate whether this approach also gives better performance compared to the original process without segmenting the customers.

## III. COLLABORATIVE FILTERING

CF is the most successful and widely used recommendation technique [14], [15]. CF utilizes a user-item matrix to make the recommendations. Suppose there is a set of *m* users  $U = \{u_1, u_2, \dots, u_m\}$  and a set of *n* items  $I = \{i_1, i_2, \dots, i_n\}$ , CF constructs an  $m \times n$  matrix *R* representing the preference of users to items. For each user, the list of relevance items can be viewed from the descending order of matrix values related to the user. If there are two users giving the same rating to an item, then it can be implied that they have the same taste or preference. As an example, from Table I, the relevance item for User\_C is Item\_3 and Item\_1, while another information is that User\_A and User\_C have the same taste or preference toward Item\_3. There are two approaches that are commonly used in CF, they are Memory-based and Model-based [16]. A study by [17] specifically compared these two approaches in the e-commerce domain. The result indicates that the Model-based is better than Memory-based not only in the accuracy and the relevancies of the recommendation but also in the computational time. Following subsections will explain more about these approaches.

TABLE I. AN ILLUSTRATION OF USERS PREFERENCE TOWARD ITEMS AS INPUT FOR USER-ITEM MATRIX

|        | Item_1 | Item_2 | Item_3 | Item_4 |
|--------|--------|--------|--------|--------|
| User_A | 5      | -      | 4      | -      |
| User_B | -      | 3      | -      | 1      |
| User_C | 2      | -      | 4      | -      |



### A. Memory-based Collaborative Filtering

This is also called Neighborhood-based CF. This is the most popular method in the recommendation system domain [15]. It generally shows that similar users have similar rating behavior, so do with similar items, they receive similar ratings [16]. The similarity can be defined among users (User-based) or items (Item-based). The distinction is that in the former case, the ratings are predicted using those of neighboring users, while in the latter case, they are predicted using the users' own ratings on neighboring or closely related items [16]. In this study, the User-based CF is chosen instead of Item-based CF since based on the experiment the former is better in term of accuracy, time, and space complexity. This might be caused by the number of items that is greater than the number of users in the case of PT XYZ. It is also known that User-based CF is one of the most widely used among CF approaches [18].

This approach uses k-Nearest Neighbor (k-NN) algorithm to find the top k similar users and predict the rating for specific items that have been bought by those k users. The similarity between users can be calculated by using distance metrics, such as Cosine Similarity or Pearson Correlation. The recommendation then will be given based on the rating calculation of items obtained from each of the k users. Generally, Memory-based or Neighborhood-based is a simple and straightforward approach yet still have an accurate prediction. Nevertheless, it has some limitation such as a scalability issue in a large matrix and the cold-start problem where the model cannot recommend a new user/item [10], [12].

### B. Model-based Collaborative Filtering

Basically, this model tries to find the hidden factors in the original/initial matrix [1]. To create the prediction model, it applies various data mining technique, such as the Decision Tree, Bayesian, and Latent Factor model. This approach is better than the Memory-based approach in several ways [19]. For example, it has better scalability since it has a good performance in a large matrix as well as the better accuracy, i.e by using Singular Value Decomposition (SVD) as the Latent Factor model [20]. SVD runs by decomposing an initial  $n$  users  $\times$   $m$  products matrix into three matrices:  $U$ ,  $\lambda$ , and  $V$  as illustrated in Fig. 2. These three matrices are updated continuously until the result of their multiplication is approaching the initial matrix. In Model-based CF approach, these matrices are equivalent to users  $U$  ( $n$  users  $\times$   $r$  concepts), concepts  $\lambda$  ( $r$  concepts  $\times$   $r$  concepts), and products  $V$  ( $r$  concepts  $\times$   $m$  products) matrices. One of the important parameters in SVD is the number of the Latent Factor, which refers to the number of concepts that are hidden in the initial matrix.

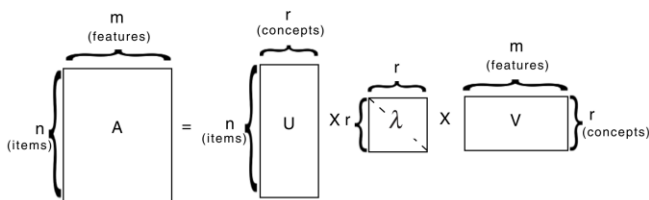


Fig. 2. Illustration of SVD [1].

## IV. RESEARCH METHODOLOGY

Fig. 3 shows the research methodology of this study. Data sources are provided by PT XYZ as one of the largest retail chains in Indonesia. PT XYZ has a unique business model, it combines retail and wholesale sales. They serve both professional customers and end-user/individual consumers. The data set contains 2.5 years of transactions data from one of their branches. This research is mainly divided into two experiments. The first experiment examines the performance of both Memory-based and Model-based CF. The former approach employs k-NN algorithm and the later uses SVD matrix factorization. This study also observes the optimal value of k and the number of the latent factor for k-NN algorithm and SVD matrix factorization respectively. The second experiment focuses on the application of customer segmentation. k-Means algorithm is then used to differentiate customer into several categories based on the RFM model. The result of both experiments is compared based on the Root-Mean-Square-Error (RMSE) metric, as it is a stable means of comparison between models [21]. The better approach based on the minimum RMSE value is then applied for the recommendation system. Following subsections explain each process in this research.

### A. Data Collecting and Clean Up

The data was taken in one of PT XYZ branch store. This data consists of information about users, products, and transactions from professional customers, including smaller retailer, hotel, restaurant, or catering, who buy products for their business necessities. This study only use the transactions from professional customers and omitting those from end-users because PT XYZ wants to test personalization to their professional customers first, before moving on to end-user later. Besides, PT XYZ already has a direct relation to professional customers by offering products via telephone. The expectation is to expand the relationship into one-to-one personalized marketing.

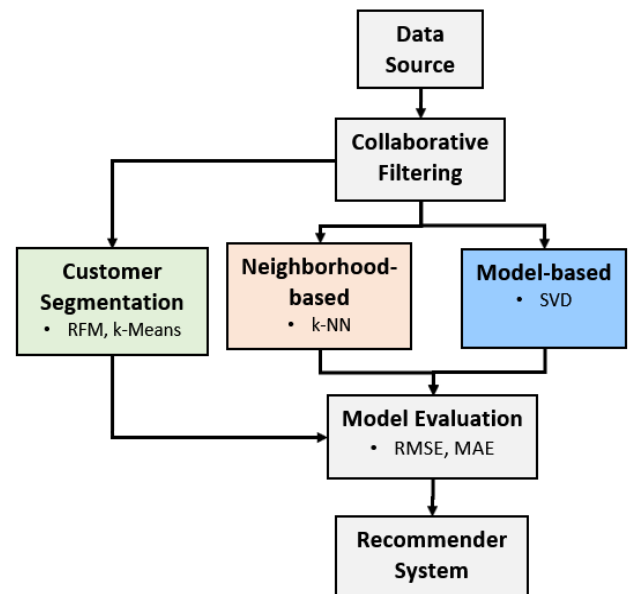


Fig. 3. Research Methodology.

The information about the user consists of customer number (cust\_no) and customer name (cust\_name). Product data also consists of product code (prod\_cd) and product name (prod\_name). While the transaction data is the purchase of an item (prod\_cd) by a customer (cust\_no) in a specific time (sale\_day). The data collected from 2.5 years of transactions from one store within time windows from January 2015 to September 2017. There are numbers of transactions that are not relevant for this study such as the internal purchase or transactions involving non-trade items such as insurance, administration fee, and shipping charges. Clean-up process is carried out by removing this type of transactions. From this data collection, a total of 8,515 customers, 23,532 items, and over 2,1 million transaction records are obtained.

**B. Data Transformation**

The Collaborative Filtering method utilizes the user-item rating matrix for making predictions. This matrix describes user *u* giving item *i* a rating value  $r_{ui}$ . Since this kind of data is not available in the case of a traditional retail store, the customers' purchase history is utilized such that it implicitly represents customers' feedback towards what they have bought. This approach is suitable for the characteristics of professional customers. Basically, they tend to continuously buy products for their day-to-day business operation. If they repeatedly buy a specific item, it is an implicit indication that they like that item because it shows that the more frequently they purchase an item, the higher the rating they give.

From the provided data, a customer did a transaction of a specific product in a specific time. From this data the number of purchases by customer for each item is extracted by using aggregate function. Table 2 presents an example of some complete transaction data, while Table 3 show the result of the aggregate process. The quantity and value of a transaction to create a level playing field between transactions in big and small companies are ignored. This is because the big company usually has bigger transaction value than the small company. From this process, one purchase of an item is considered as one transaction regardless to its total value of purchase.

The user-item matrix come from a purchase frequency matrix consists of user *u* buy item *i* for  $f_{ui}$  times. This information is converted into a rating value by using a min-max scaling algorithm as shown in Equation (1). The variable  $f_{ui}$  represents the frequency of purchase while  $f_{min}$  and  $f_{max}$  represent the minimum and maximum purchase respectively for each item. The variable  $r_{ui}$  represent the transformed rating value ranged from 1 ( $r_{min}$ ) to 5 ( $r_{max}$ ).

$$r_{ui} = \left[ \left( \frac{f_{ui} - f_{min}}{f_{max} - f_{min}} \right) \times (r_{max} - r_{min}) \right] + r_{min} \quad (1)$$

Table 4 shows the example of original data, consist of the frequency of purchasing an item by a customer. The definition of the value of  $f_{min}$  and  $f_{max}$  are based on the information of each column that represent each item instead of considering the minimum and maximum frequency of the whole matrix. Therefore, the sale rate of items can be analyzed whether they are fast-moving items or slow-moving items. Then Table 5 shows the transformation result, which is the implicit rating value of an item given by a user.

In some cases, the purchase frequency of specific items is too high, far from the purchase frequency on average case. This can result in the skewness of the rating value gather in a lower value. To overcome this condition, the data with too high purchase frequency are removed. The removal process is conducted by modelling the purchase distribution of frequency of each item as a normal distribution. When the purchase frequency ( $f_{ui}$  value) exceeding the normal distribution limit resulted from Equation (2), then this outlier is removed. The limit ( $T_i$  value) is obtained from the mean value of purchase frequency of item *i* ( $\mu_i$ ) added by 3 times the standard deviation ( $\sigma_i$ ) value of item *i*'s purchase frequency.

$$T_i = \mu_i + (3 \times \sigma_i) \quad (2)$$

**C. Model Development**

In this study, two collaborative-filtering approaches is compared, they are Memory-based by using k-NN and Model-based through SVD matrix factorization. These models have some adjustable parameters to obtain optimal performance as described below. The model for both k-NN and SVD is developed by using the optimal parameter values found in this experiment and evaluate their performance.

TABLE II. AN EXAMPLE OF COMPLETE TRANSACTIONS

| cust_no       | prod_cd    | freq | time   |
|---------------|------------|------|--------|
| 6001000580465 | 1029504000 | 1    | 2015Q1 |
| 6001000580465 | 1029504000 | 1    | 2015Q1 |
| 6090001779264 | 1020248000 | 1    | 2016Q2 |
| 6090001779264 | 1020248000 | 1    | 2016Q3 |
| 6090001779264 | 1020248000 | 1    | 2016Q4 |
| 6090001779264 | 0032080000 | 1    | 2017Q2 |

TABLE III. THE RESULT OF THE AGGREGATE PROCESS

| cust_no       | prod_cd    | freq |
|---------------|------------|------|
| 6001000580465 | 1029504000 | 2    |
| 6090001779264 | 1020248000 | 3    |
| 6090001779264 | 0032080000 | 1    |

TABLE IV. AN EXAMPLE OF USER-ITEM PURCHASE FREQUENCY MATRIX

|        | Item_1 | Item_2 | Item_3 | Item_4 | Item_5 |
|--------|--------|--------|--------|--------|--------|
| User_A | 2      | -      | -      | 7      | -      |
| User_B | -      | 1      | -      | -      | 5      |
| User_C | 8      | -      | 3      | 2      | 1      |
| User_D | -      | 20     | -      | 3      | -      |
| User_E | 10     | -      | -      | -      | -      |

TABLE V. THE DATA TRANSFORMATION RESULT: USER-ITEM RATING MATRIX

|        | Item_1 | Item_2 | Item_3 | Item_4 | Item_5 |
|--------|--------|--------|--------|--------|--------|
| User_A | 1      | 0      | 0      | 5      | 0      |
| User_B | 0      | 1      | 0      | 0      | 5      |
| User_C | 4      | 0      | 5      | 1      | 1      |
| User_D | 0      | 5      | 0      | 2      | 0      |
| User_E | 5      | 0      | 0      | 0      | 0      |

Evaluation for all experiments are based on the Root-Mean-Square-Error (RMSE) error metric, as shown in Equation (3). In this equation,  $\hat{R}$  is the collection of rating prediction, while  $r_{ui}$  is the actual rating in testing data set and  $\hat{r}_{ui}$  is the predicted rating from the model.

$$RMSE = \sqrt{\frac{1}{|\hat{R}|} \sum_{\hat{r}_{ui} \in \hat{R}} (r_{ui} - \hat{r}_{ui})^2} \quad (3)$$

1) In general,  $k$ -NN algorithm finds  $k$  most similar users ( $N_i^k(u)$ ) based on the previous buying or rating pattern. The similarity between two users ( $sim(u, v)$ ) is calculated through cosine similarity as shown in Equation (4) and the rating prediction ( $\hat{r}_{ui}$ ) is computed by using Equation (5). The variable  $I_{uv}$  represents the item that is rated by user  $u$  and  $v$ ,  $r_{ui}$  is the rating prediction value from user  $u$  for item  $i$ , and  $r_{vi}$  is the actual rating value from user  $v$  for item  $i$ . One of the important parameters for  $k$ -NN is the number of nearest/similar users (the value of  $k$ ). Fig 4 presents the result of an experiment to find the optimal value of  $k$ . The experiment is performed with training data collected from 2015 Q1 to 2017 Q2 and use the testing data from transactions in 2017 Q3. RMSE metric is used to compare the result.  $k$ -NN delivers best result when the value of  $k$  is greater or equals to 80.

$$sim(u, v) = \frac{\sum_{i \in I_{uv}} r_{ui} \cdot r_{vi}}{\sqrt{\sum_{i \in I_{uv}} r_{ui}^2} \cdot \sqrt{\sum_{i \in I_{uv}} r_{vi}^2}} \quad (4)$$

$$\hat{r}_{ui} = \frac{\sum_{v \in N_i^k(u)} sim(u, v) \cdot r_{vi}}{\sum_{v \in N_i^k(u)} sim(u, v)} \quad (5)$$

2) SVD decompose original user-item matrix into user, concept, and item matrices. As explained in [1], the decomposition or the factorization process maps the users and items into latent factors space. This latent space explains ratings by characterizing both items and users on factors that are inferred from users' feedback. Rating prediction that describes the overall interest of the user in characteristics of the item is computed within Equation (6). Each item  $i$  is associated with a vector  $q_i$  that measure the extent of item  $i$  possesses those factors. While each user  $u$  is associated with a vector  $p_u$  that measure the interest of user  $u$  towards factors. This study perform an experiment with SVD to evaluate the optimal number of factor, and the result is shown in Fig 5. The best result is achieved by using factors of less than or equals to 10. The result shows that SVD is able to find few factors in the original matrix, presumably because of the heterogeneous nature of items offered.

$$\hat{r}_{ui} = q_i^T p_u \quad (6)$$

After obtaining the optimal parameter values, then evaluation towards the use of Model-based CF by using  $k$ -NN and Memory-based CF through SVD matrix factorization in

product recommendation is conducted. Each model is built based on those previously found parameter values, i.e.  $k = 80$  for the implementation of  $k$ -NN algorithm and the SVD matrix factorization is trained by using the number of factors = 10.

The data set contains time-series data and divided into quarter year data. There are eleven quarters from 2015 Q1 until 2017 Q3. The transaction data from 2015 Q1 to 2017 Q2 are used as the training data set, while the transaction data in 2017 Q3 are used as the testing data set. The variant of each training set is made by removing data from the oldest quarter. This approach is used to determine the amount of data and the extent to which optimal training data points required. Is it sufficient to train only the newest data from the last quarter of the year? Or is it necessary to train as much data as possible? The models process the user-item rating matrix and produce rating prediction for every pair (item  $u$ , user  $i$ ) in the testing data set that is 2017 Q3 transactions (105,177 records). The rating prediction from each model is then compared to the factual sales data in this testing data set.

The better CF approach will be used to apply customer segmentation based on customer lifetime value using the RFM model. This to answer whether customer segmentation has a positive impact on the development of the product recommendation system. The segmentation process employs  $k$ -Means clustering algorithm with  $k = 3$ . The model is then tested by using the training data starting from 2015 Q1.

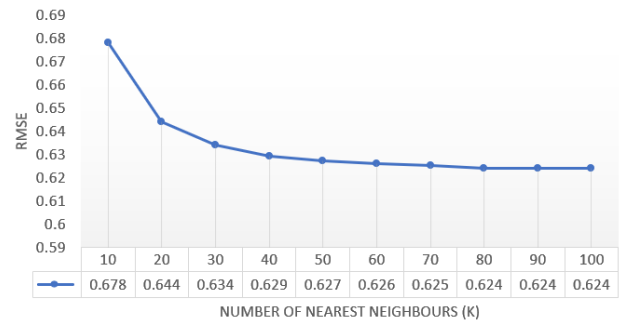


Fig. 4. The Experiment Result of Finding the Optimal Value of  $k$ .

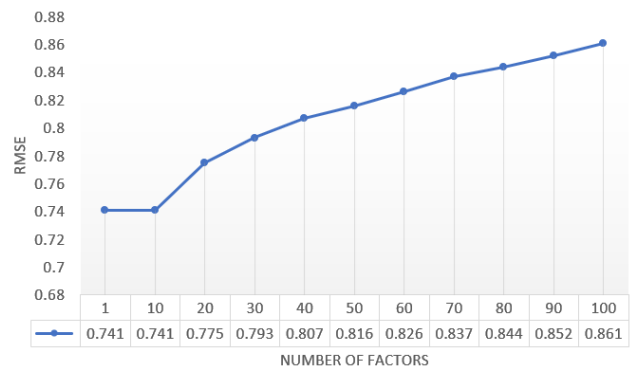


Fig. 5. The Experiment Result of Finding Optimal Value of the Latent Factors.

## V. RESULT AND DISCUSSION

### A. Memory-based vs Model-based Collaborative Filtering

The experiment in comparing the performance of using Memory-based by using k-NN and Model-based through SVD matrix factorization generally shows that the former is slightly outperforms the latter. As presented in Table 6, by using k-NN a better prediction is obtained compared to the use of SVD matrix factorization. Among the ten variants of data training, the error value by using k-NN is always below the error value of SVD, as long as the number of training data being used come from more than two quarters. This result is caused by the size of the matrix that is not exceptionally high and is still sufficient to be performed in the Memory-based approach.

Meanwhile, Table 6 also shows that the smallest value of RMSE for both approaches is obtained when all of the ten quarters training data are utilized. By comparing the result from each training data set, as the number of training data is decreased, the RMSE values are increasing consistently. From this point, it implies that the older the starting point of the training data or the more data available for training process, the less error value will be obtained, so the prediction result would be better.

### B. Customer Segmentation

This experiment tries to differentiate the customers into several segments before applying the better model, which is Memory-based by using k-NN. Table 7 presents the characteristics of each segment and the result of this experiment. It implies that Recency is the most important variable since most of the transactions come from customers who have low Recency or recently purchased items. These customers tend to have high Frequency and Monetary value because they are buying products frequently (repeat buying). This customer segment is defined as “active”. On contrary, customers who have high Recency are rarely done transactions. This is based on their level of Frequency and Monetary that are ranged from medium to low. They called as “inactive” customers, and the remaining cluster between the “active” and “inactive” ones are called “semi-active” customers.

The other finding is that based on the RMSE value, the use of training data that only come from “active” customer results in a slightly better performance of the recommendation system compared to those from both “semi-active” or “inactive” customer segments. While the use of training data that come from both of these last two segments, yield the significantly decreasing performance compared to both of the use of only “active” customer or without clustering process.

Although the data from “active” customer yield good performance, the error value is still higher than the excluding of the segmentation process. This result implies that segmenting customers into several segments is failed to improve the performance of the models. One factor that causes this is the size of the data training. By dividing the data based on a specific customer, each cluster of training data has a

fewer amount of data than the combined data. This result strengthens the previous verdict from the first experiment that the more data training being used, the better model will result.

### C. Implementations

Experiment in this study basically shows that it is better to use Memory-based CF with the use of as many as training data. The recommendation provided by this model consists of the list of relevant products related to specific customers. This list is generated by sorting the prediction rating value obtained from the model. From this list of recommendations, the additional steps are inserted in the promotion flow of PT XYZ as the offline retailer so this would help them to give better promotion result as well as increasing their profit.

Fig. 6 shows the suggested business model for PT XYZ. The sales division promote their products by phone but based on the recommended products obtained from the output of the recommendation system. If these products are included in the promo items, then they can be immediately offered to the customers. The combination of recommended products with specific discounts will attract customers’ intention to buy them. This can improve the success rate of the promotion. If customers have already purchased a specific item on the recommendation list, the sales team should skip or remove this item then move to the next item on the list. This approach can improve the variety of items to be purchased by the customers.

Another advantage of the recommendation systems’ output is for the marketing division. The actions are similar to the previous promotion by phone. Furthermore, they can utilize PT XYZs’ personal mail to create personalization towards each customer. The content of the mail or brochure consists of relevant products for the specific customer. In the normal pipeline, they send this mail to the customer once in two weeks with the same content in one segment. Now they can adjust the content and the frequency of the mailing. This approach will decrease the number of human resources needed to make promotions by phone, in case there are only a few numbers of resource available. Sending this mail to all customers also effective to improve the sale rate.

TABLE VI. MODEL EVALUATION RESULT BASED ON RMSE VALUE

| Training Data     | Number of Quarter | RMSE  |       |
|-------------------|-------------------|-------|-------|
|                   |                   | k-NN  | SVD   |
| 2015 Q1 – 2017 Q2 | 10                | 0.624 | 0.750 |
| 2015 Q2 – 2017 Q2 | 9                 | 0.639 | 0.753 |
| 2015 Q3 – 2017 Q2 | 8                 | 0.653 | 0.762 |
| 2015 Q4 – 2017 Q2 | 7                 | 0.673 | 0.770 |
| 2016 Q1 – 2017 Q2 | 6                 | 0.699 | 0.788 |
| 2016 Q2 – 2017 Q2 | 5                 | 0.727 | 0.802 |
| 2016 Q3 – 2017 Q2 | 4                 | 0.762 | 0.824 |
| 2016 Q4 – 2017 Q2 | 3                 | 0.818 | 0.850 |
| 2017 Q1 – 2017 Q2 | 2                 | 0.890 | 0.887 |
| 2017 Q2 – 2017 Q2 | 1                 | 0.974 | 0.926 |

TABLE VII. THE RESULT OF CUSTOMER SEGMENTATION EXPERIMENT

| Cluster           | Recency | Frequency    | Monetary     | Number of customers | % of transactions | RMSE  |
|-------------------|---------|--------------|--------------|---------------------|-------------------|-------|
| 1 – “active”      | Low     | High, Medium | High, Medium | 3,805               | 83%               | 0.636 |
| 2 – “semi-active” | Medium  | Medium, Low  | Medium, Low  | 2,513               | 12%               | 1.454 |
| 3 – “inactive”    | High    | Medium, Low  | Medium, Low  | 2,198               | 5%                | 1.515 |
| No cluster        | -       | -            | -            | 8,515               | 100%              | 0.625 |

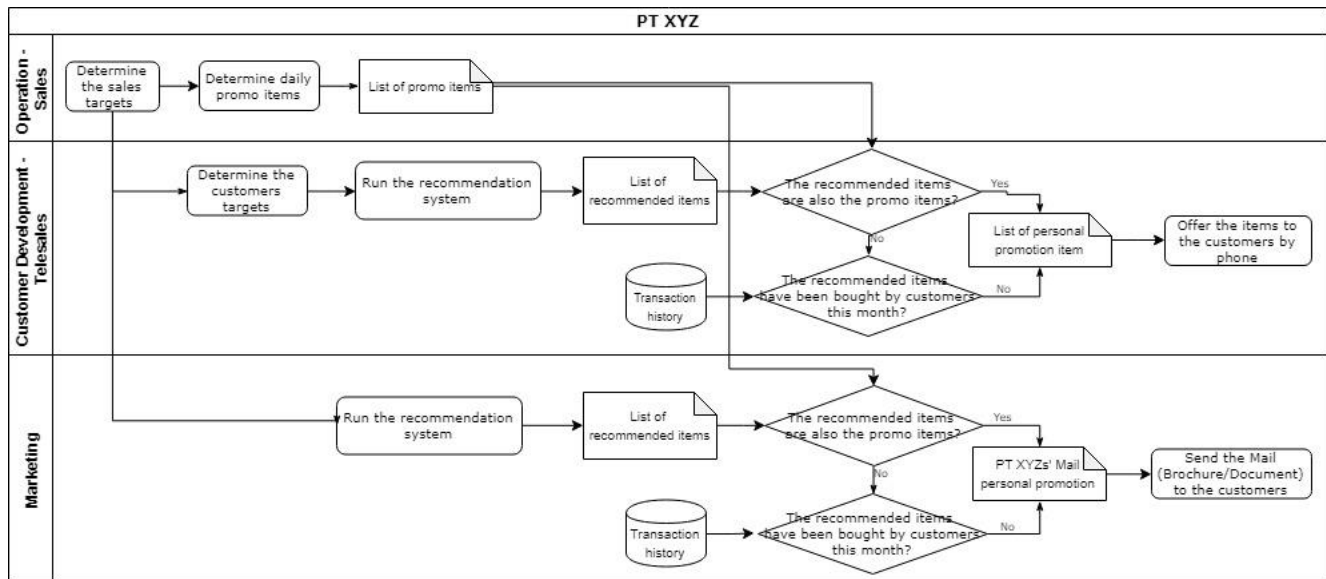


Fig. 6. Suggested Pipeline for the Promotion Process in PT XYZ.

## VI. CONCLUSION

The goal of this research is to develop a better approach to be implemented in a recommendation system for a traditional or offline retail store. Before suggesting the new pipeline for the promotion process in PT XYZ, some experiments are conducted. The first, related to one of the problems in an offline retail store is the unavailability of customers' rating/feedback data towards products. This study deals with this problem by constructing a user-item matrix based on the number of purchases by the user as an implicit feedback score. The higher number of purchases implies the higher feedback or rating being given by a customer. Nevertheless, this approach only covers the data from customers who join the membership of PT XYZ.

The second, the widely used Collaborative Filtering approach is applied in this study. An experiment to find a better approach by using Memory-based and Model-based Collaborative Filtering is conducted to predict the rating given by the target customer. The result shows that Memory-based CF with k-NN outperforms the Model-based CF through SVD. Regarding the amount of training data, it can be concluded that more data training is always resulting in a better prediction. This result is concluded by conducting an experiment to find the optimal value of the parameters being used in k-NN algorithm and SVD matrix factorization. In k-NN model. The neighborhood size (the value of k) is directly proportional with the model accuracy. The higher the neighborhood size, the better accuracy is reached. On

contrary, the number of factors that are used in SVD matrix factorization is inversely proportional, where best performance is found on the use of the fewest number of factors.

The last, related to the data scalability problem, some previous studies to cope with this problem by creating a smaller size of data is adapted. One of the approaches in previous studies is using customer segmentation based on specific criteria. Their studies show that it was better to add this segmentation process. This study applied this approach to differentiate the training data based on RFM model before running the recommendation system model, and the result implies that in the term of accuracy, it is not necessary to do this since the performance of the model is failed to be improved. Related to the case being used in this study, the more training data is needed to build a better model. While by segmenting the customers into three categories (active, semi-active, inactive) reducing the size of the training data.

As the contribution for the promotion process in PT XYZ, the best model to create a personalized list of recommended products for each targeted customer is utilized. The recommendation list can be generated by sorting the product by rating prediction. This list then can be integrated into current sales and marketing strategy, for example: creating promotional products, clearing stock or cross-selling. By offering the relevant product, a retailer can reap the benefit of the recommendation system such as personalized marketing, improving customer loyalty, and increasing cross-selling. The

main difference from previous pipeline is the list of recommended products has been adjusted to the target customers. It is more useful for customers since they are likely to receive recommendation of products they need.

#### ACKNOWLEDGMENT

This study is fully supported by Universitas Indonesia under PITTA Grant 2020 with the contract number "PUTI Q2 Nomor: NKB-1475/UN2.RST/HKP.05.00/2020".

#### REFERENCES

- [1] F. Ricci, L. Rokach and B. Shapira, *Recommender System Handbook*, Springer, 2011.
- [2] J. Lu, D. Wu, M. Mao, W. Wang and G. Zhang, "Recommender system application developments: A survey," *Decision Support Systems*, vol. 74, pp. 12-32, 2015.
- [3] S. Sivapalan, A. Sadeghian, H. Rahnama and A. M. Madni, "Recommender Systems in E-Commerce," in *World Automation Congress*, Waikoloa, HI, USA, 2014.
- [4] F. E. Walter, S. Battiston, M. Yildirim and F. Schweitzer, "Moving recommender systems from on-line commerce to retail stores," *Information System and Business Management*, vol. 10, no. 3, pp. 367-393, 2012.
- [5] M. Kang, D.-H. Shin and T. Gong, "The role of personalization, engagement, and trust in online communities," *Information Technology and People*, vol. 29, no. 3, pp. 580-596, 2015.
- [6] T. Donkers, B. Loepp and J. Ziegler, "Explaining Recommendations by Means of User Reviews," in *Workshop on Explainable Smart Systems (ExSS) 2018*, 2018.
- [7] J. Hanke, M. Hauser, A. Dürr and F. Thiesse, "Redefining the Offline Retail Experience: Designing Product Recommendation Systems for Fashion Stores," in *26th European Conference on Information Systems*, Portsmouth, United Kingdom, 2018.
- [8] H. Hwangbo, Y. S. Kim and K. J. Cha, "Recommendation system development for fashion retail e-commerce," *Electronic Commerce Research and Applications*, vol. 28, pp. 94-101, 2018.
- [9] F. Rodrigues and B. Ferreira, "Product Recommendation based on Shared Customer's Behaviour," *Procedia Computer Science*, vol. 100, pp. 136-146, 2016.
- [10] Z. Shi, Z. Wen and J. Xia, "An Intelligent Recommendation System based on Customer Segmentation," *International Journal of Research in Business Studies and Management*, vol. 2, no. 11, pp. 78-90, 2015.
- [11] M. Singh, "Scalability and sparsity issues in recommender datasets: a survey," *Knowledge and Information Systems*, pp. 1-43, 2018.
- [12] G. K. Kishore and D. S. Babu, "Recommender System based on Customer Behaviour for Retail Stores," *IOSR Journal of Computer Engineering*, vol. 19, no. 3, pp. 06-17, 2017.
- [13] S. M. Rezaenia and R. Rahmani, "Recommender system based on customer segmentation (RSCS)," vol. 45, no. 6, pp. 946-961, 2016.
- [14] Y. Shi, M. Larson and A. Hanjalic, "Collaborative filtering beyond the user-item matrix: A survey of the state of the art and future challenges," *ACM Computing Surveys*, vol. 47, no. 1, 2014.
- [15] J. Bobadilla, F. Ortega, A. Hernando and A. Gutiérrez, "Recommender systems survey," *Knowledge-Based Systems*, vol. 46, pp. 109-132, 2013.
- [16] C. C. Aggarwal, *An Introduction to Recommender Systems*. In *Recommender System*, Springer, 2016.
- [17] P. H. Aditya, I. Budi and Q. Munajat, "A Comparative Analysis of Memory-based and Model-based Collaborative Filtering on the Implementation of Recommender System for E-commerce in Indonesia : A Case Study PT X," in *International Conference on Advanced Computer Science and Information System*, Malang, Indonesia, 2016.
- [18] J. Park and K. Nam, "Group recommender system for store product placement," *Data Mining and Knowledge Discovery*, vol. 33, pp. 204-229, 2019.
- [19] S. G. Moghaddam and A. Selamat, "A Scalable Collaborative Recommender Algorithm based on User Density-Based Clustering," in *The 3rd International Conference on Data Mining and Intelligent Information Technology Applications*, Macao, China, 2011.
- [20] S. Chan, P. Treleaven and L. Capra, "Continuous hyperparameter optimization for large-scale recommender systems," in *International Conference on Big Data*, Silicon Valley, CA, USA, 2013.
- [21] P. Cremonesi, F. Garzotto, S. Negro, A. V. Papadopoulos and R. Turrin, "Looking for "Good" Recommendations: A Comparative Evaluation of Recommender Systems," in *IFIP Conference on Human-Computer Interaction*, 2011.

# Using Wearable Sensors for Human Activity Recognition in Logistics: A Comparison of Different Feature Sets and Machine Learning Algorithms

Abbas Shah Syed<sup>1</sup>  
University of Louisville,  
USA

Zafi Sherhan Syed<sup>2</sup>  
Mehran University,  
Pakistan

Muhammad Shehram Shah<sup>3</sup>  
RMIT University,  
Australia

Salahuddin Saddar<sup>4</sup>  
Mehran University,  
Pakistan

**Abstract**—The topic of human activity recognition has gained a lot of attention due to its usage for exercise monitoring, smart health and assisted living. Even though the aforementioned domains have received significant interest by researchers, activity recognition for industrial settings has received little attention in comparison. Industry 4.0 involves the assimilation of industrial workers with robots and other equipment used in the industry and necessitates the development of recognition methodologies for activities being performed in industries. In this regard, this paper presents a comparison in performance of various time/frequency domain features and popular machine learning algorithms for use in activity recognition in a logistics scenario. Experiments were conducted on inertial measurement sensor data from the recently released LARA dataset which involved three feature sets being used with four machine learning algorithms; Support Vector Machines, Decision Trees, Random Forests and Extreme Gradient Boost (XGBoost). The best result achieved in the experiments was an average accuracy of 78.61% using the XGBoost classifier while using both time and frequency domain features. This work serves as a baseline for activity recognition in logistics using IMU sensors and enables the development of solutions to support fulfillment of Industry 4.0 goals.

**Keywords**—Human Activity Recognition (HAR); inertial sensors; LARA dataset; smart industry

## I. INTRODUCTION

Human activity recognition (HAR) has been a very popular application target for the development of mobile smart devices as it brings healthcare to the home. HAR involves the determination of activities that a person performs in their daily life such as walking, standing, sitting, jogging, etc. Wearable sensors such as accelerometers, gyroscopes, magnetometers can be worn on the body and collect movement data on a person while they are performing activities. Although, there are various modalities that can be used for activity recognition, such as videos [1] or using environmental sensors (for e.g. pyroelectric infrared sensors [2]), wearable sensors provide the benefit of being nonrestrictive to movement, are cost effective and easy to ‘carry’, thereby making them suitable for use in human activity recognition tasks.

While human activity recognition has attracted wide interest in the domains of smart health, ambient assisted living and more [3], the area of activity recognition in an industrial setting has received much less attention. This, even after the fact that the Industry 4.0 vision involves workers in a factory to be equipped with sensors and other smart devices

to work fully integrated with robots and other devices [4]. It is therefore necessary to develop algorithms that are able to perform activity recognition in such an environment that allows seamless integration of workers with machines in the industry and also facilitate optimization of processes and protocols. Moreover, from a health perspective, activity recognition might help in pointing to work related injuries or avoiding them altogether if sensor signatures are not as expected.

This paper investigates the performance of four machine learning algorithms, Support Vector Machines (SVM), Decision Trees (DT), Random Forests (RF) and Extreme Gradient Boost (XGBoost) for the use of activity recognition in logistics using inertial sensors. Different time and frequency domain features are used in this work which have been utilized for activity recognition previously and their performance is discussed on the recently published LARA dataset [5]. It was found that XGBoost performs the best among the chosen algorithms using both time and frequency domain features. The paper is organized as follows, section II provides a discussion of the literature for activity recognition in the industry, section III provides an introduction to the dataset, section IV elucidates on the methodology used in this work, the way the experiments have been set up, section V provides a discussion of the results obtained for the experiments while a conclusion is presented in section VI and future work is discussed in section VII.

## II. LITERATURE REVIEW

One of the early works for activity recognition in the industry was carried out by Ward. et al. in [6] who use a combination of microphones and an accelerometer to identify different activities performed in assembly tasks in wood shops as a component for augmented reality/computer guided assembly jobs. Data is collected from subjects performing the activities with two microphones attached to them, one on the wrist and the other on the upper arm. The accelerometer is also placed on the wrist. Data from the microphones and accelerometer is segmented and they individually vote for the activity being performed. For classification on the microphone side, the authors compute the FFT, then use Linear Discriminant Analysis for dimensionality reduction followed by the computation of euclidean distance to samples in a training set for determining the current activity. For the accelerometer, they make use of Hidden Markov models to predict the activity being performed. Activity recognition for construction was presented in [7] who use data from five Inertial Measurement

Units (IMUs) placed on a worker (thigh, back of head, calf, upper arm and chest) to identify activities being performed on a construction site targeting increased productivity and reduced risk of injury. To do this, they compute various time and frequency domain features on the accelerometer and gyroscope measurements from the IMUs and after performing feature selection, compare three different ML algorithms, SVM, KNN and C4.5 decision trees for classification. They find that SVM performed the best from the models considered. The authors in [8] target behavioral modeling for assembly line workers using accelerometer data to enable task performance analysis as well as development of computer guided task instruction systems. An accelerometer is placed on a workers wrist and a two step recognition process is used, first windowed data is classified using a KNN classifier which is then passed to a state machine to identify completed tasks from the activity sequences.

In [9], human activity recognition has been performed using convolutional neural networks (CNN) on accelerometer data from the Skoda dataset[10]. Their network consists of one convolutional and pooling layer for each accelerometer signal axis followed by two fully connected layers and a softmax layer for classification. They are able to achieve an accuracy of 88.19% using their CNN. Targeting process optimization, the authors in [11] present a dataset which consists of triaxial accelerometer, gyroscope and magnetometer data from three IMU sensors of subjects performing activities in a picking process. One IMU is placed on the upper chest while the other two are placed on the right and the left wrist respectively. The authors also use this data to perform classification between activities using statistical features and three different ML algorithms, SVM, Naive Bayes and Random Forest from which Random Forest performed the best. In [12], the authors use a convolutional neural network on inertial measurement sensor data on the dataset in [11] to enable optimization in regard to Industry 4.0. Due to class imbalance, the authors use data augmentation and then pass the IMU data for each activity to a CNN with four convolutional layers, two pooling layers, one fully connected layer and a softmax layer for classification. They compare the performance of their CNN with a baseline determined using statistical features and three ML algorithms, SVM, Naive Bayes and Random Forests. The CNN outperforms the other methods by achieving an accuracy of 73.9% as the best result. Another approach which employs deep learning for human activity recognition in an industrial context is proposed in [4] who use CNNs with accelerometer data to differentiate between different activities from the dataset provided in [13]. They compare the performance of multiple preprocessing methods (raw data, spectrogram and its variants) for use with a CNN for classification. They achieve the best results when using raw values.

The authors in [14] make use of semantic representations from motion data collected in a manual picking scenario to perform human activity recognition. They use the MoCAP dataset which consists of a multichannel time series of pose information recorded by 38 cameras. This data is labeled with three different attribute representations, two by experts and one by a nonexpert, and is passed to the convolutional neural network architecture described in [15]. The attribute representations differ in the granularity of the sequences used to describe the picking process, for e.g. representation 1 uses

less attributes described for the picking process compared to representation 2. The CNN achieves a higher accuracy for representation 1 compared to the other two representations at 75%. The authors do note that the representation process is subjective as annotations are expert dependent.

As observed from the literature discussed, activity recognition in the industry is important for achieving Industry 4.0 goals of optimization, computer guided worker instruction, increased productivity and also to enable a safer work environment. However research in this direction has not been up to speed with activity recognition for other domains (smart health, assisted living etc) due to the absence of a large publicly available dataset. Fortunately, the recent introduction of the publicly available LARa dataset which contains recordings of activities performed in a logistics scenario opens up various opportunities for research in this domain. This paper uses the LARa dataset to perform activity recognition in logistics using IMU data in this work.

### III. DATASET

The LARa dataset is a novel dataset that presents multi-modal data for developing algorithms for activity recognition for logistics activity recognition. The dataset has been provided by the 'Innovationlab Hybrid Services in Logistics' at TU Dortmund and follows up from their previous research in this area [16]. The dataset consists of 14 individuals performing three different tasks in a logistics scenario, two related to picking and one related to packing. Each of these activities were recorded using an Optical Marker-based Motion Capture (OMoCap) system which measures movements of the participants as markers, RGB camera to capture videos of the participants and inertial measurement units to track participant movements while performing their activities. There are a total of 758 minutes of data in the LARa dataset which have been annotated for eight intra-activities which in certain sequences constitute the three tasks performed. In addition to this, they also provide 19 binary semantic annotations called attributes for the three scenarios which describe intra-activities in a different manner too. The activities annotated include standing, walking, cart (participant is walking with the cart), handling upwards (participant has atleast one hand raised upward to shoulder height), handling centered (participant can handle things without bending, lifting their arms or needing to kneel), handling downwards (participant has hands below his knees while kneeling or otherwise), synchronization (waving motion before each recording) and a set of samples which were unrecognizable by the annotators and have been marked as *None*.

This dataset is the first of its kind in that it provides an opportunity for researchers to develop automated algorithms for recognizing activities in the context of logistics operations as a public dataset. This work focuses on utilizing data from inertial measurement units for logistics activity recognition. Three types of inertial measurement units were used for recording this data and recordings from seven subjects is contained within the dataset. These units are used to measure accelerometer and gyroscope sensor readings on both the legs, arms and the chest/mid-body. Table I presents a summary of IMU measurements in the dataset. There are a total of 14 trials of scenario 1, 99 trials of scenario 2 and 95 trials of scenario



3. For a detail on the sequences of each activity that form the three tasks the reader is referred to [5].

TABLE I. SUMMARY OF IMU MEASUREMENTS IN THE LARA DATASET

| Subject ID   | Gender | Age | Scenario 1 | Scenario 2 | Scenario 3 |
|--------------|--------|-----|------------|------------|------------|
| S07          | M      | 23  | 2          | 13         | 14         |
| S08          | F      | 51  | 2          | 13         | 14         |
| S09          | M      | 35  | 2          | 14         | 13         |
| S10          | M      | 49  | 2          | 13         | 12         |
| S11          | F      | 47  | 2          | 12         | 0          |
| S12          | F      | 23  | 0          | 6          | 14         |
| S13          | F      | 25  | 2          | 14         | 14         |
| S14          | M      | 54  | 2          | 14         | 14         |
| <b>Total</b> |        |     | 14         | 99         | 95         |

#### IV. METHODOLOGY

The methodology for this study follows a typical ML pipeline as shown in Fig. 1 where the first stage is preprocessing (windowing in this case), followed by feature extraction and then the use of ML to perform classification. Each of these steps are discussed in detail in subsequent sections. Three different tests were performed with the four machine learning algorithms, first using both frequency and time domain features, the second using time domain features only, and the third using only frequency domain features.

##### A. Preprocessing Stage

The sensor measurements from the inertial measurement units in the dataset are recorded with a sampling frequency of 100 Hz. Data collection takes place as the participants carry out each of the three scenarios. Moreover, each triaxial accelerometer and gyroscope reading has been annotated as belonging to one of the eight activity classes described in section III.

Since the data consists of tasks determined by the performance of a number of sequential activities, during preprocessing contiguous segments are extracted from accelerometer and gyroscope measurements having the same label. For e.g. for scenario 1, Fig. 2 depicts the sequence of events in the carrying out of this task by subject seven (trial number 1) as extracted from the labeled activities of IMU data. The figure also shows the business model of the activity performed [5] for context. The windowing process extracts the contiguous samples for each of these activities, in this case there were 28 windows extracted for the intra-activities that constitute the task in scenario 1 in terms of activities: standing (0), walking (1), cart (2), handling upwards (3), handling centered (4), handling downwards (5), synchronization (6) and *None* (7). After performing ‘windowing’ for all samples, segments from annotations belonging to the categories *None* and *synchronization* were removed as they are not intended to be taken in to account [5]. The remaining segments for six activities were passed on to the feature extraction stage.

##### B. Feature Computation

Feature extraction is the process of representing data in a meaningful format so as to make it more adaptable for use in computational processes such as regression, classification or other forms of decision making. The field of human activity recognition using wearable sensor data has utilized various

types of feature extraction mechanisms such as wavelets [17], time and frequency domain computations [18] and also CNNs [19]. In this work, different time and frequency domain features have been used to represent the information contained in the extracted windows of the accelerometer and gyroscope sensors. This choice is motivated by the works of [20], [21] who achieve very good results for human activity recognition, we compute twelve time domain and four frequency domain features in the feature extraction process. These are listed in Table II. These parameters are computed for each of the segments extracted in the preprocessing stage. Moreover, for each sensor in each segment, feature values are normalized across the three axes to ensure that different scales/units of the sensors do not affect classification performance.

TABLE II. FEATURES COMPUTED FROM IMU DATA

| Domain    | Feature set                   |
|-----------|-------------------------------|
| Time      | Variance                      |
|           | Mean                          |
|           | Median                        |
|           | Standard Deviation            |
|           | Maximum                       |
|           | Minimum                       |
|           | Delta                         |
|           | 25th Percentile               |
|           | 75th Percentile               |
|           | Interquartile range           |
|           | Kurtosis                      |
|           | Skew                          |
| Frequency | Power Spectral Density Mean   |
|           | Power Spectral Density Median |
|           | Power Spectral Density RMS    |
|           | Power Spectral Entropy        |

##### C. Classification

For classification, five different algorithms have been chosen to test their efficacy for HAR in the logistics scenario. The algorithms chosen are SVM, XGBoost, Random Forests and Decision Tree. Tree based ensemble schemes have been chosen as they have been useful in [22] for fall detection purposes. Moreover, SVM has been successfully used in [23] to perform human activity recognition using inertial sensors. Each of the chosen algorithms were tuned by performing a grid parameter search. The details of the grid search for each of the tested algorithm is given in Table III.

TABLE III. PARAMETER VALUES FOR TUNING ML CLASSIFIERS

| Algorithm     | Parameter            | Grid Search values                         |
|---------------|----------------------|--|
| SVM           | Kernel               | Linear, RBF, Sigmoid, Poly                 |
|               | C                    | 0.1,0.2,0.4,0.5,1.0,1.5,1.8,2.0,2.5,3.0,10 |
| Random Forest | Number of estimators | 20, 50, 100                                |
|               | Criterion            | Gini, Entropy                              |
|               | Max Depth            | 4,5,6,7,8                                  |
| Decision Tree | Number of estimators | 20, 50, 100, 200                           |
|               | Max Depth            | 4,5,6,7,8                                  |
|               | Criterion            | Gini, Entropy                              |
|               | Min Child Weight     | 1,5,10                                     |
| XGBoost       | Gamma                | 0.5,1,1.5,2,5                              |
|               | Max Depth            | 4,5,6,7,8                                  |
|               | Number of estimators | 20, 50, 100                                |

#### V. EXPERIMENTATION, RESULTS, AND DISCUSSION

In order to test the efficacy of the features used to represent IMU data for logistics activity recognition, we perform three different experiments. These experiments were conducted to

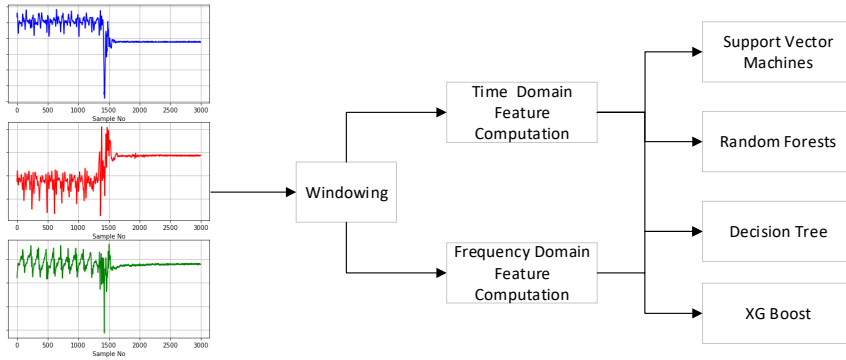


Fig. 1. Methodology for HAR Recognition in Logistics

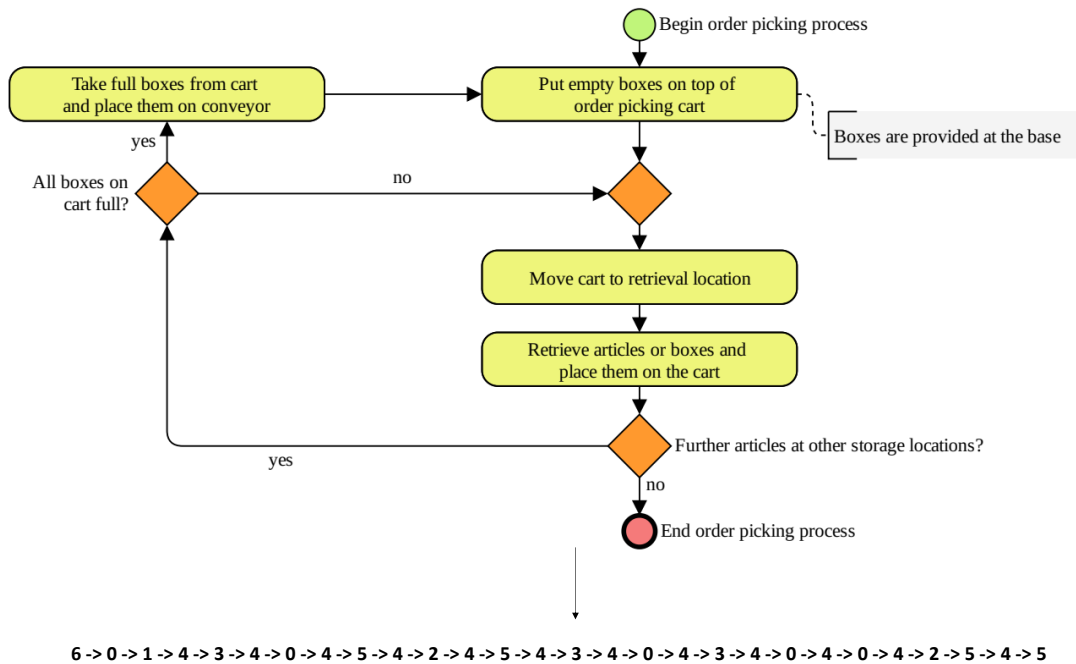


Fig. 2. Translation of Scenario 1 (LARA Dataset) to a Sequence of Activities for Recognition

ascertain the most appropriate algorithm to be used for the targeted task as well as the most appropriate feature set to be used. It is important to note that sensor data from all five locations on the body were utilized in these experiments. The implementations of support vector machine (SVM), and random forest (RF), decision tree (DT) and Extreme Gradient Boost (XGBoost) present in the scikit-learn toolkit<sup>1</sup> have been used. Training is performed with a train-test split of 75-25 using five fold cross validation.

#### A. Experiment with Time and Frequency Features

In this experiment both frequency and time domain features were used as input to the classification algorithms. Tests were conducted using both accelerometer and gyroscope values. For the case of using features from both the time and frequency domains, a total of sixteen parameters were computed for each sensor modality present as given in Table II. The length of the feature vector for this experiment was 480.

Table IV presents the results from this experiment. This work reports on the accuracies achieved for activity classes. From the table, it can be observed that the best performing algorithm is XGBoost with a mean accuracy of 78.61%, followed by 76.67% for SVM. The mean accuracy achieved by

<sup>1</sup><https://scikit-learn.org>

TABLE IV. SUMMARY OF RESULTS: TIME AND FREQUENCY DOMAIN FEATURES

| Activity         | Algorithm (Accuracy %) |          |       |                 |
|------------------|------------------------|----------|-------|-----------------|
|                  | SVM                    | RF       | DT    | XGB             |
| Stand            | 78                     | 74       | 60    | 80              |
| Walk             | 83.33                  | 66.66    | 50    | 75              |
| Cart             | 66.66                  | 50       | 50    | 83.33           |
| Hand Up          | 75                     | 56.25    | 37.5  | 75              |
| Hand Center      | 73.07                  | 80.76    | 66.66 | 74.35           |
| Hand Down        | 84                     | 84       | 68    | 84              |
| Average Accuracy | 76.67667               | 68.61167 | 55.36 | <b>78.61333</b> |

using RF and DT are 68.61% and 55.36% respectively which are significantly less than the best performing results. This result is in agreement with the works of [22] who also find gradient boosted trees to work well for use in human activity recognition applications.

### B. Experiment with Time Domain Features Only

The second experiment consisted of using only time domain features for logistics activity recognition. A total of twelve time domain features were computed from the sensor readings. This resulted in a total of 360 features being used for the classification stage. The results for each classifier are shown in Table V. Taking from average accuracy, SVM, RF and XGBoost provide similar performances with mean accuracies of around 72% with DT performing very poorly (a mean accuracy of only 39.61%). Another point to note from Table V is that all classifiers have performed poorly for the activity *Cart*. This indicates that the time domain features used in this work are unable to represent this activity well.

TABLE V. SUMMARY OF RESULTS: TIME DOMAIN FEATURES ONLY

| Activity         | Algorithm (Accuracy %) |          |          |               |
|------------------|------------------------|----------|----------|---------------|
|                  | SVM                    | RF       | DT       | XGB           |
| Stand            | 66                     | 74       | 52       | 76            |
| Walk             | 91.66                  | 83.3     | 16.66    | 75            |
| Cart             | 50                     | 50       | 0        | 50            |
| Hand Up          | 75                     | 56.25    | 37.5     | 75            |
| Hand Center      | 70.51                  | 80.76    | 61.53    | 73.07         |
| Hand Down        | 80                     | 88       | 68       | 88            |
| Average Accuracy | 72.195                 | 72.05167 | 39.28167 | <b>72.845</b> |

### C. Experiment with Frequency Domain Features Only

In this experiment, only frequency domain features were used for activity recognition. This resulted in a total of four features being computed for each sensor modality present. This resulted in a feature vector size of 90 for the classifiers. The results of the experiment are depicted in Table VI. The results indicate that the XGBoost performs the best among all the classifiers tested with a mean accuracy of 56.11% with the SVM achieving a mean accuracy of 54.495%. These results are a significant reduction from the results of experiments 1 and 2. Another point to note here is that the activity *Hand up* has the least individual performance followed by the activity *Cart*. This indicates that the frequency domain based features used here might not be enough to appropriately represent these activities for logistics activity recognition.

From the three experiments conducted, the most suitable combination of feature set and classifier is a combination of time and frequency domain features along with an XGBoost

TABLE VI. SUMMARY OF RESULTS: FREQUENCY DOMAIN FEATURES ONLY

| Activity         | Algorithm (Accuracy %) |       |          |                 |
|------------------|------------------------|-------|----------|-----------------|
|                  | SVM                    | RF    | DT       | XGB             |
| Stand            | 54                     | 46    | 70       | 48              |
| Walk             | 75                     | 66.66 | 58.33    | 75              |
| Cart             | 33.33                  | 6.66  | 16.66    | 50              |
| Hand Up          | 25                     | 0     | 6.25     | 18.75           |
| Hand Center      | 75.64                  | 76.92 | 48       | 76.923          |
| Hand Down        | 64                     | 68    | 28       | 68              |
| Average Accuracy | 54.495                 | 44.04 | 37.87333 | <b>56.11217</b> |

classifier. Table VII lists the precision, recall and F1 scores for the individual activities when using the best performing combination for activity recognition in a logistics scenario. It can be observed that the best scores are achieved for the activities *Cart* and *Hand Down* and the least scores are for the activities *Walk* and *Hand Center*.

TABLE VII. SUMMARY OF BEST RESULTS: TIME + FREQUENCY DOMAIN FEATURES AND XGBOOST CLASSIFIER

| Activity    | Accuracy | Precision | Recall   | F1       |
|-------------|----------|-----------|----------|----------|
| Stand       | 80       | 0.8       | 0.645161 | 0.714286 |
| Walk        | 75       | 0.75      | 0.9      | 0.818182 |
| Cart        | 83.33    | 0.833333  | 1        | 0.909091 |
| Hand up     | 75       | 0.75      | 0.666667 | 0.705882 |
| Hand Center | 74.35    | 0.74359   | 0.816901 | 0.778523 |
| Hand Down   | 84       | 0.84      | 1        | 0.913043 |

## VI. CONCLUSION

In this study the problem of activity recognition in a logistic scenario is addressed. In this regard, this research makes use of IMU sensor data from the novel LARa dataset which contains OMOCap, video and IMU sensor recordings for individuals performing three different scenarios in an industrial setting. The experiments conducted in this work make use of several time and frequency domain based features which have been used in activity recognition/fall detection using wearable sensors along with popular machine learning frameworks which have also proved to perform well in such applications. From the conducted experiments, the XGBoost algorithm performed the best when used with the considered time and frequency domain features and the highest mean accuracy achieved was 78.61%.

## VII. FUTURE WORK

This work establishes a baseline for logistics human activity recognition using inertial sensors on a novel dataset and can be used for optimization of logistics operations. Future work in this area will include using Deep Learning algorithms such as Convolutional Neural Networks and Recurrent Neural Networks which are able to capture information in sensor readings more intricately. Another scope of research is the use of sensor fusion of OMOCap and/or Video data with wearable sensor data for classification between the different activities for such logistics applications.

## REFERENCES

- [1] T. Singh and D. K. Vishwakarma, "Human activity recognition in video benchmarks: A survey," in *Advances in Signal Processing and Communication*. Springer, 2019, pp. 247–259.

- [2] X. Luo, Q. Guan, H. Tan, L. Gao, Z. Wang, and X. Luo, "Simultaneous indoor tracking and activity recognition using pyroelectric infrared sensors," *Sensors*, vol. 17, no. 8, p. 1738, 2017.
- [3] F. Demrozi, G. Pravadelli, A. Bihorac, and P. Rashidi, "Human activity recognition using inertial, physiological and environmental sensors: a comprehensive survey," *arXiv preprint arXiv:2004.08821*, 2020.
- [4] X. Zheng, M. Wang, and J. Ordieres-Meré, "Comparison of data preprocessing approaches for applying deep learning to human activity recognition in the context of industry 4.0," *Sensors*, vol. 18, no. 7, p. 2146, 2018.
- [5] F. Niemann, C. Reining, F. M. Rueda, N. R. Nair, J. A. Steffens, G. A. Fink, and M. t. Hompel, "Lara: Creating a dataset for human activity recognition in logistics using semantic attributes," *Sensors*, vol. 20, no. 15, p. 4083, 2020.
- [6] J. A. Ward, P. Lukowicz, G. Troster, and T. E. Starner, "Activity recognition of assembly tasks using body-worn microphones and accelerometers," *IEEE transactions on pattern analysis and machine intelligence*, vol. 28, no. 10, pp. 1553–1567, 2006.
- [7] J. Zhao and E. Obonyo, "Towards a data-driven approach to injury prevention in construction," in *Workshop of the European Group for Intelligent Computing in Engineering*. Springer, 2018, pp. 385–411.
- [8] H. Koskimäki, V. Huikari, P. Siirtola, and J. Röning, "Behavior modeling in industrial assembly lines using a wrist-worn inertial measurement unit," *Journal of Ambient Intelligence and Humanized Computing*, vol. 4, no. 2, pp. 187–194, 2013.
- [9] M. Zeng, L. T. Nguyen, B. Yu, O. J. Mengshoel, J. Zhu, P. Wu, and J. Zhang, "Convolutional neural networks for human activity recognition using mobile sensors," in *6th International Conference on Mobile Computing, Applications and Services*. IEEE, 2014, pp. 197–205.
- [10] P. Zappi, C. Lombriser, T. Stiefmeier, E. Farella, D. Roggen, L. Benini, and G. Tröster, "Activity recognition from on-body sensors: accuracy-power trade-off by dynamic sensor selection," in *European Conference on Wireless Sensor Networks*. Springer, 2008, pp. 17–33.
- [11] S. Feldhorst, M. Masoudenijad, M. ten Hompel, and G. A. Fink, "Motion classification for analyzing the order picking process using mobile sensors," in *Proc. Int. Conf. Pattern Recognition Applications and Methods*, 2016, pp. 706–713.
- [12] R. Grzeszick, J. M. Lenk, F. M. Rueda, G. A. Fink, S. Feldhorst, and M. ten Hompel, "Deep neural network based human activity recognition for the order picking process," in *Proceedings of the 4th international Workshop on Sensor-based Activity Recognition and Interaction*, 2017, pp. 1–6.
- [13] T. Szttyler and H. Stuckenschmidt, "On-body localization of wearable devices: An investigation of position-aware activity recognition," in *2016 IEEE International Conference on Pervasive Computing and Communications (PerCom)*. IEEE, 2016, pp. 1–9.
- [14] C. Reining, M. Schlangen, L. Hissmann, M. ten Hompel, F. Moya, and G. A. Fink, "Attribute representation for human activity recognition of manual order picking activities," in *Proceedings of the 5th international Workshop on Sensor-based Activity Recognition and Interaction*, 2018, pp. 1–10.
- [15] F. M. Rueda and G. A. Fink, "Learning attribute representation for human activity recognition," in *2018 24th International Conference on Pattern Recognition (ICPR)*. IEEE, 2018, pp. 523–528.
- [16] C. Reining, F. Niemann, F. Moya Rueda, G. A. Fink, and M. ten Hompel, "Human activity recognition for production and logistics—a systematic literature review," *Information*, vol. 10, no. 8, p. 245, 2019.
- [17] G. Bhat, R. Deb, V. V. Chaurasia, H. Shill, and U. Y. Ogras, "Online human activity recognition using low-power wearable devices," in *2018 IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*. IEEE, 2018, pp. 1–8.
- [18] P. Sarcevic, S. Pletl, and Z. Kincses, "Comparison of time-and frequency-domain features for movement classification using data from wrist-worn sensors," in *2017 IEEE 15th International Symposium on Intelligent Systems and Informatics (SISY)*. IEEE, 2017, pp. 000 261–000 266.
- [19] A. Ignatov, "Real-time human activity recognition from accelerometer data using convolutional neural networks," *Applied Soft Computing*, vol. 62, pp. 915–922, 2018.
- [20] A. A. Badawi, A. Al-Kabbany, and H. A. Shaban, "Sensor type, axis, and position-based fusion and feature selection for multimodal human daily activity recognition in wearable body sensor networks," *Journal of Healthcare Engineering*, vol. 2020, 2020.
- [21] S. Rosati, G. Balestra, and M. Knafnitz, "Comparison of different sets of features for human activity recognition by wearable sensors," *Sensors*, vol. 18, no. 12, p. 4189, 2018.
- [22] N. Zurbuchen, P. Bruegger, and A. Wilde, "A comparison of machine learning algorithms for fall detection using wearable sensors," in *2020 International Conference on Artificial Intelligence in Information and Communication (ICAIIIC)*. IEEE, 2020, pp. 427–431.
- [23] I. Bustoni, I. Hidayatulloh, A. Ningtyas, A. Purwaningsih, and S. Azhari, "Classification methods performance on human activity recognition," in *Journal of Physics: Conference Series*, vol. 1456, no. 1. IOP Publishing, 2020, p. 012027.

# Real-Time Healthcare Monitoring System using Online Machine Learning and Spark Streaming

Fawzya Hassan<sup>1</sup>, Masoud E. Shaheen<sup>2</sup>  
Department of Computer Science,  
Faculty of Computers and Information,  
Fayoum University, Egypt

Radhya Sahal<sup>3</sup>  
Faculty of Computer Science and Engineering  
Hodeidah University - Yemen

**Abstract**—The real-time monitoring and tracking systems play a critical role in the healthcare field. Wearable medical devices with sensors, mobile applications, and health cloud have continuously generated an enormous amount of data, often called streaming big data. Due to the higher speed of the streaming data, it is difficult to ingest, process, and analyze such huge data in real-time to make real-time actions in case of emergencies. Using traditional methods that are inadequate and time-consuming. Therefore, there is a significant need for real-time big data stream processing to guarantee an effective and scalable solution. So, we proposed a new system for online prediction to predict health status using Spark streaming framework. The proposed system focuses on applying streaming machine learning models (i.e. streaming linear regression with SGD) on streaming health data events ingested to spark streaming through Kafka topics. The experimental results are done on the historical medical datasets (i.e. diabetes dataset, heart disease dataset, and breast cancer dataset) and generated dataset which is simulated to wearable medical sensors. The historical datasets have shown that the accuracy improvement ratio obtained using the diabetes disease dataset is the highest one with respect to the other two datasets with an accuracy of 81%. For generated datasets, the online prediction system has achieved accuracy with 98% at 5 seconds window size. Beyond this, the experimental results have proved that the online prediction system can online learn and update the model according to the new data arrival and window size.

**Keywords**—Online machine learning; streaming data; Apache Spark; Apache Kafka; spark streaming machine learning

## I. INTRODUCTION

Nowadays, the era has been described as the era of big data where all data is digitalized and becomes of great importance in such beautiful fields. An enormous amount of data has been gathered and produced from different sectors, many sources like sensor networks, wireless machines, mobile applications, and from different fields [1]. Especially in the healthcare field, a big data collected in real-time by a remote sensing device, Wearable Medical Devices, and other data gathering tools, which produce new challenges that focus on data size and the fast growth rate of these data [2]. One of the most important challenges in data analytics is dealing with inappropriate technological tools to store, process, visualize, and extract knowledge with large and varied data types. In addition, exploring a new way to obtain valuable information for many users. However, the digital record of the patient's medical history is the primary health care data as it is obtained from various types of health care data sources in both clinical and non-clinical settings. Occasionally, these digital data are not available for research.

In the past, communication between doctors and patients was done by bounded visits, tele, and text communications. Consequently, doctors and hospitals could observe their patients' health constantly, and even more, they couldn't make recommendations accordingly. Thanks to IoT devices represented in wearable medical devices like heart rate monitoring cuffs, blood pressure, glucometer, etc. These devices can determine the space needed for each device and also determine the degree of interaction of people with these devices to provide health care solutions through continuous tracking of health conditions and giving patients access to personalized attention. However, the wearable medical devices continuously generate data; the amount of data stored and processed becomes a vital problem in real life. Furthermore, lately, many citizens and the elderly suffer from chronic diseases, and looking at the disadvantages of traditional health services is a very important thing. Therefore, the medical IoT is used to observe and take the required actions in real time for emergencies, like people with heart disease and diabetes[3]. Consequently, processing these enormous data generated by the sensors and implementing procedures in real time in critical cases is a very important challenge contribution.

Therefore, proposing a system that can handle big data faces three main challenges: First, collecting data from distributed systems is a complicated process due to the enormous amount of data. Second, storing that big and heterogeneous data is a major problem; therefore, the need for a big data storage system with efficient and effective performance is very essential. Third, this challenge relates to data analysis, especially big data processing in real time, including modeling, visualization, prediction, and optimization. Therefore, these challenges require new processing systems to deal with heterogeneous data or big data processing in real time. We will address the challenge related to data analytic in real time in the healthcare domain.

Recently, big data streaming computing has been used as an effective role in big data analytics to investigate the importance of big data in real-time healthcare. For example, a real-time system for flu and cancer monitoring is produced by applying twitter data mining in [4]. A model for real-time medical big data analysis is introduced in [5]. The model is performed by applying Spark Streaming [6] and Apache Kafka<sup>1</sup> using a stream of healthcare data. In [7] a real-time health status prediction system is proposed, this work focuses on applying machine learning especially Decision Tree algorithm on data

<sup>1</sup><https://kafka.apache.org>

streams received from socket streams using Apache Spark<sup>2</sup>. However, researchers and developers face problems due to distributed data sources for healthcare (i.e., distributed queuing management technologies) like Kafka, and RabbitMQ<sup>3</sup>. The aggregated health-based streaming data is analyzed using big data platforms for streaming processing such as Apache Spark, Apache Storm<sup>4</sup>, Apache Samza<sup>5</sup>, Apache Flink<sup>6</sup>, relational databases, analytics systems, and other search systems. Most recent research relies on machine learning, but streaming big data that needs to apply machine learning in real time is not handled. In particular, the previous studies just applied traditional machine learning algorithms to analyze and predict health status for patients using historical data. These studies focused mainly on Hadoop, the batch-oriented computing system. For this reason, the important challenge is applying machine learning to streaming data because conventional machine learning systems are not effective in dealing with real-time streaming data.

To the best of our knowledge, no study has been done oriented to online predicting disease in real-time. This motivates us to introduce a new online prediction system that can predict health status in real time, using Kafka data streaming, Spark Streaming, and Spark MLlib. Consequently, the research in this paper works on three important case studies: Pima Indians diabetes, Cleveland heart disease and breast cancer Coimbra because a large percentage of people have been injured with these diseases, and leading to death. So the online prediction for these diseases can decrease the mortality rate. The proposed system is used to achieve high accuracy using streaming data, i.e. Kafka producers produce stream messages of data continuously and then apply an online model to the online prediction in real time by classifying the stream of data as containing disease indication or not. The proposed system has four phases: 1) Data ingestions from the input stream within the data source; 2) Streaming Process Pipeline; 3) Online Prediction; and 4) Output Stream. The paper contributions can be summarized as:

- Developing an online prediction system to the possibility of disease using streaming data from historical medical datasets (i.e., diabetes dataset, heart disease dataset, and breast cancer dataset) and from real-time data in the form of wearable medical devices.
- Generating streaming data from simulated to wearable medical sensors and then capturing a stream of data by Kafka topic provides name to the various diseases.
- Applying StreamingLinearRegressionWithSGD algorithm to classify streaming data.
- Evaluating the result for historical medical dataset and simulated wearable sensor generator to compare the accuracy for different window sizes.

The remainder of this paper is organized as follows: The related work is presented in Section II. The proposed system of online prediction is introduced in Section III. The experimental

results are discussed in Section IV. Finally, conclusions are presented in Section V.

## II. RELATED WORKS

In last years, big data analytics concerning healthcare analytics become an important issue for many research areas like machine learning, data mining with data from healthcare as well as the available information from inside hospitals. The progress of the data collection process is due to the huge development in technology in the field of health care, where data is collected through three main stages of the flow of digital data resulting from the patient's clinical records, health research records, and organization operations records [3].

In [8] discuss an overview about the healthcare data sources. This study analyzes health care data played a very important role in many systems such as disease prediction, methods of prevention, medical guidance, and urgent medical decision-making in order to improve the standard of health care, reduce costs and increase efficiency. Also, Archenaa, J. and Anita [9] explain how can we apply Apache Spark to apply healthcare analytics through in-memory computations that can handle a large amount of structured, unstructured patient data and patients streaming data from their social network activities.

Multiple machine learning models on healthcare data using spark have been introduced in previous studies. For example In [10] Introduced a real time health status prediction system that uses spark machine learning streaming Big Data, The system was tested on the user tweets with his health attributes and the system receives these tweets, extracts the parameters and perform decision tree algorithm for user's health status predict, and finally send a direct message to him/her to take the appropriate action.

Moreover, Alotaibi, Shoayee, et al. [11] proposed a Sehaa which is a big data analytics tool for Arabic healthcare Twitter data in the Kingdom of Saudi Arabia (KSA). The system uses two different ML algorithms, including Naive Bayes, Logistic Regression, and applying multiple feature extraction methods to detect various diseases in the KSA. In [12] the system that can able to predict real-time heart disease based on Apache Spark that applied machine learning on streaming data by using memory computations are explained. The system are developed with two main stages, the first one is streaming processing which use spark MLlib with Spark streaming by applying classification algorithms on data to heart disease prediction. The second stage is data storage and visualization which uses Apache Cassandra to store a huge volume of generated data.

Most of this studies relies on specific healthcare data sources and applying processing on the offline system, but in reality, the sources of health data are various and constantly produce various data of different sizes. In addition, real-time healthcare analytical involves real-time streaming data processing, streaming machine learning algorithms, and analyzing real time to build an online electronic system to deal with the stream of healthcare data. therefore, we developed an online prediction healthcare system for streaming data coming from IoT devices to predict health status for the patient in real time.

<sup>2</sup><https://spark.apache.org>

<sup>3</sup><https://www.rabbitmq.com>

<sup>4</sup><http://storm.apache.org/>

<sup>5</sup><http://samza.apache.org>

<sup>6</sup><https://flink.apache.org>

### III. THE ONLINE PREDICTION SYSTEM

The online health status prediction system is a data analytic monitoring model that uses Kafka streaming and Spark streaming. The architecture of the proposed system consists of four phases, as shown in Fig. 1. In the first phase, data ingestions from different data sources; Kafka producers continuously generate a stream of data messages that are captured by Kafka streaming from different topics name correlated to various diseases. Second, Streaming Process Pipeline in which Spark streaming receives a stream of medical data with the attributes which related to each disease and then applies a streaming machine learning model to predict health status. Third, Online Prediction receives batches of input data and responsible for online learning and updating our deployed model according to the new data arrival. At the final phase, the Output predicted results to data sinks in which the output stream sends back again to Kafka to be consumed by other data sinks such as web service, alarm systems, dashboard, mobile application, and hospital social networks.

#### A. Data Ingestions

For analyzing healthcare data, Apache Kafka has been used as the tool for ingestion of the individual's health data from distributed sources of historical data and real-time data. The historical data is ingested to the proposed system using multiple disease datasets; diabetes disease, heart disease, and breast cancer disease, each dataset uses different Kafka's topics named "diabetes\_disease\_dataste," "heart\_disease\_dataset," and "breast\_cancer\_dataset", respectively. The real-time dataset is collected by a simulated wearable sensor generator, which generates diabetes disease data on Kafka's topic named "Diabetes\_Generator\_Dataset." It is challenging to manage this data with Spark itself; therefore, Kafka is designed specifically for streaming data management[1]. Hence, it has been integrated into our system.

#### B. Streaming Data Processing over Apache Spark Streaming

The online prediction system is a data analytic model that uses a spark streaming machine learning library (MLlib) i.e. streaming linear regression with the SGD algorithm, which requires training data. Spark streaming receives a stream of records from historical data or real-time data, that is used as training data. Spark streaming data processing uses streaming computation by applying data decomposition ( see Fig. 1). In the first, spark streaming transforms the input data stream into a Discretized stream (DS). After that the DSs are converted into Resilient Distributed Datasets (RDD). Therefore, it is urgent to transform and perform the RDD to be able to fulfill streaming processing. Furthermore, spark streaming can process the RDD data based on MLlib; the online prediction system uses the StreamingLinearRegressionWithSGD algorithm, which will be described in the following section.

1) *Spark MLlib Streaming linear regression algorithm:* Streaming linear regression uses Spark Streaming technology to train or predict a linear regression model based on streaming data. It applies SGD for each new batch of data coming from DStream to update the model. Each batch of data is expected to be an RDD of LabeledPoints. The spark streaming linear regression algorithm takes four parameters: the number of

iteration, step size (learning rate), mini-batch fraction time, and initial weights vector. The number of iteration is needed to finish the gradient descent. The learning rate determines how slow or fast the algorithm can update the optimal regression coefficients. The initial weight vector must be provided; the default initial weight is 0.0. The mini-batch fraction time is used to set the batch time. The batch time parameter estimates the window time for spark streaming. Spark Streaming linear regression has a latency of many seconds, because of mini-batch time. Conversely, this mini-batch time efficiently ensures that each stream data will be processed exactly once [13].

#### C. Online Prediction

The online prediction phase is divided into two main components; the deployed model and the online learning mode. The online learning model takes the batches of input to train the model then send the training queue to the deployed model to learn and predict the result using the StreamingLinearRegressionWithSGD algorithm. The online learning algorithm takes the first batch result model and then picks up one to one observation from the training queue and recalibrates the weights on each input parameter. The deployed model is continuously learning, and it updates parameters for each batch result, which is close to "learning-on-the-fly". It helps to learn variations in distribution as quickly as possible and improve accuracy in many cases. The online prediction is found to be relatively faster than their batch equivalent methods.

#### D. Output Predicted Results to Data Sinks

In this phase, the output stream is sent back again to Kafka to be consumed by other data sinks. As the proposed system suppose to work in real-time, these data sinks could be any online data consumer within the real-world healthcare application in the industrial setting such as 1) web service for healthcare monitoring [14], [15], 2) alarm systems which connected to the hospital emergency department [16], [17], 3) real-time dashboard [18], 4) hospital medical records which stored in big data cloud storage (e.g., HDFS, MongoDB) [19], [20], and 5) hospital social networks considering patients' privacy.

## IV. EXPERIMENTAL RESULTS AND DISCUSSION

In this section, the experimental evaluation of the proposed online healthcare monitoring system is presented, and the results are discussed. Experiments were conducted on two different data sources of healthcare 1) historical medical data and real-time data. The historical data is ingested into the proposed system using three medical datasets chosen from UCI machine learning repository; Pima Indians diabetes [21], Cleveland heart disease [22] and breast cancer Coimbra [23]. The real-time dataset is collected by a simulated wearable sensor generator which generates diabetes disease data. Further details will be extensively discussed in the next subsection.

#### A. Experimental Setup

The proposed system has been implemented on top of Apache Spark using scala. Our experiments are conducted through Apache Kafka and Spark Streaming for data ingestion and data processing, respectively. The online machine learning

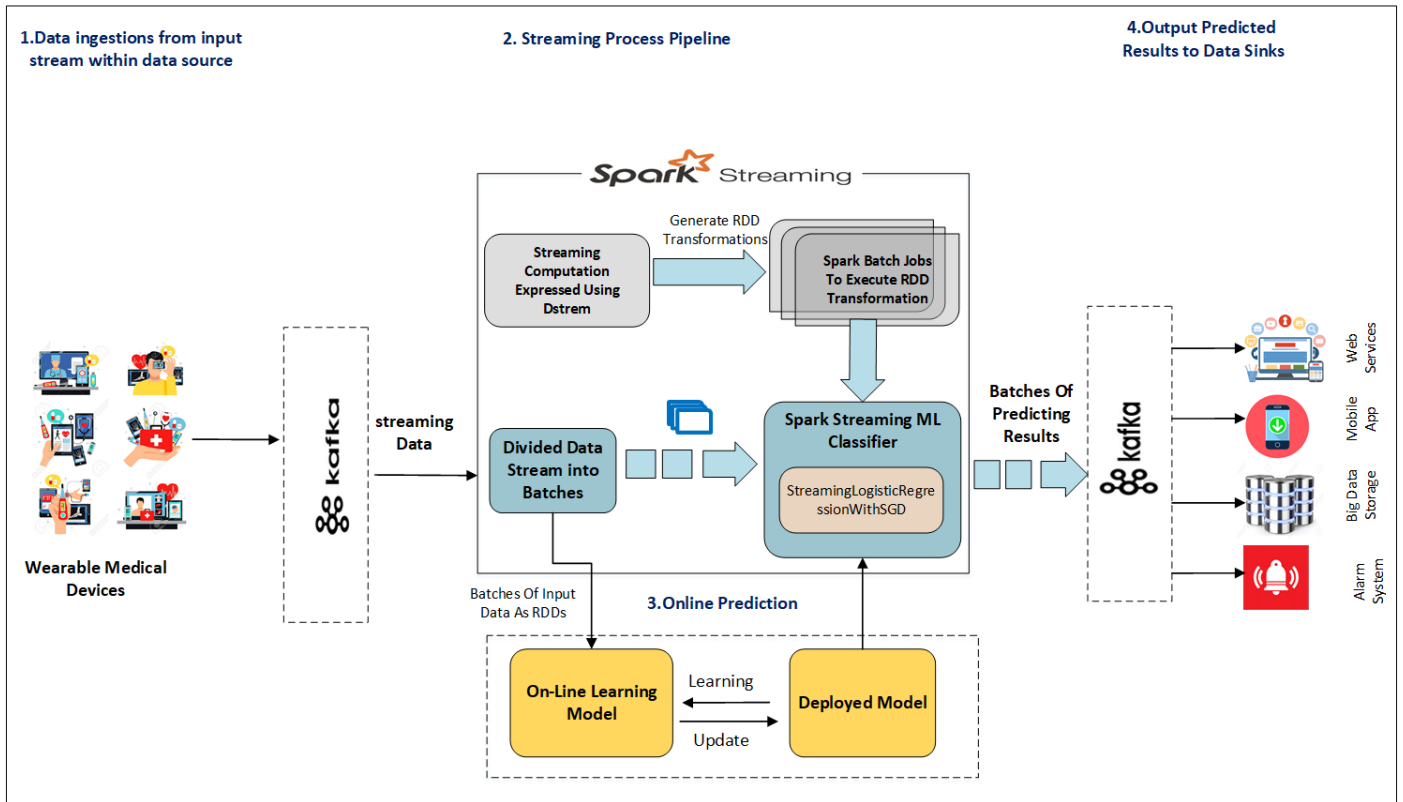


Fig. 1. The Online Health Status Prediction System.

classifier, which is the StreamingLinearRegressionWithSGD algorithm, is applied to train on streaming data and then predict patient status. The historical medical data have been read from the chosen datasets and then ingested into Kafka topic while the online data have been generated using the developed simulated generator and then ingested into Kafka topic as well. For hardware specification, the experiments have been performed using Spark cluster version 2.3.1, consisting of one master node and two nodes for workers. Table I illustrates the characteristics of the master node and the worker nodes. For the performance metrics, the performance evaluation of the classification is done through various performance measures such as accuracy, precision, recall, and F1-score. Moreover, the performance comparisons are done according to the number of batches and window size. Accordingly, the performance evaluation of online prediction techniques is done through eight experiments. The first, second and third experiments were conducted using historical medical datasets, and the rest five experiments were conducted using generated datasets and five window sizes such as 1, 2, 3, 4, and 5 seconds.

TABLE I. CLUSTER NODES CHARACTERISTICS.

| Parameter        | Master         | Worker         |
|------------------|----------------|----------------|
| Processor        | Core i7        | Core i7        |
| Cores            | 4              | 4              |
| Memory           | 20 GB          | 20 GB          |
| Operating System | Ubuntu 18.04.2 | Ubuntu 18.04.2 |

### B. Performance Analysis of the Historical Medical Datasets

To evaluate the efficiency of the proposed system, various experiments were conducted on different datasets that contain medical records describing the patient’s health information in terms of a set of attributes and the corresponding patient health status. For this research work, three medical datasets are used; Pima Indians diabetes dataset [24], Cleveland heart disease dataset [24], and breast cancer Coimbra disease dataset. Table II presents the description of the used datasets in terms of the number of samples, number of attributes, names of attributes, and labels. The reason behind using a real-time streaming method for historical data analysis is to assess the ability of the online prediction for the proposed system using benchmark datasets. Basically, Apache Spark reads data from a file and converts data to an RDD, then the continuous DStreams of data come from Spark Streaming. According to this work, each dataset is read from a CSV file and then sent to Spark streaming. Each RDD in a Dstream splits ingested data according to the window size and the size of the dataset. For the evaluation taken in this work, the datasets have been split into an 80% training set and a 20% testing set.

1) *Results of Pima Indians diabetes dataset:* In the first experiment, the Pima Indians diabetes dataset is performed. The number of samples of the Pima Indians diabetes dataset is 768, which has been read and sent to Spark Streaming. The window size has been configured into 2 seconds. According to the dataset size ( i.e., the number of samples) and the configured window size, Dstream splits ingested data into two batches; the first batch and the second batch denoted by 1st



TABLE II. THE MEDICAL DATASETS' DESCRIPTIONS.

| Dataset Name                  | Number of Samples | Numbers of Attributes | Names of Attributes  | Labels                               |
|-------------------------------|-------------------|-----------------------|--|--------------------------------------|
| Pima Indians diabetes dataset | 768               | 8                     | Pregnancies, Glucose, Blood Pressure, Insulin, Skin fold thickness, body mass index, diabetes pedigree function, Age   | Classes:<br>0- absence<br>1- present |
| Cleveland heart disease       | 303               | 13                    | Age, sex, chest pain type, resting blood pressure, serum cholestorol, fasting blood sugar, resting electrocardiographic results, maximum heart rate achieved, exercise induced angina, ST depression induced by exercise relative to rest, number of major vessels (0-3) colored by flourosopy, thal | Classes:<br>0- absence<br>1- present |
| Breast cancer coimbra dataset | 116               | 10                    | age, BMI, glucose, insulin, HOMA, leptin, adiponectin, resist in and MCP   | Classes:<br>1-healthy<br>2- patient  |

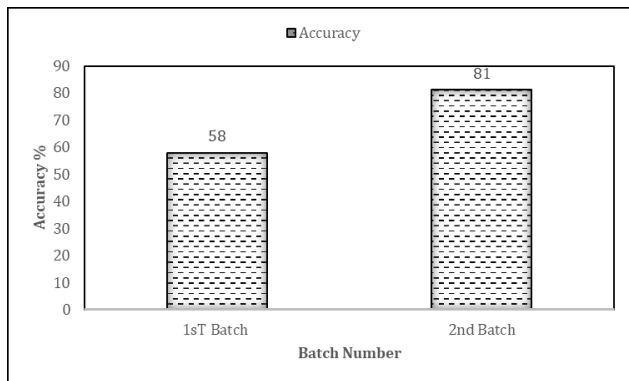


Fig. 2. The Accuracy of Online Prediction using the Pima Indians Diabetes Dataset.

Batch and 2nd Batch respectively. Table III shows the results of the Pima Indians dataset. It is noted that the second batch has obtained higher performances with respect to the first batch. We attribute this behavior to the more data ingested to the model, the performance metrics become higher. In particular, when the number of samples increases, the model learns and updates itself by time. For example, the obtained accuracy in the first batch is 58%, which increases to 81% for the second batch (see Table III and Fig. 2). As can be seen, other performance metrics, including precision, have improved with time as well in the second batch ( precision of 83%, recall of 81%, and F1-score of 80%). The improvement ratios for the second batch for the performance metrics with respect to the first batch are 28%, 22%, 28%, and 31% for accuracy, precision, recall, and F1-Score, respectively.

TABLE III. THE PERFORMANCE RESULTS OF ONLINE PREDICTION USING THE PIMA INDIANS DIABETES DATASET.

| Batch No  | Accuracy | Precision | Recall | F1-Score |
|-----------|----------|-----------|--------|----------|
| 1st Batch | 58       | 65        | 58     | 55       |
| 2nd Batch | 81       | 83        | 81     | 80       |

2) Results of Cleveland heart disease dataset: In the second experiment, the Cleveland heart disease dataset is performed. The number of samples of the Cleveland heart disease dataset is 303, which has been read and sent to Spark Streaming. Similar to the first experiment, we have configured the window size to 2 seconds. Consequently, Spark Dstream splits ingested heart disease dataset into two batches; the first batch and the second batch denoted by 1st Batch and 2nd Batch, respectively. Table IV shows the Cleveland heart disease dataset results, as can be seen, that the performances of the second batch are higher with those obtained by the first batch. For instance, the obtained accuracy in the first batch is 58%, which increases to 81% for the second batch (see Table IV and Fig. 3). Also, other performance metrics have improved in the second batch ( precision of 82%, recall of 79%, and F1-score of 78%). The reason behind this behavior is that the online prediction model learns and updates its performances by time. Based on this experiment using the Cleveland heart disease dataset, the improvement ratios for the second batch for the performance metrics with respect to the first batch are 18%, 13%, 18%, and 23% for accuracy, precision, recall, and F1-Score, respectively.

TABLE IV. THE PERFORMANCE RESULTS OF ONLINE PREDICTION USING THE CLEVELAND HEART DISEASE DATASET.

| Batch No  | Accuracy | Precision | Recall | F1-Score |
|-----------|----------|-----------|--------|----------|
| 1st Batch | 65       | 71        | 65     | 60       |
| 2nd Batch | 79       | 82        | 79     | 78       |

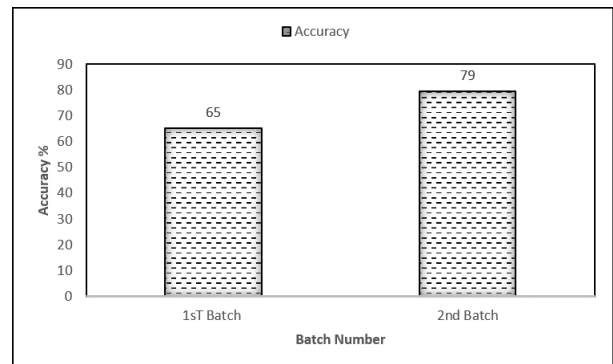


Fig. 3. The Accuracy of Online Prediction using the Cleveland Heart Disease Dataset.

3) Results of breast cancer Coimbra disease dataset: In the third experiment, the breast cancer Coimbra disease dataset is performed. The number of samples of the breast cancer dataset is 116, which has been read and sent to Spark Streaming. Similar to the previously conducted experiments, we have configured the window size to 2 seconds. Consequently, Spark Dstream splits ingested heart disease dataset into two batches; the first batch and the second batch denoted by 1st Batch and 2nd Batch, respectively. Table V shows the results of the breast cancer Coimbra disease dataset. It is noticed that the second batch outperforms the first batch. In particular, the performances of the second batch are higher with those obtained by the first batch. For instance, the obtained accuracy in the first batch is 63%, which increases to 67% for the second batch (see Table V and Fig. 4). Also, other performance metrics have improved in the second batch ( precision of 76%,

recall of 74%, and F1-score of 73%). The reason behind this behavior is that the online prediction model learns and updates its performances by time. Based on this experiment using the Cleveland heart disease dataset, the improvement ratios for the second batch for the performance metrics with respect to the first batch are 17%, 9%, 15%, and 19% for accuracy, precision, recall, and F1-Score respectively.

TABLE V. THE PERFORMANCE RESULTS OF ONLINE PREDICTION USING THE BREAST CANCER COIMBRA DISEASE DATASET.

| Batch No  | Accuracy | Precision | Recall | F1-Score |
|-----------|----------|-----------|--------|----------|
| 1st Batch | 63       | 70        | 63     | 60       |
| 2nd Batch | 76       | 76        | 74     | 73       |

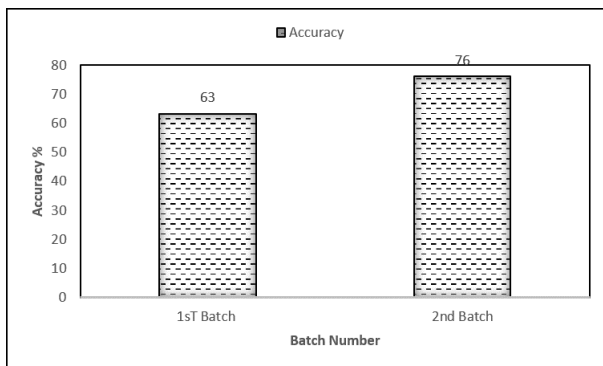


Fig. 4. The Accuracy of Online Prediction using the Breast Cancer Coimbra Disease Dataset.

4) Discussion: In the performance analysis of the historical medical datasets, three datasets have been evaluated using the proposed systems. Fig. 5 shows the accuracies of the second batch for the three medical datasets. As can be seen that the diabetes disease dataset has achieved the highest accuracy at 81% with respect to the heart and breast cancer datasets. The heart disease dataset has recorded the second rank on the average of the accuracy while breast cancer dataset has recorded the third rank (accuracy at 79% and 76% for heart and breast cancer dataset, respectively).

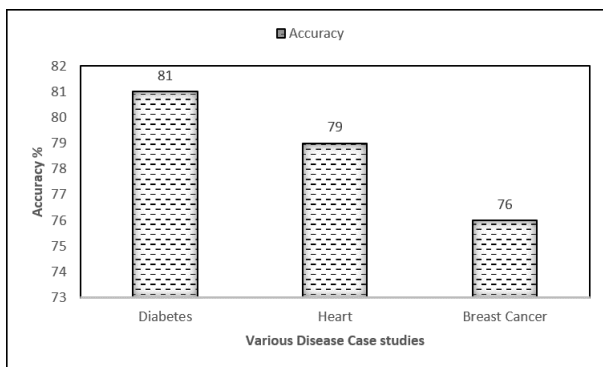


Fig. 5. The Accuracy of Online Prediction for the Second Batch using the Three Medical Datasets.

From the results obtained in our experiments, Fig. 6 depicts deeply the empirical results showing the improvement ratio of the accuracies for the second batch with respect to the first

batch for the three medical datasets. It can be noticed that the accuracy improvement ratio, which has been obtained using the diabetes disease dataset, is the highest one with respect to the other two datasets. However, the achieved accuracy of the first batch using the diabetes disease dataset was 58%, which is lower than the accuracies, which have been achieved by the other two datasets in the first batch ( accuracy at 65% and 63% for heart and breast cancer dataset, respectively). We attribute this behavior to the online prediction, which is responsible for the batching of input data. The online prediction can online learn and update the model according to the new data arrival. For the diabetes disease dataset, the arrival of real-time training data is larger compared to datasets because of the large number of samples, 768. Based on these results, it can be tentatively concluded that using a larger number of samples for online prediction will improve the accuracy of the proposed system over time. Particularly, the larger number of samples can lead to updating of the model in real-time and further improving the accuracy of real-time prediction. For this reason and due to lack of large datasets, this motivates us to develop a data generator that simulates the medical sensors data to generate a large number of samples, which would improve the online prediction performance.

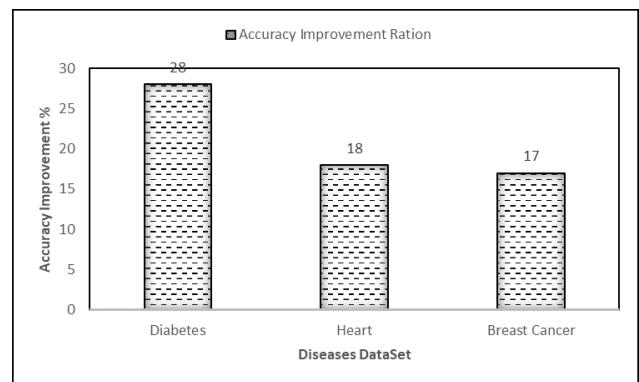


Fig. 6. The Improvement Ratios of the Online Prediction Accuracies for the Second Batch with respect to the First Batch using the Three Medical Datasets.

### C. Performance Analysis for Real-time Streaming Dataset

After evaluating the proposed system using the historical medical datasets, we have noticed that the large dataset achieves higher performance for online prediction. Therefore, to assess the efficiency of the online prediction system, we have developed a simulated data generator to generate a large dataset that can be trained for achieving higher performances using multiple DStreams. In particular, the simulated data generator has been developed to generate streaming diabetes samples as JSON format (see Fig. 7). The multiple data streams samples are sent via Apache Kafka and then processed by Spark Streaming. The rule of thumb of generated samples of diabetes dataset is introduced in [25]. The diabetes disease depends on three factors, which represented as attributes; A1c, Fasting Plasma Glucose (FPG), and Oral Glucose Test (OGT) (see Table VI). A1c tests the blood trail of a person for recent months. FPG tests a fasting plasma glucose level to recognize diabetes. OGT describes the oral glucose to analyze diabetes.

TABLE VI. DIABETES CONDITIONS

| A1c      | Fasting Plasma Glucose(0,199) | Oral Glucose Test(0,846) | Label |
|----------|-------------------------------|--------------------------|-------|
| A1c >5.7 | Glucose >100                  | Insulin >140             | 1     |
| A1c <5.7 | Glucose <= 99                 | Insulin <= 139           | 0     |

```

{
  "sensor_id": "1",
  "Label": "0",
  "observationTimestamp": "2020-07-26 05:47:47.246363",
  "Insulin": "51",
  "A1c": "4.8",
  "Glucose": "97"
},
{
  "sensor_id": "2",
  "Label": "1",
  "observationTimestamp": "2020-07-26 05:48:52.459688",
  "Insulin": "230",
  "A1c": "5.8",
  "Glucose": "118"
}
    
```

Fig. 7. JSON-like Generated Streaming Diabetes Data.

1) *Comparison using different window sizes:* To evaluate the efficiency of the proposed online prediction system, five experiments were conducted by fine-tuning the window sizes such as 1,2, 3, 4, and 5 seconds and using the generated diabetes dataset. The obtained results from diabetes prediction of the first experiment using 1 sec are shown in Table VII. It can be seen that the Spark streaming has split the generated data into 14 batches. Also, it can be noticed that the accuracy of the online prediction has increased linearly and improved by time (see Fig. 8). The highest accuracy obtained with the 1st batch is 59%, while the highest accuracy achieved by the 14th batch is 86%.

Similarly, as we have configured window size in 2 seconds for the second experiment, Spark streaming has split the generated data into 12 batches (see Table VIII). Fig. 9 depicted the accuracy of the online prediction for the 12 batches, where the accuracy of the online prediction has increased linearly and improved by time. The highest accuracy obtained with the 1st batch is 60%, while the highest accuracy achieved by the 12th batch is 88%. For the third experiment, the window size has been configured to 3 seconds, which makes Spark streaming to split the ingested generated data into 9 batches. Consequently, the performances of the online prediction using 3 sec as a window size have been shown in Table IX. Also, Fig. 10 presents the obtained accuracy, which increases linearly from the 1st batch to the 9th batch starting by 64% and ending by 90%.

Table X describes the performance of the online prediction using 4 seconds window size. The performances have been obtained among 6 batches, which are split by Spark streaming. The obtained performances have increased by time similar to the previous experiments. For instance, the accuracy has increased from 71% for the 1st batch and then 78% for the 2nd batch and so on (see Fig. 11). For the fifth experiment, the window size has been set to 5 seconds, which leads to 3 batches. Table XI presents the corresponding performances which have been obtained using a 5 second window size. Also, Fig. 12 depicts the accuracies of the online prediction, which are 85%, 93%, and 98% for the 1st, 2nd, and 3rd batch, respectively. It can be seen that the three obtained accuracies

have increased, and the improvement has grown faster by time.

TABLE VII. THE PERFORMANCE RESULTS OF ONLINE PREDICTION USING GENERATED DIABETES DATASET AND 1 SEC WINDOW SIZE.

| Batch No   | Accuracy | Precision | Recall | F1-score |
|------------|----------|-----------|--------|----------|
| 1st Batch  | 59       | 66        | 59     | 56       |
| 2nd Batch  | 60       | 68        | 60     | 57       |
| 3rd Batch  | 64       | 70        | 64     | 59       |
| 4th Batch  | 66       | 73        | 66     | 62       |
| 5th Batch  | 68       | 77        | 68     | 66       |
| 6th Batch  | 70       | 78        | 70     | 67       |
| 7th Batch  | 71       | 78        | 71     | 68       |
| 8th Batch  | 73       | 81        | 73     | 70       |
| 9th Batch  | 76       | 76        | 76     | 75       |
| 10th Batch | 77       | 77        | 77     | 77       |
| 11th Batch | 79       | 79        | 79     | 79       |
| 12th Batch | 81       | 82        | 82     | 81       |
| 13th Batch | 83       | 83        | 83     | 83       |
| 14th Batch | 86       | 85        | 85     | 85       |

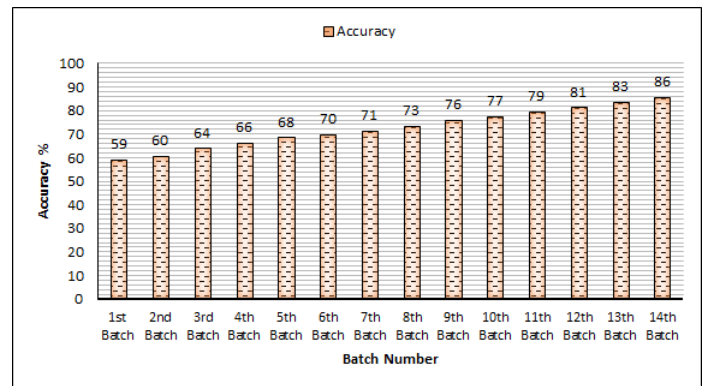


Fig. 8. The Accuracy of Online Prediction using the Generated Diabetes Dataset and 1 sec Window Size.

TABLE VIII. THE PERFORMANCE RESULTS OF ONLINE PREDICTION USING GENERATED DIABETES DATASET AND 2-SEC WINDOW SIZE.

| Batch No   | Accuracy | Precision | Recall | F1-score |
|------------|----------|-----------|--------|----------|
| 1st Batch  | 60       | 67        | 60     | 57       |
| 2nd Batch  | 62       | 69        | 61     | 59       |
| 3rd Batch  | 66       | 73        | 66     | 61       |
| 4th Batch  | 68       | 75        | 68     | 63       |
| 5th Batch  | 72       | 81        | 72     | 70       |
| 6th Batch  | 75       | 84        | 75     | 73       |
| 7th Batch  | 78       | 86        | 78     | 75       |
| 8th Batch  | 80       | 80        | 80     | 80       |
| 9th Batch  | 83       | 83        | 83     | 82       |
| 10th Batch | 85       | 85        | 85     | 84       |
| 11th Batch | 87       | 86        | 87     | 86       |
| 12th Batch | 88       | 88        | 88     | 88       |

2) *Discussion:* We analytically and experimentally summarize the performance gained from different window sizes. Fig. 13 depicts the window size tuning and their superiority over the time. In particular, if the window size is greater, the online prediction performances will be slightly improved. Consequently, a higher window size causes rapid prediction accuracy. For example, the 5 second window size has recorded 98% which is the highest accuracy among the other window sizes because the model learns and updates itself using three batches. Similarly, the 4 second, 3 second, 2 second and 1 have registered 95%,90%,88% and 86% using 6,9,12 and 14 batches, respectively. It can be seen that the larger window sizes allows the online prediction to process large data sizes,

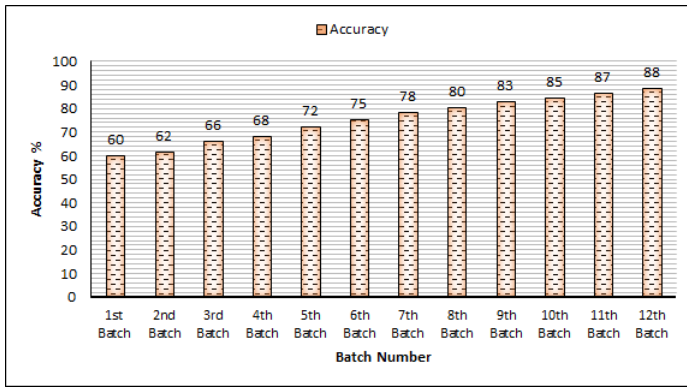


Fig. 9. The Accuracy of Online Prediction using the Generated Diabetes Dataset and 2-sec Window Size.

TABLE IX. THE PERFORMANCE RESULTS OF ONLINE PREDICTION USING GENERATED DIABETES DATASET AND 3 SEC WINDOW SIZE.

| Batch Number | Accuracy | Precision | Recall | F1-score |
|--------------|----------|-----------|--------|----------|
| 1st Batch    | 64       | 71        | 64     | 60       |
| 2nd Batch    | 70       | 76        | 70     | 67       |
| 3rd Batch    | 74       | 81        | 74     | 71       |
| 4th Batch    | 78       | 78        | 78     | 78       |
| 5th Batch    | 80       | 80        | 80     | 80       |
| 6th Batch    | 84       | 86        | 84     | 84       |
| 7th Batch    | 86       | 86        | 86     | 85       |
| 8th Batch    | 88       | 88        | 88     | 87       |
| 9th Batch    | 90       | 91        | 90     | 90       |

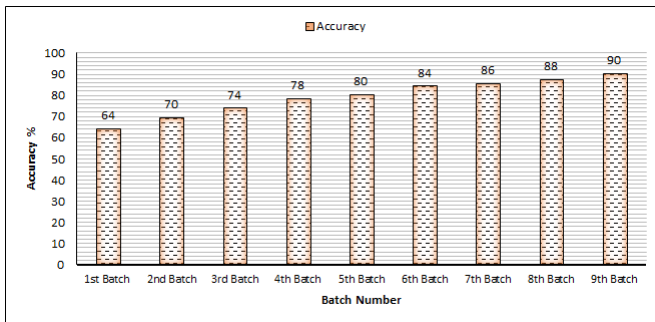


Fig. 10. The Accuracy of Online Prediction using the Generated Diabetes Dataset and 3-sec Window Size.

TABLE X. THE PERFORMANCE RESULTS OF ONLINE PREDICTION USING GENERATED DIABETES DATASET AND 4-SEC WINDOW SIZE.

| Batch Number | Accuracy | Precision | Recall | F1-score |
|--------------|----------|-----------|--------|----------|
| 1st Batch    | 71       | 78        | 71     | 69       |
| 2nd Batch    | 78       | 78        | 78     | 77       |
| 3rd Batch    | 83       | 83        | 83     | 82       |
| 4th Batch    | 88       | 89        | 88     | 88       |
| 5th Batch    | 92       | 93        | 92     | 92       |
| 6th Batch    | 95       | 95        | 95     | 94       |

TABLE XI. THE PERFORMANCE RESULTS OF ONLINE PREDICTION USING GENERATED DIABETES DATASET AND 5 SEC WINDOW SIZE.

| Batch Number | Accuracy | Precision | Recall | F1-score |
|--------------|----------|-----------|--------|----------|
| 1st Batch    | 85       | 86        | 85     | 84       |
| 2nd Batch    | 93       | 94        | 93     | 93       |
| 3rd Batch    | 98       | 98        | 98     | 97       |

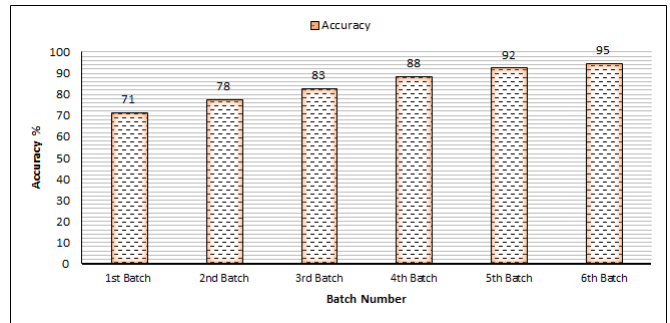


Fig. 11. The Accuracy of Online Prediction using the Generated Diabetes Dataset and 4-sec Window Size.

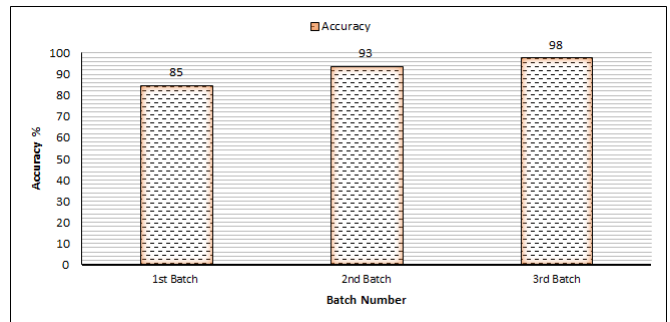


Fig. 12. The Accuracy of Online Prediction using the Generated Diabetes Dataset and 5 sec Window Size.

even less batch number which improves the prediction accuracy. We can conclude that the window size has a significant impact on the processing rate of Spark Streaming in terms of training a larger number of samples.

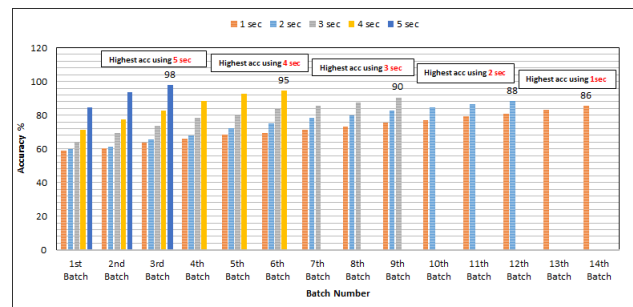


Fig. 13. Comparison of Accuracy of Online Prediction using Different Window Sizes.

## V. CONCLUSION

In this paper, we have presented an online prediction system to predict real-time health status. The proposed system has been developed using Spark Streaming, Apache Kafka, Apache Spark, and streaming machine learning algorithm named streaming linear regression with SGD. It has applied to two distributed data sources; historical medical data sources (diabetes, heart and breast cancer) and simulated wearable sensor generator which generates diabetes dataset. The diabetes dataset has achieved the highest accuracy at 81% with respect to the heart and the breast cancer datasets. The generated

diabetes dataset has achieved the highest accuracy at 98% using 5-second window size comparing to other window sizes: 1, 2, 3 and 4 seconds. The experimental results have shown that the larger window sizes allow the online prediction to process large amounts of data sizes, even fewer batch numbers, which improve prediction accuracy.

#### REFERENCES

- [1] R. Sahal, J. G. Breslin, and M. I. Ali, "Big data and stream processing platforms for industry 4.0 requirements mapping for a predictive maintenance use case," *Journal of Manufacturing Systems*, vol. 54, pp. 138 – 151, 2020. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0278612519300937>
- [2] G. Manogaran and D. Lopez, "Health data analytics using scalable logistic regression with stochastic gradient descent," *International Journal of Advanced Intelligence Paradigms*, vol. 10, no. 1-2, pp. 118–132, 2018.
- [3] A. Ed-daoudy and K. Maalmi, "A new internet of things architecture for real-time prediction of various diseases using machine learning on big data environment," *Journal of Big Data*, vol. 6, no. 1, p. 104, 2019.
- [4] K. Lee, A. Agrawal, and A. Choudhary, "Real-time disease surveillance using twitter data: demonstration on flu and cancer," in *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*, 2013, pp. 1474–1477.
- [5] U. Akhtar, A. M. Khattak, and S. Lee, "Challenges in managing real-time data in health information system (his)," in *International Conference on Smart Homes and Health Telematics*. Springer, 2016, pp. 305–313.
- [6] A. Spark, "Spark streaming," <https://spark.apache.org/docs/2.3.0/streaming-programming-guide.html/>, 2020.
- [7] A. Ed-daoudy and K. Maalmi, "Application of machine learning model on streaming health data event in real-time to predict health status using spark," in *2018 International Symposium on Advanced Electrical and Communication Technologies (ISAECT)*. IEEE, 2018, pp. 1–4.
- [8] J. Sun and C. K. Reddy, "Big data analytics for healthcare," in *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*, 2013, pp. 1525–1525.
- [9] J. Archenaa and E. M. Anita, "Interactive big data management in healthcare using spark," in *Proceedings of the 3rd International Symposium on Big Data and Cloud Computing Challenges (ISBCC-16)*. Springer, 2016, pp. 265–272.
- [10] L. R. Nair, S. D. Shetty, and S. D. Shetty, "Applying spark based machine learning model on streaming big data for health status prediction," *Computers & Electrical Engineering*, vol. 65, pp. 393–399, 2018.
- [11] S. Alotaibi, R. Mehmood, I. Katib, O. Rana, and A. Albeshri, "Sehaa: A big data analytics tool for healthcare symptoms and diseases detection using twitter, apache spark, and machine learning," *Applied Sciences*, vol. 10, no. 4, p. 1398, 2020.
- [12] A. Ed-Daoudy and K. Maalmi, "Real-time machine learning for early detection of heart disease using big data approach," in *2019 International Conference on Wireless Technologies, Embedded and Intelligent Systems (WITS)*. IEEE, 2019, pp. 1–5.
- [13] B. Akgün and Ş. G. Ögüdücü, "Streaming linear regression on spark mllib and moa," in *Proceedings of the 2015 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining 2015*, 2015, pp. 1244–1247.
- [14] G. B. Laleci, A. Dogac, M. Olduz, I. Tasyurt, M. Yuksel, and A. Okcan, "Saphire: a multi-agent system for remote healthcare monitoring through computerized clinical guidelines," in *Agent technology and e-health*. Springer, 2007, pp. 25–44.
- [15] W. N. Robinson, "Monitoring web service requirements," in *Proceedings. 11th IEEE International Requirements Engineering Conference, 2003*. IEEE, 2003, pp. 65–74.
- [16] M. Bransby and J. Jenkinson, *The management of alarm systems*. Citeseer, 1998.
- [17] R. L. Wears and S. J. Perry, "Human factors and ergonomics in the emergency department," *Annals of emergency medicine*, vol. 40, no. 2, pp. 206–212, 2002.
- [18] A. Franklin, S. Gantela, S. Shifarrow, T. R. Johnson, D. J. Robinson, B. R. King, A. M. Mehta, C. L. Maddow, N. R. Hoot, V. Nguyen *et al.*, "Dashboard visualizations: Supporting real-time throughput decision-making," *Journal of biomedical informatics*, vol. 71, pp. 211–221, 2017.
- [19] M. Fazio, A. Celesti, A. Puliafito, and M. Villari, "Big data storage in the cloud for smart environment monitoring," 2015.
- [20] R. Nachiappan, B. Javadi, R. N. Calheiros, and K. M. Matawie, "Cloud storage reliability for big data applications: A state of the art survey," *Journal of Network and Computer Applications*, vol. 97, pp. 35–47, 2017.
- [21] N. I. of Diabetes, Digestive, and K. Diseases, "Pima indians diabetes," <https://www.kaggle.com/uciml/pima-indians-diabetes-database>, 2020.
- [22] "Heart disease uci," <https://www.kaggle.com/ronitf/heart-disease-uci>, 2020.
- [23] "Breast cancer coimbra data set," <https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Coimbra>, 2020.
- [24] H. Kahramanli and N. Allahverdi, "Design of a hybrid system for the diabetes and heart diseases," *Expert systems with applications*, vol. 35, no. 1-2, pp. 82–89, 2008.
- [25] N. Sneha and T. Gangil, "Analysis of diabetes mellitus for early prediction using optimal features selection," *Journal of Big data*, vol. 6, no. 1, p. 13, 2019.

# Fundamental Capacity Analysis for Identically Independently Distributed Nakagami-q Fading Wireless Communication

Siam Bin Shawkat<sup>1</sup>, Md. Mazid-Ul-Haque<sup>2</sup>, Md. Sohidul Islam<sup>3</sup>, Borshan Sarker Sonok<sup>4</sup>  
Department of Computer Science  
American International University-Bangladesh

**Abstract**—With the advancement in technology, decent transfer rate of data for fast communication is an exigency. Different distributions on different wireless communication channels have been used previously to model them and to do performance analysis on the systems. In this work, capacity analysis of identically independently distributed Nakagami-q fading single-input multiple-output (SIMO) wireless communication is presented. The derivation of channel capacity with the analytical solution have been conducted using small limit argument approximation. Where the small limit argument approximation corresponds to the low signal-to-noise ratio (SNR) regime. SIMO channel capacity behavior with respect to number of receiver antennas and with respect to SNR have been explored in depth. The improvement of capacity is depicted rigorously. It has been found that using Nakagami-q distribution, capacity of the system increases as number of receiver antenna increases. It is also found that the capacity of this SIMO wireless system can be further improved through changing of certain parameters.

**Keywords**—Wireless communication; SIMO channel capacity; Nakagami-q fading; Hoyt distribution; low SNR regime

## I. INTRODUCTION

Wireless technologies have become an indivisible part of the daily life of a person in this digitized era. Wireless network with high capacity and high data rates is always wanted by everyone. A very important consideration is how fast data can be sent over a channel. Previously the main focus was to overcome the distortion of the received wireless signal. As technologies are increasing day by day, flawless high data transaction rate is now a big concern. Good channel capacity is one of the challenging issues in wireless communication. As a major need, channel capacity has to be improved as much as possible.

Shanon [1] developed the theorem which tells the maximum rate at which information can be transmitted over a communications channel of a specified bandwidth in the presence of noise. This theorem also specifies many important factors of communication channels. Since then many researchers have focused on analyzing and improving channel capacity on various distributions. These researches also involved systems like single-input single-output (SISO), single-input multiple-output (SIMO), multiple-input single-output (MISO), multiple-input multiple-output (MIMO). Authors in [2], discussed and showed that using multiple antennas at transmitter and receiver sides, higher capacity with efficiency can be attained. They showed that capacity increase has linear relation with the increase of antennas. In [3], authors presented comparisons between SISO

and MIMO system channels and also analyzed capacities in various distributions which involved Uniform distribution, Chi-square distribution, Gaussian distribution in the SISO system. Authors worked in Rayleigh fading single-input multiple-output (SIMO) channel system considering the presence of eavesdropper and found optimal power allocation at the side of transmitters in [4]. They have also described their proposed mathematical expression of secure outage probability. In [5], authors have examined both Broadcast channel and SISO channel in light communication system and used entropy power inequity and Lagrangian function to develop closed form lower and upper bound. They have also investigated the beamforming design problem of BC and VLC based on the utilization of obtained closed-form expression [5]. In [6], authors focused on increasing the data rate capacity of multiple-input multiple-output (MIMO) system channels through comparison of SISO, SIMO and MISO systems. Experimental simulated results of SISO and MIMO are analyzed, discussed and compared [6]. Based on given error probability and block length, authors in [7], investigated maximal achievable data rate over quasi static SIMO channels. They have also verified the fast convergence to outage capacity through zero dispersion when there is an increase in the block length. Authors in [8], analyzed discrete time Rayleigh fading capacity in SISO & MIMO systems, presented and discussed their results obtained from their findings. In [9], authors studied SISO and SIMO channels and measured frequency response of line-of-sight (LOS) and non-line-of-sight (NLOS) channels. Their measurement result indicates that high data rates can be supported by channel capacity and it depends on the distance of the transmitter and the receiver. They also found in their measurement in [9], that the relationship between SIMO system channel capacity and number of receiving antennas is not linear. Considering the presence of an external eavesdropper, authors in [10], studied confidential communication in gaussian MIMO channel with a number of receivers. They also focused on proving the secrecy capacity of MIMO channels which are degraded based on some important factors. Researchers focus on various distribution models for analyzing and developing SISO, SIMO, MISO, MIMO channels. Gaussian or normal distribution,

Rician, Nakagami-m, Nakagami-n, Rayleigh are very popular distribution models. Beside these models, Nakagami-q is one of them and has a wide area of research interest. A fading distribution, Nakagami-q or also known as Hoyt distribution, serves as a suitable decent model under incontrovertible conditions. Using this distribution model, authors in [11],

presented a simulation program where envelopes of received signals can be modelled. They found that this modeling can be done without presiding over NLOS [11]. In [12], authors showed that it is possible to construct Hoyt or Nakagami-q distribution from a conditional exponential distribution model. They proposed a method that is able to analyze performances of any wireless link under Nakagami-q fading in a very convenient approach. Based on the derivation of the squared Hoyt distribution formula, authors in [13], derived simple expressions of secondary link capacity on various scenarios of their research interests. For severe fading, they showed that the capacity of ST-SR link increases in the presence of severe fading. Authors used this Hoyt distribution to model and analyze limit of SISO wireless communication channel's data rate in [14].

In this work, capacity analysis of SIMO wireless system over the identically independently distributed Nakagami-q fading wireless channel is presented. This study follows the derivation of the channel capacity as described in [14] and modifies the system model and channel capacity equations to fit this SIMO wireless system of this study. First, the derivation of SIMO channel capacity equation is done for low SNR regime. Then the analysis of capacity is made with respect to the number of antennas at the receiver end while also depicting an improvement in the capacity. Capacity corresponding to the instantaneous SNR of this system and a comparison of the capacity of SIMO system and SISO system are also shown.

The rest of the paper is structured as follows. In Section 2, Nakagami-q/Hoyt distribution. In Section 3, SIMO system model for this research work. Calculation of capacity for low SNR regimes for this Nakagami-q fading channels in Section 4. Analysis and results is Section 5. At the end, Section 6 includes the conclusions and future work of this study.

## II. NAKAGAMI-Q/HOYT DISTRIBUTION

Throughout the long periods of remote correspondences, and relying upon framework working situations, an extraordinary number of channel models have been proposed to depict the measurements of the plentifulness and period of multipath signals. Nakagami-q is one of the popular proposed distribution models. This model was presented by Nakagami as an estimation for Nakagami-m distribution within the scope of fading that stretches out from one-sided Gaussian model to the Rayleigh model [15]. The distribution model is being utilized all the more habitually in execution investigation and different examinations identified with mobile radio interchanges. In [16], authors found that this model can be easily used in mobile communication channels. This fading is generally seen in satellite connection subject to amazing ionospheric flash and strongly shadowed conditions [16]. Nakagami-q or Hoyt model is normally used to delineate the transient sign subject to fading in variety of certain mobile communication channels [17]. It's probably distribution function (PDF) [18] is

$$p_{\gamma}(\gamma) = \frac{(1+q^2)\gamma}{2q\bar{\gamma}} e^{-\frac{(1+q^2)^2\gamma}{4q^2\bar{\gamma}}} I_0\left(\frac{(1-q^4)\gamma}{4q^2\bar{\gamma}}\right), \quad \gamma \geq 0 \quad (1)$$

where, q is the fading parameter of Nakagami-q distribution with the value range from 0 to 1. Instantaneous SNR is represented as  $\gamma$  and average SNR is represented as  $\bar{\gamma}$ .  $I_0(\cdot)$  is represented as the modified Bessel function of the first kind of zeroth-order [14]. When  $q=0$ , it represents the one sided Gaussian fading and Rayleigh fading is represented by  $q=1$ .

## III. SIMO SYSTEM MODEL

The single-input multiple-output (SIMO) frequently used to empower a collector framework that gets signals from various free sources to scrap the impacts of fading. It has been utilized for a long time with short wave listening/accepting stations to cope with the impacts of ionospheric impedance and fading. The single-input multiple-output (SIMO) is the model where the transmitter includes a single antenna and the receiver end has numerous receiving antennas.

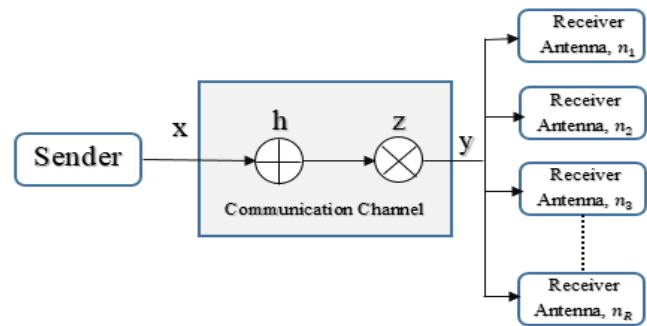


Fig. 1. Nakagami-q Fading SIMO Wireless Channel System Model.

Fig. 1 represents the SIMO system model for this research work considering the fact that the receiving power of each antennas are identical and the receiver antennas are mutually independent. Here, other considering factors include  $P_t$  as transmitted signal power, signal getting corrupted by additive white gaussian noise (AWGN) at the receiving end and the transmission is continued with Nakagami-q distribution. Based on [19] received signal vector of this system is given by,

$$r = hx + z \quad (2)$$

In Eq. (2), channel gain from the transmitter to the receiver end is defined by h which is a vector, x as transmitted signal and complex additive white gaussian noise is represented by the vector z [19]. For SIMO systems, [20] states, Shannon channel limit is the most aloft responded data between the signal that is being sent and being received one. From [20], the equation which defines capacity is given as following,

$$C = W \log_2(1 + \rho) \quad (3)$$

In Eq. (3), transmitted signal to noise ratio is represented as  $\rho = P_t/\sigma^2$  and transmission bandwidth is represented as W.

In [21] and [22] it has been said that, the channel is limited by power in the insight that,  $P_t = E\{|x|^2\}$ . Here, the  $E\{\cdot\}$  denotes the expectation operator. This expectation operator can be assessed by the probability distribution function (PDF) of the above mentioned vector  $h$ .

Here, for low SNR regime, low limit argument approximation is considered as it is described by the author that, using low limit argument approximation, the zeroth-order modified Bessel function of the first kind can be estimated as  $I_0 \approx 1$  in [17].

So, using  $I_0 \approx 1$  in Eq. (1), it becomes,

$$p_\gamma(\gamma) = \frac{(1+q^2)\gamma}{2q\bar{\gamma}} e^{-\frac{(1+q^2)^2\gamma}{4q^2\bar{\gamma}}}, \quad \gamma \geq 0 \quad (4)$$

So, for low SNR regime, Nakagami-q distribution under the condition of low limit argument approximation is represented by Eq. (4). By changing fading parameter to  $q=1$ , Eq. (4) turns to the Rayleigh fading as described in the Nakagami-q/Hoyt Distribution section.

For SIMO channel, all the receiver antennas will be receiving fading signals, so the Nakagami-q fading parameter becomes,  $q = n_R * q_0$ , where number of antennas at the receiver end are represented as  $n_R$  and when  $n_R = 1$  then the fading parameter,  $q$  becomes,  $1 * q_0 = q_0$ , so the system becomes single-input single-output.

So, for SIMO channel Eq. (4) becomes,

$$p_\gamma(\gamma) = \frac{(1+(n_R q_0)^2)\gamma}{2n_R q_0 \bar{\gamma}} e^{-\frac{(1+(n_R q_0)^2)^2\gamma}{4(n_R q_0)^2 \bar{\gamma}}}, \quad \gamma \geq 0 \quad (5)$$

For simplifying Eq. (5) let us consider

$$x = \frac{1+(n_R q_0)^2}{2n_R q_0 \bar{\gamma}}, y = \frac{(1+(n_R q_0)^2)^2}{4(n_R q_0)^2 \bar{\gamma}} \quad (6)$$

Using the simplifications of Eq. (6), Eq. (5) becomes,

$$p_\gamma(\gamma) = x\gamma e^{-y\gamma}, \quad \gamma \geq 0 \quad (7)$$

Here, Eq. (7) represents Nakagami-q distribution for for low SNR regime under low limit argument approximation [14].

#### IV. CALCULATION OF CAPACITY

The calculation of capacity in low SNR region for this SIMO wireless channel is done in this section. The capacity can be acquired a by using Eq. (3) [14],

$$C = \int_0^\infty p_\gamma(\gamma) \log_2(1 + \rho\gamma) d\gamma \quad (8)$$

Performing substitution of Eq. (7) into the Eq. (8), the channel capacity of SIMO channel becomes,

$$C = \int_0^\infty x\gamma e^{-y\gamma} * \log_2(1 + \rho\gamma) d\gamma \quad (9)$$

An advanced computing system, Mathematica [23] of Wolfram Research is used here to elucidate and to validate Eq. (9) in order to analyze the capacity of SIMO channel.

The channel capacity, using Eq. (9) results in,

$$C_{SIMO} = \frac{x}{\rho^2 \text{Log}[2]} \left( G_{2,3}^{3,1} \left( \frac{y}{\rho} \middle| \begin{matrix} -2, -1 \\ -2, 0 \end{matrix} \right) \right) \quad (10)$$

if  $\text{Re}[\rho] > 0$  and if  $\text{Re}[b] > 0$ .

Here,  $x = \frac{1+(n_R q_0)^2}{2n_R q_0 \bar{\gamma}}$ ,  $y = \frac{(1+(n_R q_0)^2)^2}{4(n_R q_0)^2 \bar{\gamma}}$  and  $G_{p,q}^{m,n}(x|_{b_q}^{a_p})$  is the Meijer Gamma function [14] and  $C_{SIMO}$  represents the capacity of SIMO channel.

In the next section, Eq. (10) is used for the analysis of the SIMO channel capacity. The analyses are conducted with respect to the number of receiver antennas, average SNR, instantaneous SNR of the channel and the comparison of the capacity of SIMO and SISO systems are depicted also.

#### V. ANALYSIS OF CAPACITY

The analyses of SIMO channel capacity are presented here. The section is divided into four subsections. Capacity behaviour with respect to the number of receiver antennas,  $n_R$  is presented in subsection one. For the improvement of the capacity of the system, capacity behaviour with respect to the number of receiver antennas,  $n_R$  while varying the average SNR,  $\bar{\gamma}$  of this system is presented in the next subsection. The following subsection includes capacity behaviour with respect to the instantaneous SNR of the system. And the last subsection contains a comparison between SIMO channel capacity of this research work and SISO channel capacity, which was depicted in [14].

For the analysis section the value for fading parameter  $q_0$  is considered 0.17 because it has been found in this analysis that when  $q_0 > 0.17$ , the conditions of Eq. (10) become false.

##### A. SIMO Channel Capacity with respect to the Number of Receiver Antennas

This subsection describes the channel capacity of the Nakagami-q fading SIMO channel. The channel capacity depicted by Eq. (10) is plotted with respect to the number of antennas at the receiver end. For this, value of  $n_R$  taken from 2 to 6. And the value of the average SNR is taken  $\bar{\gamma}=5$ . The verification of the correctness of Eq. (10) is done also which has been procured from Eq. (9).

It is observed in Fig. 2 that the Eq. (10) is the correct solution of Eq. (9) as both equations have been plotted considering the same values and both the equations yields the exact same plot. So, the correctness of the SIMO channel capacity equation is justified.

From Fig. 2, it can be seen that the capacity is increasing with the increasing of the antennas at the receiver end. When there are 2 antennas at the receiver end the value of capacity of



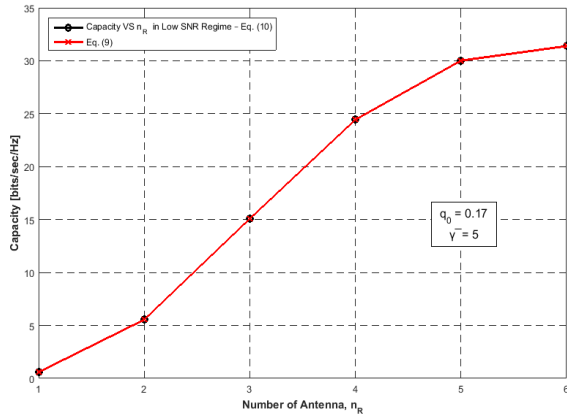


Fig. 2. SIMO Channel Capacity [bits/sec/Hz] with respect to the Number of Receiver Antennas,  $n_R$ .

is 5.54478 bits/sec/Hz, when  $n_R = 4$  the capacity is 24.4357 bits/sec/Hz and at last for  $n_R = 6$ , the channel capacity is at it's highest which is 31.3864 bits/sec/Hz.

So, it is clear that with the increasing of the antennas at the receiver end the capacity of the this system is increasing.

### B. Improvement of SIMO Channel Capacity

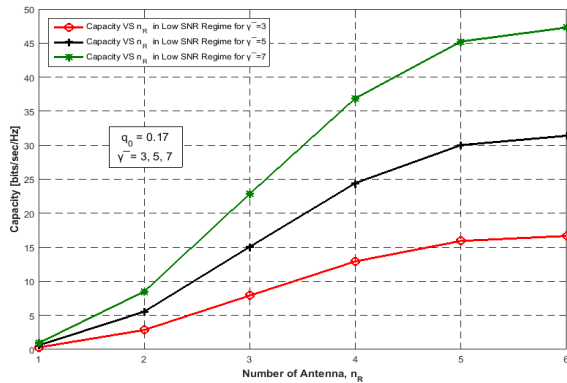


Fig. 3. Capacity [bits/sec/Hz] with respect to Number of Receiver Antennas,  $n_R$  when  $\bar{\gamma}$  varies

In this subsection again the capacity behaviour with respect to the number of receiver antennas are presented but varying the average SNR,  $\bar{\gamma}$  of the system. As in the previous subsection here also the number of receiver antennas are from 2 to 6. The value of average SNR of the system is taken,  $\bar{\gamma} = 3, 5, 7$  respectively for the analysis of this subsection.

It can be observed in Fig. 3 that for all the three cases the channel capacity is increasing when the number of receiver antenna increases. But, it can also be observed that the capacity increase is even higher when there is an increase in the average SNR.

When average SNR,  $\bar{\gamma} = 3$  the maximum capacity of the system is 16.6751 bits/sec/Hz, when average SNR,  $\bar{\gamma} = 5$  the

maximum capacity is 31.3864 bits/sec/Hz, at last, when  $\bar{\gamma} = 7$ , maximum capacity is 47.2828 bits/sec/Hz.

It is clearly seen that the capacity of this SIMO channel is on the rise as average SNR of the system increases. But, the capacity of the system is the highest when the average SNR,  $\bar{\gamma} = 7$ . So, the capacity can be improved when the average SNR of the system is considered higher.

### C. SIMO Channel Capacity with respect to the Instantaneous SNR

In previous subsection it has been found that there is an improvement in the capacity of of the system when the average SNR is higher. So, in this subsection  $\bar{\gamma} = 7$  is considered for analysis the SIMO channel capacity behaviour with respect to the instantaneous SNR of the system. And the number of receiver antennas are considered,  $n_R = 2, 4$  and 6.

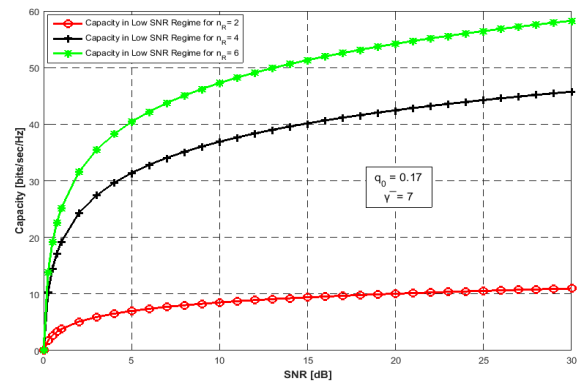


Fig. 4. SIMO Channel Capacity [bits/sec/Hz] with respect to the Instantaneous SNR [dB].

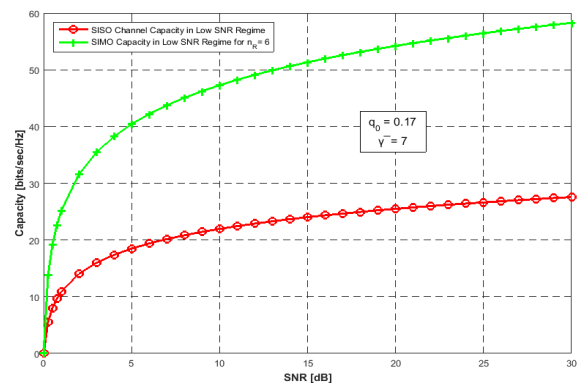


Fig. 5. Comparison of SISO Channel Capacity [14] and SIMO Channel Capacity

The increase of capacity with the increasing of instantaneous SNR for all the three numbers if receiver antennas can be observed from Fig. 4.

For receiver antenna,  $n_R = 2$ , the maximum capacity when SNR is at 30 dB is 10.9601 bits/sec/Hz. When, the number receiver antennas, for  $n_R = 4$ , the maximum capacity

is 45.7309 bits/sec/Hz and last, for  $n_R = 6$ , the maximum capacity at 30 dB is 58.2786 bits/sec/Hz.

The capacity of this system is higher when there are more receiver antennas, in this case it highest when there are  $n_R = 6$  number of receiver antennas.

### D. Comparison of SIMO and SISO Channel Capacity

In the previous subsection C it is found that when  $n_R = 6$  the capacity of this SIMO system is highest with respect to instantaneous SNR. So, in this subsection the SIMO channel capacity is considered for value of  $n_R = 6$ .

The SISO channel capacity is taken with respect to instantaneous SNR from [14] for conducting the comparison with the SIMO channel capacity of this work.

In Fig. 5, the red line represents the maximum SISO channel capacity with respect to the instantaneous SNR [14] in low SNR regime. The green line in Fig. 5 represents the maximum SIMO channel capacity of this work also in the low SNR regime. Maximum channel capacity of the SISO channel is 27.5641 bits/sec/Hz [14] and maximum capacity of SIMO channel observed in this work is 58.2786 bits/sec/Hz. There is about 111.43% increase in the capacity the SIMO system of this work.

So, it can be said that the capacity in SIMO system is much more higher than the capacity of SISO system.

## VI. CONCLUSION AND FUTURE WORK

In this work, the capacity of identically independently distributed Nakagami- $q$  fading SIMO channel in low SNR regime is studied. For this study, the mathematical expression for the capacity in low SNR regime of this SIMO system is derived using small limit argument approximation. This paper presented in-depth analyses of the channel capacity of the SIMO wireless system. It has been found in this work that the capacity is increasing with the increasing of number of antennas at the receiver end and the capacity is higher when the average SNR of the system is higher. With respect to instantaneous SNR, it is seen that there is exponential growth in the capacity. The maximum capacity of this SIMO wireless channel is obtained when there are 6 number of receiver antennas at the receiver end in this case. From the comparison of the channel capacity between Nakagami- $q$  fading SIMO and SISO wireless channels it has been found that the capacity of SIMO system outperforms the capacity of SISO system.

Though the derived equation can accurately measure the channel capacity, the measurement is only for low SNR regime. Moreover, as high SNR region is also an important part to consider when doing capacity performance analysis, the equations of this work can be re-used and modified to derive novel equation for high SNR region.

## REFERENCES

[1] C. Shannon, "Communication theory of secrecy systems," Bell System Technical Journal, vol. 29, pp. 656-715, 1949.  
[2] G. Foschini, "Layered space-time architecture for wireless communication in a fading environment when using multi-element antennas," Bell labs technical journal, vol. 1, no. 2, pp. 41-59, 1996.

[3] Arif Khan and Rein Vesilo, "A Tutorial on SISO and MIMO Channel Capacities," in Department of Electronics, Macquarie University NSW, Sydney Australia.  
[4] T. Ratnarajah, "Secrecy capacity and secure outage performance for Rayleigh fading SIMO channel," in 2011 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2011.  
[5] Shuai Ma, Ruixin Yang, Hang Li, Zhi-Long Dong, Huaxi Gu and Shiyin Li, "Achievable rate with closed-form for SISO channel and broadcast channel in visible light communication networks," Journal of Lightwave Technology, vol. 35, no. 14, pp. 2778-2787, 2017.  
[6] Nimay Ch Giri, Anwesha Sahoo, J. R. Swain, P. Kumar, A. Nayak and P. Debogowami, "Capacity and performance comparison of SISO and MIMO system for next generation network (NGN)," International Journal of Advanced Research in Computer Engineering and Technology (IJARCET), vol. 3, no. 9, pp. 30131-3035, 2014.  
[7] Wei Yang, Giuseppe Durisi, Tobias Koch and Yuri Polyanskiy, "Quasi-static SIMO fading channels at finite blocklength," In 2013 IEEE International Symposium on Information Theory, pp. 1531-1535, 2013.  
[8] Vignesh Sethuraman, Ligong Wang, Bruce Hajek and Amos Lapidoth, "Low-SNR capacity of noncoherent fading channels," IEEE Transactions on Information Theory, vol. 55, no. 4, pp. 1555-1574, 2009.  
[9] B. Nkakanou., G. Y. Delisle, N. Hakem and Y. Coulibaly, "UWB-SIMO channel capacity in an underground mine," In 2012 International Conference on Wireless Communications in Underground and Confined Areas, 2012.  
[10] Ersen Ekrem and Sennur Ulukus, "The secrecy capacity region of the Gaussian MIMO multi-receiver wiretap channel," IEEE Transactions on Information Theory, vol. 57, no. 4, pp. 2083-2114, 2011.  
[11] Sandro Adriano Fasolo and Renan Stihel Duque, "Fading channel simulator for Hoyt distribution," in Instituto Nacional de Telecomunicacoes, Brasil.  
[12] Romero-Jerez, Juan M. and F. Javier Lopez-Martinez, "A new framework for the performance analysis of wireless communications under hoyt (Nakagami- $q$ ) Fading," IEEE Transactions on Information Theory, vol. 63, no. 3, pp. 1693-1702, 2017.  
[13] Romero-Jerez, Juan M. and F. Javier Lopez-Martinez, "Fundamental capacity limits of spectrum-sharing in Hoyt (Nakagami- $q$ ) fading channels," in In 2016 IEEE 84th Vehicular Technology Conference (VTC-Fall), 2016.  
[14] Md. Mazid-Ul-Haque and Md. Sohiful Islam, "Data Rate Limit in Low and High SNR Regime for Nakagami- $q$  Fading Wireless Channel," International Journal of Advanced Computer Science and Applications(IJACSA), 11(7), 2020.  
[15] Nakagami and Minoru, "The  $m$ -distribution—A general formula of intensity distribution of rapid fading," in Statistical methods in radio wave propagation, Pergamon.  
[16] Youssef, Neji, Cheng-Xiang Wang, and Matthias Patzold, "A study on the second order statistics of Nakagami-Hoyt mobile fading channels," IEEE Transactions on Vehicular Technology, 54(4), pp. 1259-65, 2005.  
[17] K. S and Kölbig, "A definite integral with modified Bessel functions," Geneva.  
[18] Simon, Marvin K., and Mohamed-Slim Alouini, "Digital communication over fading channels," New York.  
[19] Amer M. Magableh and Mustafa M. Matalgah, "Capacity of SIMO systems over non-identically independent Nakagami- $m$  channels," in 2007 IEEE Sarnoff Symposium, Princeton, NJ, 2007.  
[20] J. G. Proakis, Digital Communication Systems, 4th edition, McGraw-Hill.  
[21] Md Sohiful Islam and Mohammad Rakibul Islam, "Ergodic Capacity of a SIMO System Over Nakagami- $q$  Fading Channel," DUET Journal, vol. 2, no. 1, 2014.  
[22] Md Sohiful Islam and Mohammad Rakibul Islam, "Positive secrecy mutual information over non-identically independently distributed Nakagami- $q$  fading wireless channel," in International Conference on Engineering, Research, Innovation and Education, 2013.  
[23] Wolfram Research, Inc., Mathematica, Version 12.0, Champaign, IL (2019).

# Unified Approach for White Blood Cell Segmentation, Feature Extraction, and Counting using Max-Tree Data Structure

Bilkis Jamal Ferdosi

Department of Computer Science and Engineering  
University of Asia Pacific  
Dhaka, Bangladesh

**Abstract**—Accurate identification and counting of White Blood Cells (WBCs) from microscopy blood cell images are vital for several blood-related disease diagnoses such as leukemia. The inevitability of automated cell image analysis in medical diagnosis results in a plethora of research for the last few decades. Microscopic blood cell image analysis involves three major steps: cell segmentation, classification, and counting. Several techniques have been employed separately to solve these three problems. In this paper, a simple unified model is proposed for White Blood Cell segmentation, feature extraction for classification, and counting with connected mathematical morphological operators implemented using the max-tree data structure. Max-tree creates a hierarchical representation of connected components of all possible gray levels present in an image in such a way that the root holds the connected components comprise of pixels with the lowest intensity value and the connected components comprise of pixels with the highest intensity value are in the leaves. Any associated attributes such as the size or shape of each connected component can be efficiently calculated on the fly and stored in this data structure. Utilizing this knowledge-rich data structure, we obtain a better segmentation of the cells that preserves the morphology of the cells and consequently obtain better accuracy in cell counting.

**Keywords**—Segmentation; feature extraction; White Blood Cell (WBC); mathematical morphology; max-tree

## I. INTRODUCTION

Microscopic blood cell image analysis is crucial for the diagnosis of several blood-related diseases. It may require complete blood count (CBC) where a complete count of red blood cells, white blood cells, and platelets is investigated. In some cases, differential blood count (DBC) may be required where five different types of white blood cells: eosinophils, basophils, monocytes, lymphocytes, and neutrophils need to be separated and counted. Blood image analysis is also crucial in the diagnosis of leukemia where lymphoblasts are needed to be separated from the healthy WBCs and counted. Manual analysis by the human experts is time-consuming, the accuracy of the result vastly depends on the expert's capability, and varying results may be obtained even if the procedure is repeated by the same expert. Thus, image-based analysis of blood cells gained much popularity in the past decades.

Image-based automated blood cell analysis poses three major challenges: segmentation, feature extraction for classification, and counting of cells from very complex blood smear images. To solve the challenging problem of cell segmenta-

tion, several approaches have been utilized in the literature. Clustering-based approaches such as expectation maximization (EM) [1], [2], K-means method [3], [4], the fuzzy C-means method [5], type-2 fuzzy logic [6], thresholding-based approach [7], edge detection based method [8], shape-based matching method [9], machine learning [10], or energy minimization [11], Gram-Schmidt orthogonalization [12], combining several image processing techniques such as thresholding, k-means clustering, and modified watershed algorithm [13], morphological operators [14], etc. to mention a few.

For classification different features such as morphological and textural features have been used [15]–[17]. Few others used genetic features extracted with the genetic algorithm [18]–[20].

Finally, different types of classified cells need to be counted. For counting, some methods that require prior cell segmentation and detection [21], few others approximated the number of cells from estimated density obtained from user annotation by compromising accuracy over speed [22], [23]. A method of cell counting based on morphological image analysis of blood cell images without requiring user annotation is reported in [24].

From segmentation to the counting of the cells widely varying techniques have been utilized in the literature. There is no unified approach that can facilitate in all three analysis steps of segmentation, feature extraction for classification, and counting. Inspired by the work in [24], this paper tries to use the full potential of Max-tree data structure which is an efficient structure for morphological connected operators. Morphological connected operators work on connected components of a gray level image known as flat zones and preserve only those flat zones that satisfy given criteria removing the rest of the flat zones [25], [26]. The criteria can be based on one or more attributes computed from the flat zones. Max-tree data structure enables the processing steps of these operators efficiently. Max-tree is a structured representation of an image where connected components with the highest intensity are in the leaves of the tree, the connected component with the lowest intensity is in the root, and the rest of the nodes hold the connected components for all threshold levels present in the image. Besides, the nodes of the tree are capable of storing a plethora of knowledge such as size and shape granulometry, texture, moment, or motion-oriented attributes. In this paper, the capability of this knowledge-rich data structure has been

utilized for cell segmentation, feature extraction, and counting. The cell counting part of the work is also reported in a conference paper [24]. In [24] only the structured representation of the image by the tree is utilized where the number of leaves in the tree is reported as the number of cells present in the image.

Piuri and Scotti proposed a system for leukocytes detection and classification based on the morphological operators in [27]. However, they used the structuring element based morphological operators for the identification of the cells. They only reported observational performance instead of any quantitative performance measure of their work. They also did not solve the cell counting problem. Besides, the structuring element based filtering is known to distort the original shape of the object. On the other hand, shape of the object is preserved in our method and thus yield better segmentation.

Moshavash et al. proposed a color-based cell segmentation technique where they converted the RGB images into CMYK color space and separate the background and red blood cells using the M-channel [28]. After the background and red blood cells are separated remaining are considered as the candidate for the WBCs.

On the other hand, this paper utilized connected morphological operators that filters without using any structuring element; rather it uses the structure of the input signal itself for filtering. Connected operators do not introduce any distortion or new structures to the resultant image. In this paper, a Max-tree data structure is used where connected components of the image are hierarchically stored thus every connected component is reachable for further processing. The main strengths of the proposed model are

- The method obtain a better segmentation of the cells that preserves the morphology of the different types of WBCs and their nuclei.
- The Method achieve better accuracy in counting cells.
- Feature extraction can be done on the fly.
- Segmentation, feature calculation, and counting are done in a unified model where Maxtree data structure plays the central role.
- Maxtree representation can be utilized in other applications such as cell tracking, cell visualization, etc.
- It is conceptually simple.
- Computationally efficient.

The rest of the paper is organized as follows: in Section 2, the concept of mathematical morphology and connected attribute filters are described. In Section 3, the basic idea, implementation detail of max-tree data structure, and attribute estimation methods are discussed. Section 4 contains a detailed discussion of the proposed method. A brief description of the data set can be found in Section 5. In Section 6, experiments and results are reported. Section 7 concludes the paper with a few directions for future work.

## II. MATHEMATICAL MORPHOLOGY AND CONNECTED ATTRIBUTE FILTERS

Mathematical Morphology [29] is popularly used in digital image analysis consists of several operators based on topological and geometrical concepts such as size, shape, contrast, etc. In mathematical morphology, grayscale images are considered in a form  $f(x) : \mathbb{Z} \rightarrow \mathbb{R}$  where  $\mathbb{Z} \subseteq \mathbb{E}$  mapping Euclidean space or grid,  $\mathbb{E}$  into  $\mathbb{R}$ . Image,  $f(x)$  is interacted with a small set called structuring elements,  $s(x)$  utilizing the order relation on  $\mathbb{R}$ . Connected attribute filters are morphological operators that can eliminate or merge the connected components or flat zones of an image where the image signal is constant [30]. Its power of simplifying images without distorting the contours makes it popular for various applications.

The notion of connectivity in digital images can be defined as the local neighborhood of pixels. If a 2D image,  $I$ , is mapped into  $m \times n$  grid and the position of a pixel,  $p$  in the grid is defined as row and column pair  $(i, j)$ . The four pairs of pixels positioned in  $(i \pm 1, j)$  and  $(i, j \pm 1)$  are 4-neighbors of  $p$ . A common choice of the local neighborhood of a pixel in the case of 2D images is 4 or 8 adjacency; and in the case of 3D is 6, 18, or 26 adjacency.

Connected operators act as a filtering tool on gray-level images that eliminates some of the connected components leaving other components unchanged. If  $I$  is the original image and  $S_k$  is the structuring element of size  $k$ , the opening of  $I$  can be defined using equation 1 which is erosion ( $\epsilon$ ) followed by dilation ( $\delta$ ):

$$J_0 = \delta S_k(\epsilon S_k(I)) \quad (1)$$

$$J_k = \delta_c(J(k-1)) \cap I \quad (2)$$

where  $c$  = structuring element that defines the connectivity

The operation in equation 1 will mark the connected components that need to be preserved and iterating equation 2 until idempotence will provide the desired result. Different types of connected operators can be obtained by the composition of any family of openings and closings by reconstruction. Connected attribute operators such as attribute openings, closing, thickenings, and thinnings can be utilized to filter connected components based on their attributes such as size, shape, contrast, etc. Simplest size oriented connected operators can be obtained by area opening and closing which is idempotent, anti-extensive, and increasing [31]. A large number of connected operators based on size attributes such as the moment of inertia, diagonal length of smallest enclosing box, etc. can be obtained for image filtering. The binary area opening,  $\Gamma$  of a binary image,  $I$  at point  $x$  with threshold parameter,  $\lambda$  obtains the connected component with an area greater or equal to  $\lambda$  and to which  $x$  belongs:

$$\Gamma_\lambda(I) = x \in I | A(\Gamma_x(I)) \geq \lambda \quad (3)$$

Apart from filtering images using size-based connected operators, shape-based filtering can be implemented using attribute thinning and thickening which is antiextensive, idempotent, and scale-invariant [32]. Being scale-invariant shape-based operators are insensitive to the size of the structures. Several shape-based attributes such as elongation [33], complexity or simplicity, motion estimation, entropy, etc. [34] can

be derived using these connected operators. Binary attribute thinning can be defined in terms of binary connected openings. The trivial thinning  $T^C$  with criterion  $C$  of a connected set  $S$  is the set that satisfies  $C$ , or empty otherwise. Thus  $T^C$  of set  $x$  with criterion  $C$  can be obtained by

$$T^C(I) = \bigcup_{x \in I} T^C(\Gamma_x(X)) \quad (4)$$

To decompose an image,  $I$ , according to size or shape, the image needs to be filtered using  $\Gamma_\lambda(I)$  or  $T^C(I)$  and after filtering the resulting image,  $I_r$  will contain structures that meet the criteria  $\lambda$  or  $C$ . The difference image,  $I - I_r$ , should contain the structures that fail to meet the desired criteria. Derivation of these size and shape based operators for grayscale images is straightforward and can be obtained from their binary counterpart.

### III. IMPLEMENTATION OF CONNECTED OPERATORS

Implementing connected operators using a structuring element for more than one dimension is difficult [35]. There are several algorithms proposed for attribute opening and closing such as the Pixel-Queue algorithm [32], the Union-Find algorithm [36], Max-Tree algorithm [34]. In [37] Meijster and Wilkinson discussed the pros and cons of these algorithms. The method proposed in this paper considered the findings of [37] and uses the Max-Tree algorithm because it requires linear time in both the number of pixels and connectivity for processing and pruning the tree and also for creating output images. Additionally, it also can be used for thinning and thickening.

#### A. Max-Tree

Max-Tree is a proficient data structure introduced by Salembier et al. in [34] for connected attribute filtering of the images. Max-tree being a rooted tree provides a hierarchy of flat zones with ordering relationships for extracting and filtering of the connected components by the operators. To describe the Max-Tree description of a few related terms is required. A connected component or flat zone,  $F_{l_i}$  at gray level  $l$  of an image  $I$  constitutes a set of pixels  $p \in E | I(p) = l$ ; a regional maxima,  $R_{l_j}$  corresponds to a flat zone at level  $l$  of which gray values of neighboring pixels are smaller than  $l$ ; a peak component,  $P_{l_k}$  is a flat zone of the thresholded image,  $T_l(I)$  at level  $l$ . In these definitions,  $i, j, k$  indicate the index of several such components. In the max-tree representation of an image, each node  $N_{l_k}$  holds only those pixels of the peak component  $P_{l_k}$  which have a gray value of  $l$ . The node  $N_{l_k}$  also contains attributes such as area of  $P_{l_k}$ . The node  $N_{0_0}$  is the root node of the tree and contains the flat zones with the lowest gray value and the hierarchical structure of the tree ensures that all of the flat zones of the highest gray value can be found in the leaf-nodes. After the max-tree representation of an image is obtained filtering of the image based on any attribute can be done by pruning the tree in an appropriate branch where all nodes have attributes such as area smaller than the threshold.

Fig. 1 depicts a pictorial description of the max-tree representation of an image containing six cells. The image has been binarized first then a Gaussian kernel is applied to

obtain a smoothed image that will produce flat zones with various gray levels such as  $G_0, G_1, G_2, G_3, G_4$  where  $G_0 < G_1 < G_2 < G_3 < G_4$  (figure 1(c)). In Fig. 1(d) max-tree representation of the flat zones is shown where the root  $A0$  represents the background of the image which is the flat zone at gray level  $G_0$ .  $B1, C1$ , and  $D1$  are the three children of the root that correspond to the three flat-zones at gray level  $G_1$ . Among them,  $B1$  has three flat zones as its children and  $D1$  has two flat zones as its children. All the flat zones are stored in the tree preserving the parent-child hierarchy and the leaves  $B04, B14, B24, C04, D04$ , and  $D14$  holds the connected components of the highest gray level  $G_4$ .

1) *Max Tree Creation*: The process of Maxtree creation can be described using the original work of Salembier et al. [25] where recursive flooding is utilized to build the tree. A hierarchical FIFO (first in first out) queue is implemented for each gray level value for the appropriate scanning and processing order of the pixels. The nodes,  $N_{G_i}$  (represents all nodes at gray level,  $G$  for all available gray levels,  $G$ ) and the links between parent and child nodes are established by storing all local background pixels of gray level,  $G$ , to the parent node and the child nodes get the pixels at each connected component with a gray level higher than  $G$ .

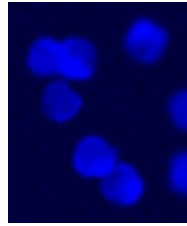
2) *Attribute Estimation and Filtering Approaches*: Different attributes such as size, shape, the moment of inertia, etc. of each node of max-tree can be calculated on the fly. During the max-tree creation required attributes of each node (for each connected component) is also calculated.

In max-tree implementation whenever a pixel is added to a node, an associated variable is increased to obtain the area of the component in terms of the number of pixels. Similarly, the ratio of the moment of inertia,  $M_I$  to the square of the area,  $A$  i.e.,  $M_I/A^2$  is easily and accurately calculated as a shape attribute.

Later, during the filtering phase based on the threshold value of a specific attribute the algorithm decides to remove or keep the node. Different decision criteria for filtering are described in [25] such as Min, Max, Viterbi, Direct, etc. Urbac et al. proposed another decision criterion: the subtractive decision rule in [33]. According to Min decision rule, a node,  $N_{l_i}$  is removed if the criterion value such as area, perimeter, or moment of inertia, etc., of that node, is less than the set threshold value or any of the ancestors of  $N_{l_i}$  is removed. Max decision rule removes a node,  $N_{l_i}$  along with its descendants if the criterion value is less than the set threshold value. Viterbi solves the filtering process as an optimization problem. The direct rule removes a node,  $N_{l_i}$  if its criterion value does not meet the threshold value and assigns the pixels of that node to a gray value of its highest ancestor which meets the threshold value while keeping the descendants unchanged. The subtractive decision rule is the same as the direct rule but it changes the gray value of the descendants by the same amount as the node,  $N_{l_i}$ .

### IV. PROPOSED METHOD

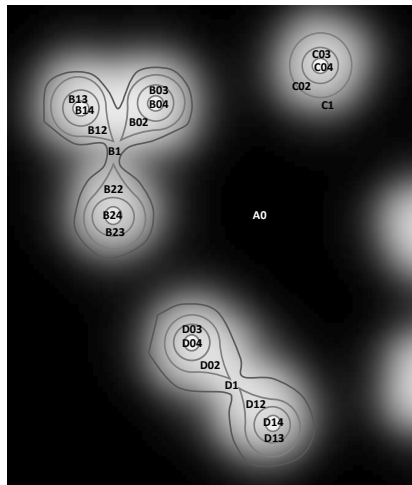
This paper proposed a method of blood cell image segmentation, feature extraction, and counting using connected morphological operators implemented using the max-tree data structure. The work is based on the principle of mathematical



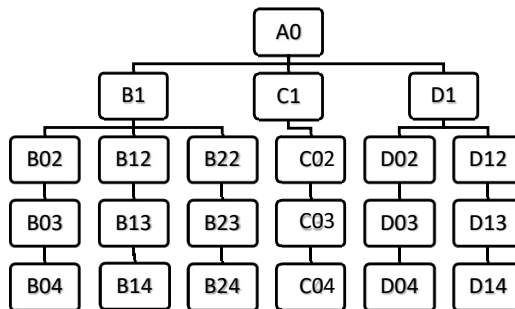
(a)



(b)



(c)



(d)

Fig. 1. (a) Cropped Region from a Microscopic Cell Image, (b) Thresholded image, (c) Different Gray Levels in Smoothed Image of (b), (d) Max-Tree Representation of the Image on (c)

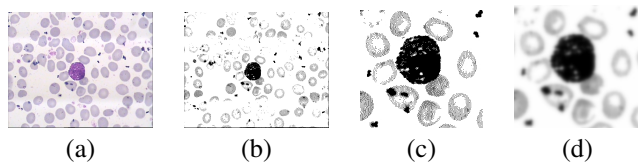


Fig. 2. (a) Stained Microscopic Cell Image, (b) Thresholded Image, (c) Cropped Region of (b), (d) Smoothed Image after Applying a Gaussian Kernel

morphology and connected attribute filters described in the previous section. Segmentation of the WBCs in the blood smear images has been done by pruning the max-tree that has been created as a representation of the image, based on area as a size attribute taking advantage of the distinct size of the RBCs and WBCs. Other attributes such as shape, moments, etc. that are computed on the fly during max-tree creation contributed to the feature vector to be used for the classification of the cells.

#### A. Morphology of the Blood Cells

The morphology of the blood cells plays a vital role in different steps of this work. There are three different major cells present in blood namely Red Blood Cell (RBC) or Erythrocytes, White Blood Cell (WBC) or Leukocytes, and Platelets (Thrombocytes). RBCs are generally smaller in size approximately  $7\mu\text{m} - 8\mu\text{m}$  and has a very thin cell membrane that allows easy oxygen diffusion. The morphology of an RBC is similar to a torus without a whole inside and does not contain any nucleus. On the other hand, WBCs are larger in size of approximately  $15\mu\text{m} - 20\mu\text{m}$  and contain a large nucleus that occupies most of the cell area. The platelets are the smallest in size approximately  $2\mu\text{m} - 3\mu\text{m}$  in diameter with a very irregular shape. In this work, the separation of RBCs and WBCs is discussed where the platelets are considered to the RBC group.

#### B. Preprocessing

Microscopic blood-stained image datasets are usually with high-quality images because of their exclusive and careful acquisition processes. Due to this assumption, there are no noise removal steps involve in this work. However, this step can be introduced if needed and several filtering approaches can be utilized.

The preprocessing step of our work is greatly influenced by the staining process of peripheral blood smear images. Staining blood cells with different colorants is a common practice in peripheral blood smear images so that various components especially the white and red blood cells can be examined microscopically. Typical components of the stains are oxidized methylene blue, azure B, and eosin Y colorants [38]. The methylene blue and azure B stains the nucleus of cells with different shades of blue to purple color and eosin Y colorants stains the cytoplasm of cells an orange to pink shades [39]. Thus, after staining torus like part of an RBC get a nonuniform shade of pink where intense pink in some areas and pale shade in some other areas. The central part of the cell gets brighter intensity. The nucleus of the WBC is intensely stained with

blue color and the cytoplasm is stained with different shades of blue.

To separate the WBCs from other blood cells the intense staining of its nucleus is utilized. Otsu's [40] method for thresholding is applied to the images and the nucleus of the WBC being intensely stained resulted into larger connected components and the RBC turned in to a combination of several small connected components almost the size of the dot because of the non-uniform staining [see Fig. 2(c)]. After binarization, a Gaussian kernel is applied to obtain a smoothed image [Fig. 2(d)]. This is done to minimize the noise introduced by the binarization process and to obtain a continuous grayscale image.

#### C. Cell and Nucleus Segmentation through Attribute Filtering

The block diagram in Fig. 3 summarizes the proposed method. The proposed method concentrates on segmenting WBCs, calculating features for the feature vector to be used in the classification of different WBCs, and counting the WBCs.

RBCs and WBCs can be distinctly identified by their size and presence of the nucleus. In the preprocessing step, WBCs due to the presence of a nucleus and its intense staining resulted in connected components with the largest area, the second-largest area is the area between the nucleus and the cell membrane i.e., the cytoplasm, and rest of the blood components resulted in much smaller connected components due to the binarization process. Therefore, filtering based on size (area) attribute results in the removal of those connected components with an area smaller than the threshold value. Filtering based on area attributes separate the candidates of WBCs but there remain several abnormal components that need to be removed. These abnormal components are not WBCs but stained similarly as a WBC. However, these components are irregularly shaped, unlike the WBCs which are mostly circular. Thus, shape-based attributes such as Eccentricity, Solidity, Compactness, etc. can be utilized to remove such abnormal components. However, in this work, the ratio of the moment of inertia to the square of the area (as described in Section III-A2), is used where the value of  $M_I/A^2$  increases as the shape deviates from the circular shape. This attribute is used for the filtering of the abnormal components.

In Fig. 4, cell and nucleus segmentation for a sample image from the ALL-IDB dataset is shown. The shape of the segmented region is the shape of the original cell or nucleus unlike structuring element based filtering where the shape of the segmented region is influenced by the shape of the structuring element. The proposed method faces difficulty if the cell boundary becomes hard to identify after the binarization process (see Fig. 5b). In that case, only the nucleus is segmented. However, in the Maxtree implementation, we also preserve the centroids of each nucleus which can be used to obtain the sub-image containing one WBC (Fig. 5d). Later using gray-level thresholding followed by hole filling and erosion resulted in cell segmentation (Fig. 5f).

#### D. Feature Extraction

Pattern spectra based on size, shape, or any other attribute or a combination can also be computed from the max-tree representation to obtain the feature vector for classification of

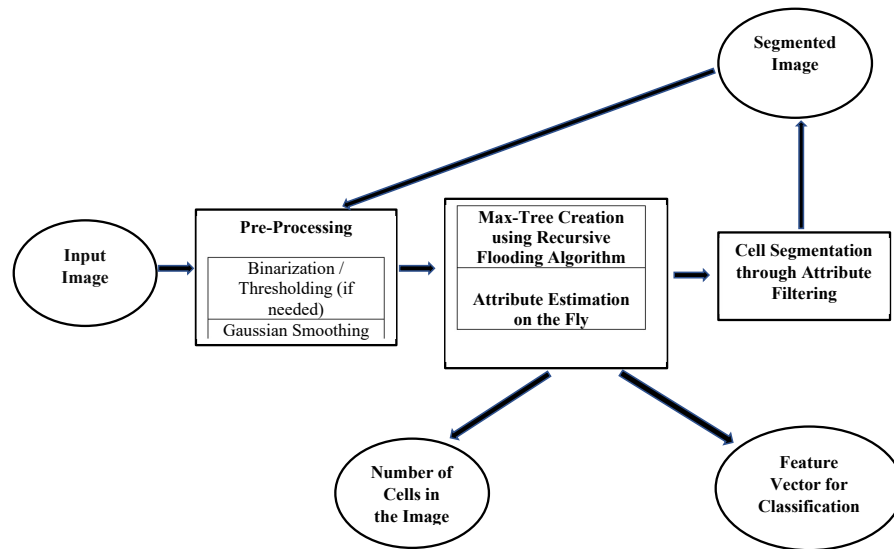


Fig. 3. Block Diagram of the Proposed Method

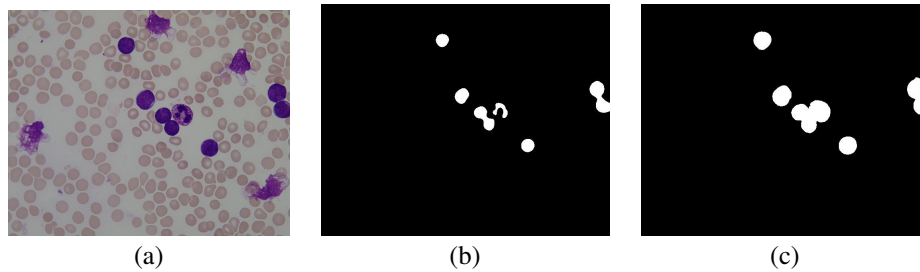


Fig. 4. (a) Image Sample from ALL-IDB 1 Dataset [41], (b) Cell Segmentation, and (c) Nuclues Segmentation by Area and Shape Attribute Filtering using Maxtree Representation.

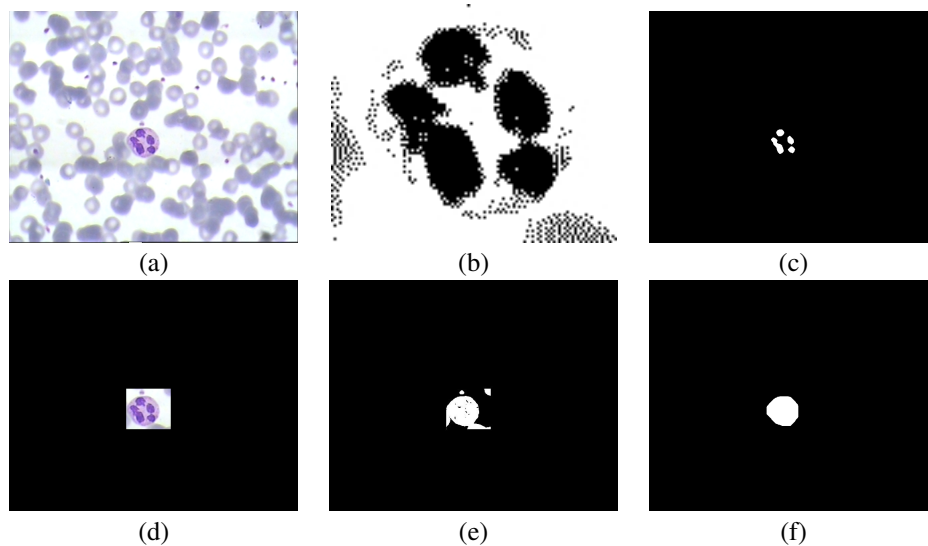


Fig. 5. Segmentation of Cells with Ill Defined Cell Boundary (a) Microscopic Image Containing Neutrophil WBC, (b) Binarized (Zoomed) Image of (a), (c) Nuclues Segmentation by Our Method, (d) Extracted Sub Image of (a) using Centroids of the Segmented Nucleus, (e) After Applying Gray Level Thresholding on (d), (f) Segmented WBC Cell.



the WBCs. Size and/or shape distribution, used for generating pattern spectra, comprises an ordered set of operators each of which removes features smaller than a particular size or shape from the image. Different sizes and shape-based features such as Area, Entropy, Moment of Inertia, Elongation, Mean x-position, Mean y-position, Eccentricity, Solidity, etc. can be calculated incrementally during the Maxtree creation. Besides, the number of lobes in the nucleus which is an important feature for the classification of WBCs can also be calculated from the Max-tree using the same technique used to count the number of cells [24]. Similarly, elongation of the connected components can be used as a shape measure. After obtaining the feature vector various state of the art classifiers can be explored for classification. In this paper, we did not explore the classifiers since a plethora of work achieved the state of the art classification performance with these features but the method of extracting the features was different than our approach. Nevertheless, we will also explore different classifiers with the features extracted using our approach in our future work.

### E. Cell Counting

In [24] we have reported the process of counting and annotation of cells using Max-tree representation of the cell images. In segmented cell images with a dark background, cells are represented with bright intensity. Therefore, in Max-tree representation cells being the extremal intensity can be found in the leaves of the tree (see Fig. 1) and thus the number of cells can be approximated by the number of leaves. The proposed method is conceptually easy, does not require any prior training or annotation on the contrary to the other state of the art approaches. Rather the proposed method provides the annotation of the cells for further use. Our Maxtree based cell counting approach is robust in case of partial overlapping of cells. Segmentation of such cells will result in a single connected component. However, filtering the binary image with a Gaussian kernel creates a region in the center of each cell which is brighter than its surrounding (see Fig. 6(c)). In Maxtree representation these regions will be in three leaves and thus counted as three different cells. However, if the cells are fully overlapped, it is not possible to identify them with our method and will be counted as one cell.

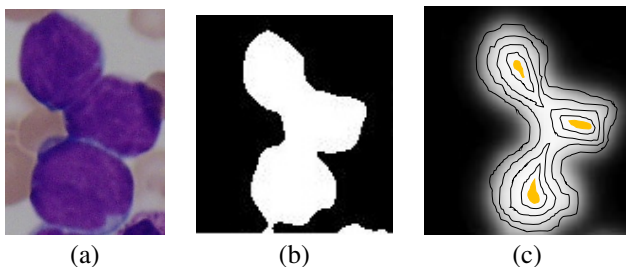


Fig. 6. (a) Partially Overlapped Cells, (b) Segmented as One Connected Component, (c) Identified as Three Different Cells in Proposed Method.

## V. DATASETS

In this work, we have tested the proposed method with two datasets: Leukocyte Images for Segmentation and Classification (LISC) database [42] and Acute Lymphoblastic Leukemia Image Database I (ALL-IDB I) [41].

### A. LISC Database

LISC Database <sup>1</sup> includes peripheral blood samples from healthy people. The Gismo-Right technique is used for smearing and staining the slides for obtaining microscopic images. The microscopic images are then digitized in BMP format with a size of (720 × 576) pixels. There are 250 images with ground truth provided with the freely distributed database. Ground truth for segmentation along with the classification of the WBCs into five classes of normal WBCs is done by the expert. There are 53 images with basophil, 39 images with eosinophil, 52 images with lymphocyte, 48 images with monocyte, and 50 images with neutrophil WBCs in this dataset.

### B. ALL-IDB I Database

The ALL-IDB I Database <sup>2</sup> includes blood samples from both Healthy Non-ALL subjects and probable ALL patients. There are 108 images in JPG format with 24-bit color depth and resolution of (2592 × 1944) pixels. It also includes the ground truth positions of the WBC cells in the images that are labeled by the experts.

## VI. EXPERIMENTS AND RESULTS

The performance of the WBC cell segmentation follows the performance of the nucleus segmentation. Segmentation of the nucleus uniquely identifies a WBC from other blood cells since other cells do not have a nucleus. Therefore, we report the performance of the nucleus segmentation and compare the performance with the performance of the method proposed by Moshavash et al. in [28]. In [28] Moshavash et al. used a similarity measure defined as in equation 5 for measuring the performance of nucleus segmentation:

$$Similarity = 100 \times \frac{(A_{algorithm} \cap A_{expert})}{\max(A_{algorithm}, A_{expert})} \quad (5)$$

LISC database contains five classes of WBCs, each of the classes has nuclei with a distinct morphology. For example, nuclei of Neutrophil are mostly multilobed, nuclei of Basophil and Eosinophil are bilobed, nuclei of Lymphocyte are eccentric, whereas the nucleus of Monocyte is almost kidney-shaped. The performance of WBC Classification is immensely influenced by the proper segmentation of the different types of the nucleus. Moshavash et al. reported that their method obtained 76% of the average similarity measure for the LISC database which is much better than other methods. However, there is no discussion of the segmentation performance of their method in different classes of WBCs. The similarity measure used in performance measurement is based on the assumption that the true segmentation is the manual segmentation done by an expert. However, in our experiments, we have observed that the segmentation made by the expert fails to achieve the complex morphology of nuclei properly in most of the cases. In Fig. 7 it is visible that the proposed method achieved better segmentation of nuclei where lobes of the nucleus are properly segmented. The number of lobes present in the nucleus is

<sup>1</sup>available at <http://users.cecs.anu.edu.au/hrezatofighi/Data/Leukocyte%20Data.htm>

<sup>2</sup>available at <https://homes.di.unimi.it/scotti/all/>

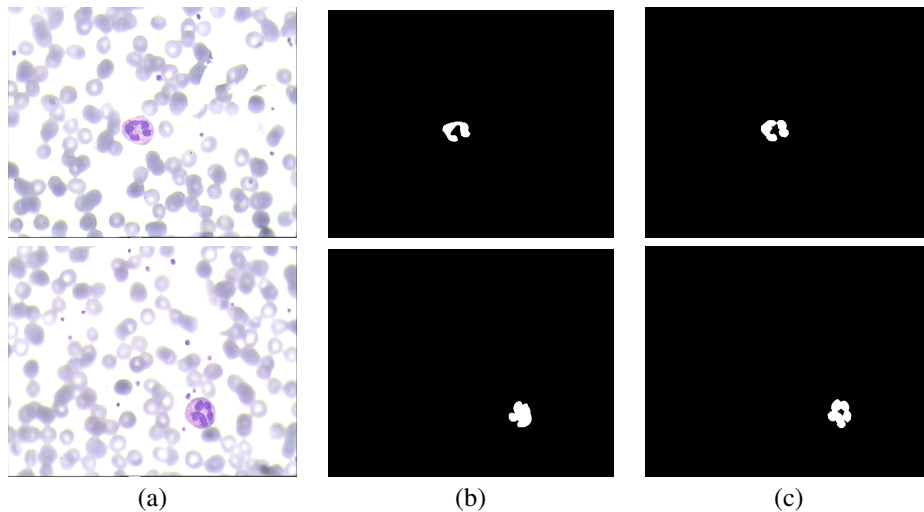


Fig. 7. (a) Microscopic Cell Image with Neutrophil WBC, (b) Segmentation of Nucleus by Expert, (c) Nucleus Segmentation by the Proposed Method where Multi-Lobed Morphology is better achieved compared to the Expert Segmentation

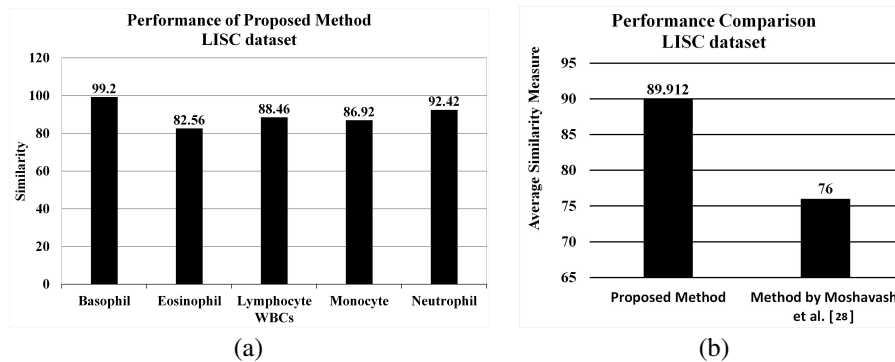


Fig. 8. (a) Performance of Proposed Method in Nucleus Segmentation of Different Types of WBCs of LISC Dataset, (b) Performance comparison of Proposed Method and Method Proposed by Moshavash et al. in [28]

one of the important features in classifying different types of WBCs. The proposed method will be able to obtain this information more accurately because of its better segmentation of the lobes. Therefore, similarity measure obtained using equation 5 may not be able to perform proper justice to the proposed method and the actual performance of the proposed method would excel.

Even though, the proposed method obtained significant improvement compared to the state-of-the-art method proposed by Moshavash et al. in [28] in segmenting nucleus (as seen in Fig. 8) where the proposed method achieved 89.912% average similarity measure. The segmentation performance of the proposed method varies among different types of the nucleus. This varied performance is mainly due to the quality of the staining of the nucleus. Cells with the properly stained nucleus are segmented more accurately since proper staining resulted in larger connected components.

In the ALL-IDB1 dataset, only the positions of the lymphoblasts are identified by the experts. However, there are lots of other healthy WBCs along with several abnormal components present in the images. Due to the absence of ground

truth position of all WBCs and ground truth segmentation of the cells, it is not possible to measure the segmentation performance using the equation 5. Therefore, the performance of this dataset is evaluated visually (Moshavash et al., also do not report any similarity measure for this data set). Few results are shown in Fig. 9.

The performance of the cell counting algorithm proposed in [24] by the author of this paper mostly depends on the performance of the segmentation. If the cells are segmented properly the performance of the proposed counting algorithm also improves. In Fig. 10 the performance of the counting algorithm on the LISC dataset can be observed.

## VII. CONCLUSION

The task of automated microscopic cell image analysis to identify the WBCs is very complex. Numerous researches have been done to solve this problem. However, to overcome the challenges of different stages such as segmentation, feature extraction, classification, and counting different techniques are utilized. In this paper, a system based on connected mathematical morphological operators implemented using Maxtree data

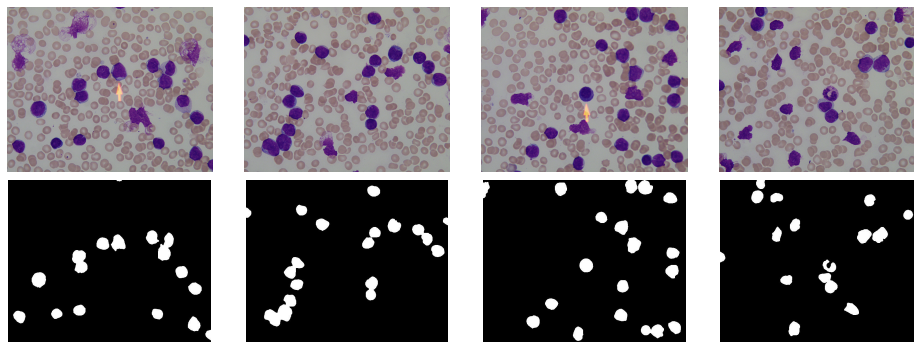


Fig. 9. ALL-IDB 1 Dataset: (Top Row) Original Image (Bottom Row) Segmented by the Proposed Method

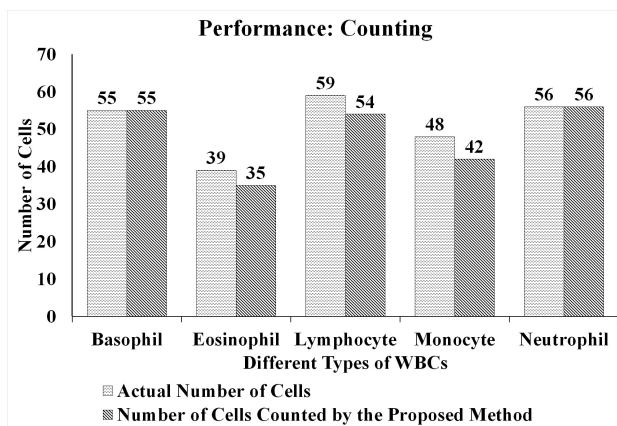


Fig. 10. Performance of the Cell Counting Algorithm on LISC Data Set

structure is proposed to solve the problem of segmentation, feature calculation, counting of the WBCs. In the proposed system Maxtree data structure plays the central role that facilitates the analysis of the cell images. It stores the connected components of every gray level present in the image along with their different attributes/features calculated on the fly during max-tree creation. These attributes are used in the segmentation stage and also can be used as a feature vector in the classification stage. Besides, the hierarchical structure of the tree enables the counting of the number of WBCs present in the image. The proposed system is conceptually easy, computationally efficient, and performs better than a state of the art method.

### VIII. FUTURE WORK

In the future, the performance of the classifiers will be explored using the features obtained from Maxtree. An interactive platform for ALL cell detection will be developed using the proposed technique.

### REFERENCES

- [1] C. Carson, S. Belongie, H. Greenspan, and J. Malik, "Blobworld: Image segmentation using expectation-maximization and its application to image querying," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 24, no. 8, p. 1026–1038, Aug. 2002. [Online]. Available: <https://doi.org/10.1109/TPAMI.2002.1023800>
- [2] S. Chen, L. Cao, Y. Wang, J. Liu, and X. Tang, "Image segmentation by map-ml estimations," *IEEE Transactions on Image Processing*, vol. 19, no. 9, pp. 2254–2264, Sep. 2010.
- [3] M. Mignotte, "A de-texturing and spatially constrained k-means approach for image segmentation," *Pattern Recogn. Lett.*, vol. 32, no. 2, p. 359–367, Jan. 2011. [Online]. Available: <https://doi.org/10.1016/j.patrec.2010.09.016>
- [4] B. J. Ferdosi, S. Nowshin, F. A. Sabera, and Habiba, "White blood cell detection and segmentation from fluorescent images with an improved algorithm using k-means clustering and morphological operators," *2018 4th International Conference on Electrical Engineering and Information & Communication Technology (iCEEICT)*, pp. 566–570, 2018.
- [5] R. Saha, M. Bajger, and G. Lee, "Spatial shape constrained fuzzy c-means (fcm) clustering for nucleus segmentation in pap smear images," in *2016 International Conference on Digital Image Computing: Techniques and Applications (DICTA)*, Nov 2016, pp. 1–8.
- [6] Z. Wang and Y. Yang, "A non-iterative clustering based soft segmentation approach for a class of fuzzy images," *Applied Soft Computing*, vol. 70, pp. 988 – 999, 2018. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S1568494617302788>
- [7] X. Chen, X. Zhou, and S. T. C. Wong, "Automated segmentation, classification, and tracking of cancer cell nuclei in time-lapse microscopy," *IEEE Transactions on Biomedical Engineering*, vol. 53, no. 4, pp. 762–766, April 2006.
- [8] C. Wählby, I.-M. Sintorn, F. Erlandsson, G. Borgefors, and E. Bengtsson, "Combining intensity, edge and shape information for 2d and 3d segmentation of cell nuclei in tissue sections," *Journal of Microscopy*, vol. 215, no. 1, pp. 67–76, 2004.
- [9] E. Türetken, X. Wang, C. J. Becker, C. Haubold, and P. Fua, "Network flow integer programming to track elliptical cells in time-lapse sequences," *IEEE Transactions on Medical Imaging*, vol. 36, no. 4, pp. 942–951, April 2017.
- [10] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *MICCAI 2015: 18th International Conference, Munich, Germany*. Springer International Publishing, 2015, pp. 234–241.
- [11] M. Maška, O. Daněk, S. Garasa, A. Rouzaut, A. Muñoz-Barrutia, and C. Ortiz-de-Solorzano, "Segmentation and shape tracking of whole fluorescent cells based on the chan-vease model," *IEEE Transactions on Medical Imaging*, vol. 32, no. 6, pp. 995–1006, June 2013.
- [12] S. H. Rezaatofghi, H. Soltanian-Zadeh, R. Sharifian, and R. A. Zoroofi, "A new approach to white blood cell nucleus segmentation based on gram-schmidt orthogonalization," in *2009 International Conference on Digital Image Processing*, March 2009, pp. 107–111.
- [13] N. Ghane, A. Vard, A. Talebi, and P.Nematollahy, "Segmentation of white blood cells from microscopic images using a novel combination of k-means clustering and modified watershed algorithm," *Journal of Medical Signals and Sensors*, vol. 7, no. 2, pp. 92–101, 2017.
- [14] L. B. Dorini, R. Minetto, and N. J. Leite, "White blood cell segmentation using morphological operators and scale-space analysis," in *Proceedings of the XX Brazilian Symposium on Computer Graphics and Image Processing*, ser. SIBGRAPI '07. USA: IEEE Computer Society, 2007, p. 294–304. [Online]. Available: <https://doi.org/10.1109/SIBGRAPI.2007.43>

- [15] N. Theera-Umpon and S. Dhompongasa, "Morphological granulometric features of nucleus in automatic bone marrow white blood cell classification," *IEEE Transactions on Information Technology in Biomedicine*, vol. 11, no. 3, pp. 353–359, May 2007.
- [16] J. . Thiran and B. Macq, "Morphological feature extraction for the classification of digital images of cancerous tissues," *IEEE Transactions on Biomedical Engineering*, vol. 43, no. 10, pp. 1011–1020, Oct 1996.
- [17] D. M. U. Sabino, L. da Fontoura Costa, E. Gil Rizzatti, and M. Antonio Zago, "A texture approach to leukocyte recognition," *Real-Time Imaging*, vol. 10, no. 4, p. 205–216, Aug. 2004.
- [18] T.-C. Lin, R.-S. Liu, Y.-T. Chao, and S.-Y. Chen, "Classifying subtypes of acute lymphoblastic leukemia using silhouette statistics and genetic algorithms," *Gene*, vol. 518, no. 1, pp. 159 – 163, 2013, proceedings of the 23rd International Conference on Genome Informatics (GIW 2012).
- [19] N. Zong, M. Adjouadi, and M. Ayala, "Artificial neural networks approaches for multidimensional classification of acute lymphoblastic leukemia gene expression samples," in *Proceedings of the 9th WSEAS International Conference on Computers*, ser. ICCOMP'05. Stevens Point, Wisconsin, USA: World Scientific and Engineering Academy and Society (WSEAS), 2005.
- [20] M. E. Ross, X. Zhou, G. Song, S. A. Shurtleff, K. Girtman, W. K. Williams, H.-C. Liu, R. Mahfouz, S. C. Raimondi, N. Lenny, A. Patel, and J. R. Downing, "Classification of pediatric acute lymphoblastic leukemia by gene expression profiling," *Blood*, vol. 102, no. 8, pp. 2951–2959, 10 2003.
- [21] J. Cheng, M. Veronika, and J. C. Rajapakse, "Identifying cells in histopathological images," in *Recognizing Patterns in Signals, Speech, Images and Videos*, D. Ünay, Z. Çataltepe, and S. Aksoy, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2010, pp. 244–252.
- [22] L. Fiaschi, U. Koethe, R. Nair, and F. A. Hamprecht, "Learning to count with regression forest and structured labels," in *Proceedings of the 21st International Conference on Pattern Recognition (ICPR2012)*, Nov 2012, pp. 2685–2688.
- [23] W. Xie, J. A. Noble, and A. Zisserman, "Microscopy cell counting and detection with fully convolutional regression networks," *Computer Methods in Biomechanics and Biomedical Engineering: Imaging & Visualization*, vol. 6, no. 3, pp. 283–292, 2018.
- [24] B. J. Ferdosi, "Microscopy cell counting and annotation using a max-tree representation of the blood cell images," in *Proceedings of the 3rd International Conference on Biomedical Signal and Image Processing*, ser. ICBIP '18. New York, NY, USA: Association for Computing Machinery, 2018, p. 61–65.
- [25] P. Salembier, A. Oliveras, and L. Garrido, "Antiextensive connected operators for image and sequence processing," *IEEE Transactions on Image Processing*, vol. 7, no. 4, pp. 555–570, April 1998.
- [26] R. Jones, "Connected filtering and segmentation using component trees," *Comput. Vis. Image Underst.*, vol. 75, no. 3, p. 215–228, Sep. 1999. [Online]. Available: <https://doi.org/10.1006/cviu.1999.0777>
- [27] V. Piuri and F. Scotti, "Morphological classification of blood leucocytes by microscope images," in *2004 IEEE International Conference on Computational Intelligence for Measurement Systems and Applications, 2004. CIMSA.*, July 2004, pp. 103–108.
- [28] H. Moshavash, Z. and Danyali and M. Helfroush, "An automatic and robust decision support system for accurate acute leukemia diagnosis from blood microscopic images," *J Digit Imaging*, vol. 31, p. 702 – 717, 2018.
- [29] J. Serra, *Image Analysis and Mathematical Morphology*. USA: Academic Press, Inc., 1983.
- [30] P. Salembier and J. Serra, "Flat zones filtering, connected operators, and filters by reconstruction," *IEEE Transactions on Image Processing*, vol. 4, no. 8, pp. 1153–1160, Aug 1995.
- [31] L. Vincent, "Morphological area openings and closings for grey-scale images," in *Shape in Picture*. Berlin, Heidelberg: Springer Berlin Heidelberg, 1994, pp. 197–208.
- [32] E. J. Breen and R. Jones, "Attribute openings, thinnings, and granulometries," *Comput. Vis. Image Underst.*, vol. 64, no. 3, p. 377–389, Nov. 1996.
- [33] E. R. Urbach and M. H. F. Wilkinson, "Shape-only granulometries and gray-scale shape filters," 2002.
- [34] P. Salembier, A. Oliveras, and L. Garrido, "Antiextensive connected operators for image and sequence processing," *IEEE Transactions on Image Processing*, vol. 7, no. 4, pp. 555–570, April 1998.
- [35] P. T. Jackway and M. Deriche, "Scale-space properties of the multiscale morphological dilation-erosion," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 18, no. 1, pp. 38–51, Jan 1996.
- [36] M. H. F. Wilkinson and J. B. T. M. Roerdink, *Fast Morphological Attribute Operations Using Tarjan's Union-Find Algorithm*. Boston, MA: Springer US, 2000, pp. 311–320.
- [37] A. Meijster and M. H. F. Wilkinson, "A comparison of algorithms for connected set openings and closings," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 24, no. 4, pp. 484–494, April 2002.
- [38] P. N. Marshall, "Romanowsky-type stains in haematology," *The Histochemical Journal*, vol. 10, no. 2, p. 1–29, 1978.
- [39] B. H. O'Connor, *A color atlas and instruction manual of peripheral blood cell morphology /*. Baltimore :: Williams and Wilkins., c1984., originally presented as the author's thesis (master's—Quinnipiac College, Hamden, Conn.).
- [40] N. Otsu, "A threshold selection method from gray-level histograms," *IEEE Transactions on Systems, Man, and Cybernetics*, vol. 9, no. 1, pp. 62–66, 1979.
- [41] R. D. Labati, V. Piuri, and F. Scotti, "All-idb: The acute lymphoblastic leukemia image database for image processing," in *ICIP*. IEEE, 2011, pp. 2045–2048.
- [42] S. H. Rezatofighi, K. Khaksari, and H. Soltanian-Zadeh, "Automatic recognition of five types of white blood cells in peripheral blood," in *Image Analysis and Recognition*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2010, pp. 161–172.

# DistB-SDoIndustry: Enhancing Security in Industry 4.0 Services based on Distributed Blockchain through Software Defined Networking-IoT Enabled Architecture

Anichur Rahman<sup>1</sup>

Department of Computer Science and Engineering  
Mawlana Bhashani Science and Technology University  
Tangail, Bangladesh  
National Institute of Textile Engineering and Research  
(NITER), Savar, Dhaka, Bangladesh

Umme Sara<sup>2</sup>

Department of Computer Science and Engineering  
National Institute of Textile Engineering and Research  
(NITER), Dhaka, Bangladesh

Dipanjali Kundu<sup>3</sup>

Department of Computer Science and Engineering  
National Institute of Textile Engineering and Research  
(NITER), Savar, Dhaka, Bangladesh

Saiful Islam<sup>4</sup>

Department of Computer Science and Engineering  
International Islamic University Chittagong  
Chittagong, Bangladesh

Md. Jahidul Islam<sup>5</sup>

Department of Computer Science and Engineering  
Green University of Bangladesh  
Dhaka, Bangladesh

Mahedi Hasan<sup>6</sup>

Department of Computer Science and Engineering  
National Institute of Textile Engineering and Research  
(NITER), Dhaka, Bangladesh

Ziaur Rahman<sup>7</sup>

Department of Information and Communication Technology  
Mawlana Bhashani Science and Technology University  
Tangail, Bangladesh

Mostofa Kamal Nasir<sup>8</sup>

Department of Computer Science and Engineering  
Mawlana Bhashani Science and Technology University  
Tangail, Bangladesh

**Abstract**—The concept of Industry 4.0 is a newly emerging focus of research throughout the world. However, it has lots of challenges to control data, and it can be addressed with various technologies like Internet of Things (IoT), Big Data, Artificial Intelligence (AI), Software Defined Networking (SDN), and Blockchain (BC) for managing data securely. Further, the complexity of sensors, appliances, sensor networks connecting to the internet and the model of Industry 4.0 has created the challenge of designing systems, infrastructure and smart applications capable of continuously analyzing the data produced. Regarding these, the authors present a distributed Blockchain-based security to industry 4.0 applications with SDN-IoT enabled environment. Where the Blockchain can be capable of leading the robust, privacy and confidentiality to our desired system. In addition, the SDN-IoT incorporates the different services of industry 4.0 with more security as well as flexibility. Furthermore, the authors offer an excellent combination among the technologies like IoT, SDN and Blockchain to improve the security and privacy of Industry 4.0 services properly. Finally, the authors evaluate performance and security in a variety of ways in the presented architecture.

**Keywords**—IoT; SDN; BC; AI; security; privacy; industry 4.0

## I. INTRODUCTION

Mostly, the monitoring and control mechanism in industries that manages specifics of commodity production, inventory data, knowledge of employees working in the supply chain, is typically time-intensive, costly, and sluggish. Industry 4.0, often referred as the Industrial Internet of Things (IIoT), is a modern step of the Industrial Revolution, focused extensively on interconnectedness, robotics, artificial intelligence, and real-time data. In addition, the smart industry has four intelligent features. Firstly, sensors that make decisions to alter behaviour based on environmental changes. Secondly, it has internet connectivity and real-time access. Then, IR 4.0 is strongly compatible with robot vision systems and AI techniques. The last one is Virtual Reality (VR) strategies that enable human-machine interaction in IR-4.0 [1]. However, the IT protection issue is the most daunting part of the adoption of industry 4.0 platforms or strategies. Further, the key issues in Industry 4.0 are the absence of granularity of knowledge and real-time tracking [2].

On the contrary, the SDN [3], [4] is an evolving technology that enables the implementation of a protected Industry 4.0 manufacturing environment [5], because low-level computing functions in SDN are far more effective, authentication functions are embedded in the network rather than being centralized in individual network components. Hence the key aim of introducing SDN is to minimize public response time and continuous availability. These technologies do have specific types of implementations in various smart technical fields, such as smart buildings, grids, healthcare, industries, and many more. Furthermore, research into security mechanisms is critical for next-generation IoT and the creation of advanced confidentiality defence schemes to tackle numerous attacks on IoT networks. To deliver influential functionality such as continuing anonymity, verification, and robustness; Blockchain technology is a secure solution [6]. In addition to this, Blockchain accounts for new technology and developments of the future. The invention of the Blockchain is a radical change of the conventional societal structure and style of service. Due to the popularity of Ethereum and Bitcoin cryptocurrency, Blockchain has gone authoritative in the field of network security. It creates immutable data structures that cannot be hacked, changed. The process of new block data appendage is performed by some proof of work and acknowledging the current blocks of data appended already in the decentralized public ledger of the Blockchain [7].

Several researchers have investigated the importance of various innovations for developing the smart industry or IR 4.0 [8], [9]. Nevertheless, such innovations cannot address the current problems of Industry 4.0 alone, such as stability, data protection, etc. On the other hand, utilizing all the new technology together, such as SDN, IoT, Blockchain, deep learning, etc. techniques would contribute to complicated structures. However, the IoT, Blockchain, SDN innovations are integrated to reach a higher degree of operating performance, profitability, transparency, protection, and privacy in Industry 4.0. Within IR 4.0 major focus is paid to SDN, Blockchain and IoT, etc.

After analyzing the above discussion, this study proposed a model based on distributed Blockchain through SDN-IoT enabled architecture to ensure adequate security which is the primary concern in industry 4.0. Moreover, the authors focus on the significant concerning issues of Industry 4.0 applications like security, privacy, confidentiality efficiently.

The main contributions of the paper are following:

- This study proposes a framework “DistB-SDoIndustry” focused on both SDN-IoT and Blockchain technologies to control more securely in Industry 4.0 services.
- We also address Blockchain technology for data validation, verification, broadcasting, and so on. Additionally, it is capable of providing data with a confidential route to enter the desired cloud locations efficiently.

The remaining sections of the paper are structured as fol-

lows. The authors present related works in Section II. Section III highlight the proposed “DistB-SDoIndustry” architecture for industry 4.0 security management. After that result analysis and discussions are also presented in Section IV. Finally, this work concludes with the significance and future ramifications in Section V.

## II. MOTIVATIONAL BACKGROUND AND LITERATURE REVIEWS

### A. Background Study

In this section, the authors discuss the IoT, SDN and background knowledge of Blockchain technology with Industry 4.0 applications briefly.

1) *IoT with SDN technologies:* The IoT includes various types of information advancing devices like routers, different environment detector which can perceive data from surroundings and pass the data to the next level of digital system [10]. IoT can choke the data collected from the natural world but not securely and conveniently. To handle the information smartly, the compound structure of IoT, including Software Defined Network (SDN) is helpful. SDN controls the information through multiple central processors and makes the system pliant with the reward of programmability [11], [12]. It could have pertained as a secret agent between the data control layer and IoT by which information is collected.

2) *Overview of Blockchain Concept:* Blockchain [13] is a system of recording data in a way that makes it difficult or impossible to change, hack, or cheat the system. A Blockchain is a decentralized, distributed, and public digital which consists of blocks. In general, every block is connected between them and sets of timestamped transactions. The chain of blocks, or Blockchain, serves as a publically accessible digital ledger. In this technology node exchange data by creating a transaction and each transaction depends on the previous transaction, where one transaction outputs are connected in another transaction as inputs hence forming chain among them. The Blockchain representation is shown in Fig. 1. The first block is called a generic block, and the rest of the blocks create by participating nodes called miners to try to solve a cryptographic puzzle named Proof of Works. Hence, participating nodes create a trusted network over untrusted participants in the network. New transactions are verified by all participating nodes that exclude the requirement of the central dependency and propose a distributed management system. Each block holds the hash of its previous block which assures the integrity of the transaction; therefore, make sure no alteration of the block in the network. If one transaction is valid, then the transaction is continuously stored in the Blockchain network that can be reached by any node. All transactions in this network are signature using public-key cryptography so that the authenticity nature of Blockchain is fulfilled [14].

3) *Blockchain for Industry 4.0:* Industry 4.0 aims to assemble, study, disc the information of individuals and assure the activities in real-time. Today’s world is having a lot of

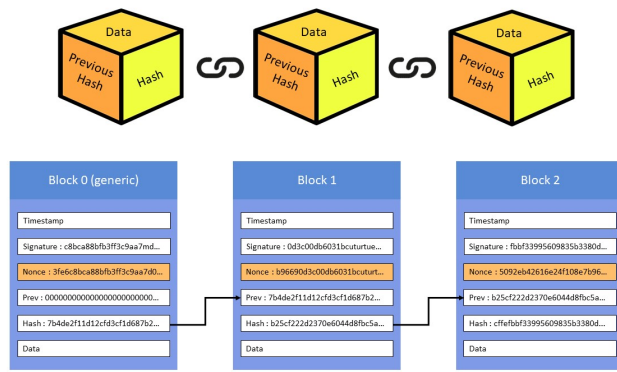


Fig. 1. Concept of Blockchain

industries, and each industry deals with a vast amount of products and customers. This information helps the industry to take the decision about manufacturing products in future. It is quite backbreaking to take care of the information of yield and clients and predicting the doings [15]. IoT with other relevant technologies extensively fires the 4th industrial revolution [16]. It can be considered as a bright set up of a factory and extraordinarily dynamic and automated production mesh. Moreover, it is to develop the manual manufacturing process into an elastic and self-coordinating production channel [17]. In a real-time machine, it is quite risky to store and control data as a lot of actuators are collecting information from different places. Blockchain comes up at this point. As mentioned earlier Blockchain can secure data through its unique structure and processing power. It is not impossible to litigate a transaction of the product without any human resources using Blockchain.

### B. Literature Review

Several researchers have been proposed in recent years based on Blockchain and SDN with IoT technologies in various purposes. Throughout this portion, we are going to present similar studies based on IoT convergence with SDN and Blockchain for Industry 4.0 applications.

In [18] through this contribution, the enhanced framework for considering the smart industry can be demonstrated by using IoT. Moreover, the authors taken an energy-efficient approach into account for Industry 4.0 but not mentioned stability. On the similar research, an IoT network protection framework by using the core technologies provided by the SDN proposed by [19] such an advanced protection strategy including IoT system authentication and authorizing approved flows will help secure IoT networks from malicious IoT devices and attacks. However for Peer to Peer (P2P) connectivity between IoT systems and SDN controllers, [20] proposed design that utilizes public and private Blockchains, to remove Proof-of-Work (POW), in addition, the proposed model has used cluster configuration routing algorithm to maximize energy usage and improve protection. On the other hand, in [6] introduced a protection architecture and applied it

as a Blockchain-based Platform-as-a-service (PaaS) paradigm to validate the data secrecy of authenticated users while applications are on the move and provide a robust solution to threat detection for usage in the IoT context. The proposed model for protection is efficient. From a virtualization viewpoint, it can be improved, maintaining certain protection features such as secrecy, transparency, etc. Again, the similar work [21] authors proposed “DistBlockBuilding” framework to handle the stable and efficient movement of data from one surface to another, also assessed the efficiency of a protected network based on IoT-SDN architecture. They suggested a cluster head selection algorithm; in addition, they included energy-saving and load balancing resources for SDN-IoT infrastructure in distributed Blockchain-based network. Further, [22] authors developed 5G network intrusion detection and mitigation technologies for the SDN/NFV cloud. Additionally, [23] focused Blockchain definition that can be converged to an SDN based IoT framework to strengthen its protection aspects further.

In another research growth, in [24] discusses the effects of Blockchain for IoT also stresses how IoT can take advantage of the Blockchain’s decentralized, arbitrary and transparent nature to improve agility in asset management. But they overlooked Blockchain’s technological qualities. Therefore, the practicalities of applying the core characteristics have not been assessed. Moreover, Rane et al. [25] suggested Man4Ware architecture with the supplementary Blockchain infrastructure. Also, the creation of smart manufacturing software compatible with the Industry 4.0 vision with “Man4Ware” support will generate enough incentives for innovative and insightful apps. But the proposed architecture does not include more sophisticated functionality and help for technological innovations and Blockchain apps. To render Blockchain entirely available and customizable [17] reviewed the current research on the applicability of Blockchain in various IIoT-specific industries. In Industry 4.0 and IIoT, they looked at the different industrial application of Blockchain to include an abstract indicator of acceptance in practice; Then, they added that in order for Blockchain to be completely functional and scalable, industry-oriented work would also tackle several of these problems, including personal data security, block system scalability, the participating organization’s data confidentiality and safety, Blockchain deployment and implementation costs, and policy regulations. Besides, [26] have defined IIoT technologies and have addressed the state-of-the-art protection vulnerabilities in Industry 4.0. Furthermore, Oztemel et al. [27] analyzed Industry 4.0, and associated innovations also [28] reviewed some of the cyber-security threats included in the advanced industrial production and Industry 4.0 path, further the relevant preventive measures currently adopted or under implementation. On the similar work to introduced an energy-efficient and QoS-aware parallel routing optimization algorithm for Software-Defined IIoT focused on healthcare systems, a scalable routing scheme for an outsized IIoT network has proposed [29] which is eight times faster than the existing programmes. Moreover, a quick parallel online routing optimization architecture is pro-

posed for SDN-enabled smart healthcare networks supported IIoT.

In summary, existing researches have focused on many branches of Blockchain, IoT-SDN technologies. But only a few number of researchers have addressed Industry 4.0 applications. As the field is still new and, is facing many threats in security correspondence. That's why, we attempt to minimize challenges such as security threats in the Industry 4.0 applications.

### III. PROPOSED "DISTB-SDOINDUSTRY" FOR INDUSTRY 4.0 APPLICATIONS

From the above discussed sections we can conclude that Industry 4.0 is vulnerable to security issues. To manage a massive amount of data that is transferred need to be maintained efficiently and securely. In addition to this the privacy of data also need to be ensured as the data transferred over the internet is highly susceptible to attackers. For managing the different applications of Industry 4.0, the authors propose a model "DistB-SDoIndustry" based on emerging technologies like SDN-IoT and Blockchain in security purposes, which shown in Fig. 2. We consider several steps for elaborating the proposed architecture, such as SDN-IoT enabled environment with perception, control and application layers, distributed Blockchain-based security in Industry 4.0 credentials, and Industry 4.0 services and security. In addition, the whole architecture is controlled by SDN and Blockchain technologies. Indeed, SDN helps to provide the programmability and data flexibility to the Industry 4.0 environment. Then, Blockchain is capable of handling the security, privacy as well as the confidentiality of the desire networks efficiently [30].

#### A. SDN-IoT Enabled Environment

Basically, the IoT sensor is capable of sensing data for desire applications. The IoT provides these data-enabled devices like firewalls, routers, switches, and so on. On the other hand, SDN helps to ensure data security and flexibility in the IoT environment. Indeed, an SDN is organized by some distinct planes such as data, network, and application planes efficiently. In addition, discussion of these layers is given below.

1) *Perception Layer*: First of all, this layer is responsible for providing data to the Industry 4.0 management. It can collect huge data for future use in the goal system perfectly. However, the initial sources of data forwarding are treated as smart sensor, actuators, and core smart grid networks. Moreover, it can perform and monitor all data effectively. In SDN platform, it estimates the data path management regionally on this layer using compromises node-to-node in real-time standings. On the other hand, requests for more extended bandwidth or source may be difficult to meet the device requirement in real-time, as this is unpredictable information subject to application and traffic Profile. Moreover, an SDN data plane manages the whole gathered data using a common gateway path; this gateway incorporates data as well as filtering this from external interruptions.

2) *Control Layer*: In the SDN paradigm, a most essential plane is the control path. This basically provides all benefits of the SDN platform like data control and data security also control the transmission of data from the edge layer to manage layer efficiently. Therefore, the workflow of the control layer in the SDN environment, several protocols like OpenFlow, OpenDayLight, and OpenStack used. However, it interfaces of two Application Program Interface (API), such as southbound and northbound APIs. Where the Southbound APIs features, provide connectivity with the switch fabric, virtualized network frameworks, or the consolidation of a decentralized network of computers. On the other hand, a northbound interface is an interface that allows the communication with a higher-level feature of a specific component of a device. Nevertheless, the entire data path can be updated if the application layer demands a change based on the gathering of data, availability of contact connections and a distribution scheme.

3) *Applications Layer*: The application layer is the topmost tier of the SDN stages, and it involves data centre, servers, storage, analysis, processing, and statements. It efficiently processes all data. Further, SDN application layer collects all filtered data from SDN controllers applying suitable instructions of the control layer. After performing all operation of data in the data and control layer, the application layer consists of all filtered data suitably. Moreover, it helps to store these data in the cloud storage using a secure communication path. Then, Industry 4.0 entertains these secured data in various applications efficiently.

#### B. Distributed Blockchain-based Security in Industry 4.0 Applications

There are many challenges in the growing Industry 4.0 applications; these challenges are cost challenge, structural, technology, security challenges, and so on. Moreover, Industry 4.0 applications are facing data integrity and data redundancy problems. In this paper, the author's primary concern is security challenges. Indeed, the authors also highlighted security improvement using the Blockchain approach for Industry 4.0 services management. In the presented model, the Blockchain approach provides the solution such a concerning issues like confidentiality, access control, authorization, and integrity in the Industry 4.0 applications. Moreover, a Blockchain is a chain based strategy, which provides the chain to connect every block to each other. Each block consists of data, timestream, hash data, and so on, as depicted in Fig. 3. In the Blockchain scenario, hash data is a unique component that can point one block to another. Every block of information is connected with one to another using the hash function properly. First of all, the genesis block creates by the system, then the rest of the block continues their activity based on the genesis block. If a new block is added to the Blockchain environment, it will need get the permission of all minor blocks. When it achieves the consent of at least 60-80 percent block, it will be added in the Blockchain network as a new member securely. This process is performed very confidentially and



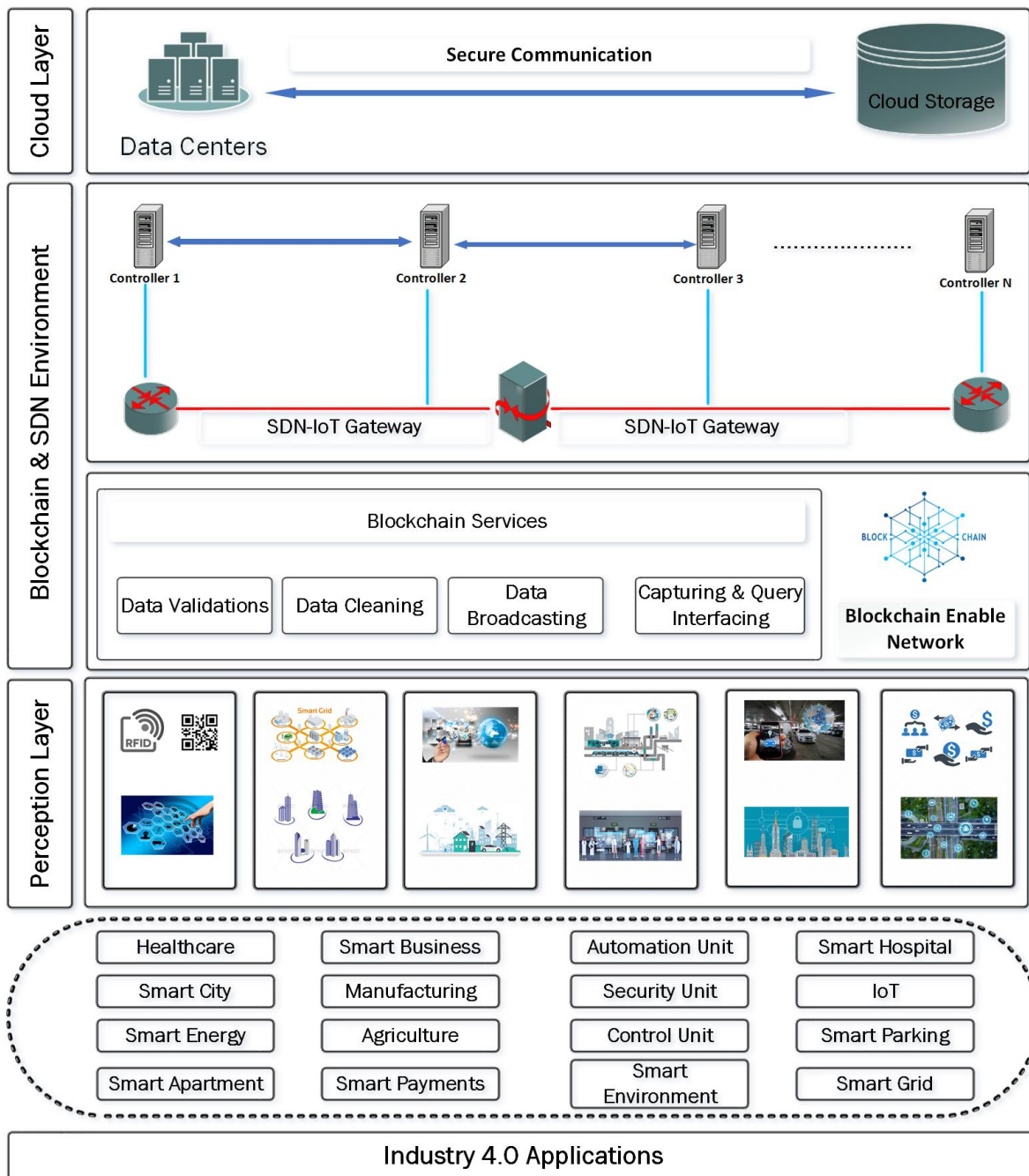


Fig. 2. Proposed Architecture of "DistB-SDoIndustry"

firmly. In contrast, there is no involvement of a third party or any intruder in this communication system. For the benefits of Blockchain, we have addressed this technology in the proposed method. In Industry 4.0 applications, all data would be safe by using Blockchain technology. Then, it provides different services like data validations, cleaning, broadcasting, as well as data capturing & query interfacing in the modern Industry 4.0 environment, as shown in Fig. 2. Furthermore, it can manage various security attacks and also able to perform a large number of operations in the Industry 4.0 applications.

### C. Industry 4.0 Services and Securities

Modern Industry 4.0 provides some benefits like efficiency in automation, the innovation of new products, costs consideration, revenues, etc. [31]. It can also offer intelligent components such as connectivity, automation, and optimization. Moreover, the IIoT consists of all sensors and machines in the Industry 4.0 platform.

In addition, automation means the digitalization of all services for Industry. After that Industry 4.0 includes some leading technologies like AI, whose main contributor is Big

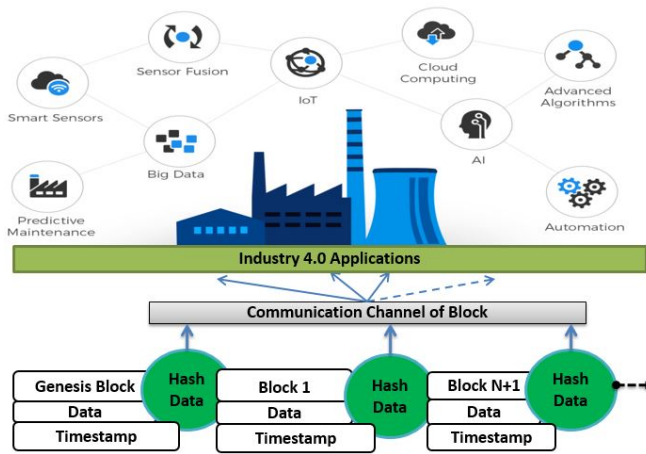


Fig. 3. Blockchain-based Security Approach in Industry 4.0

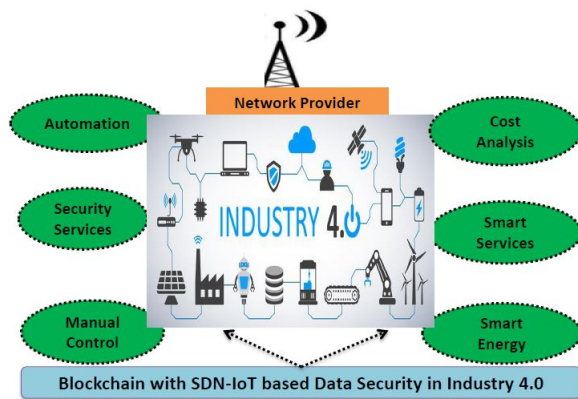


Fig. 4. Industry 4.0 System and Services

Data and Analytics, which utilizes machine learning and AI techniques efficiently. Based on these services of Industry 4.0, the authors mention different services such as security services, smart services, as well as smart energy, as shown in Fig. 4. Indeed, the authors also address various types of services such as smart cities, smart apartments, smart health-care, energy, manufacturing, agriculture, payments, hospital, business, smart grid, and so on, as shown in Fig. 2 based on Blockchain technology with SDN-IoT architecture to enhance the Industry 4.0 applications exceedingly.

#### IV. RESULTS ANALYSIS AND DISCUSSIONS

##### A. Environment Setup

In this section, we have analyzed to set up the proposed model with the Emulator (Mininet), Mininet-Wifi simulation tools for measuring the environmental activities of the SDN. Again, the OpenFlow protocol is used in the SDN context to accomplish the goal outcomes. Further, different types of packet sizes are used, such as bytes 256, 800, and 1024. Nodes turn in the rectangular field according to the configuration of the random waypoint. In which, at some random location, each node is positioned the rectangular area at simulation initialization. Throughout the simulation period, all nodes shift

in line with the movement set out in the scenario script. The nodes are allowed to move in the dimension 3000m x 3000m rectangular area. In addition, Ubuntu (GNU/Linux), x86 (2.20GHz), 8GB RAM, 2TB ROM, and other external memory used to test our desired performance. Importantly, the Wireshark platforms have been adequately utilized to visualize the IoT network performance based on the SDN platform.

##### B. Performance Analysis

In this segment, the authors have considered three parameters, such as throughput analysis, secured rate analysis and packets failure rate comparison, to evaluate the execution of the proposed system efficiently.

1) *Throughput Analysis:* The authors have analyzed the throughput based on the number of packets transmission, as shown in Fig. 5.

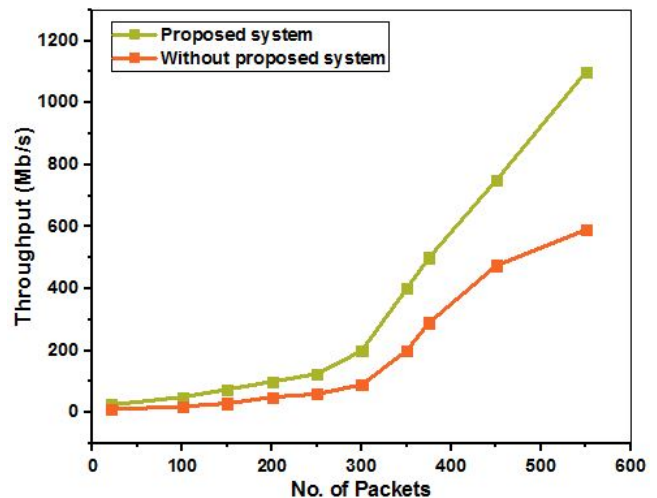


Fig. 5. Throughput Comparison

However, it displays the throughput comparisons between existing centralized model and the suggested model diagrammatically. Further, we have regarded that when the number of packets is less, then the throughput is almost the same with each other. But when the number of packets is progressing, then the throughput is also growing. After performing a particular time, we have also noticed that due to less engagement of attacks or undesired components like noise and intruder our proposed architecture “DistB-SDoIndustry” shows a much better performance than the centralized core model performance effectively.

2) *Security Rate Analysis:* Then, we have observed the security rate according to the number of data packets, as shown in Fig. 6. After that, it shows the comparisons of the protection rate between the core model and the proposed system “DistB-SDoIndustry”. Moreover, we have also noticed that when the number of data packets is less, the security rate is nearly the equivalent. But when the number of packets increases, then both of the performance is also expanding. Furthermore, due to the rise in the number of data flow after performing a specific

time, the authors have remarked that the proposed “DistB-SDoIndustry” model shows better security few involvements of attacks than the performance of the existing network.

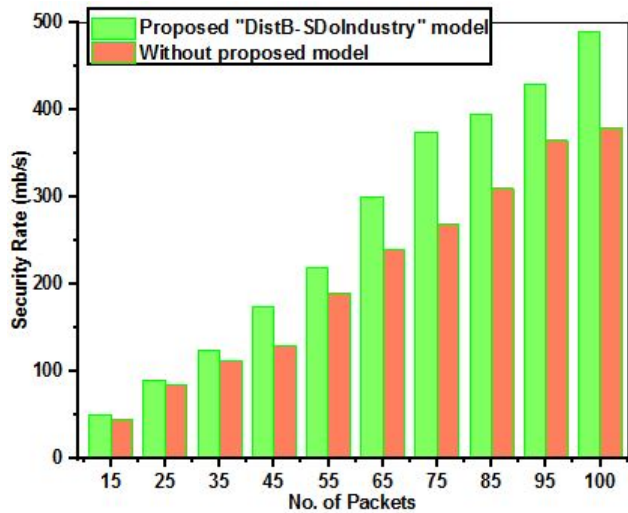


Fig. 6. Security Comparison

3) *Node Failure Rate Analysis:* On the other hand, in Fig. 7 shows the data failure rate based on the number of nodes between our proposed system and the core system—actually, node failure rate occurs based on different types of attacks. From the above analysis, it is clear that the proposed design performs the minimal node failure rate due to less impact of attacks. Further, the authors have notified that the node failure rate is initially less for both. However, with the increase in the number of nodes, the failure rate also progresses. Besides, the authors have also observed that primarily node failure rate is 2% or 3% only for both. But after completing a few times in the core scheme’s node failure rate is 90% or above, on the contrary, the proposed scheme’s failure rate is 38% to 43 % only. Thus the above analysis shows that the offered secured system “DistB-SDoIndustry” overcomes failure rate significantly than the traditional core model.

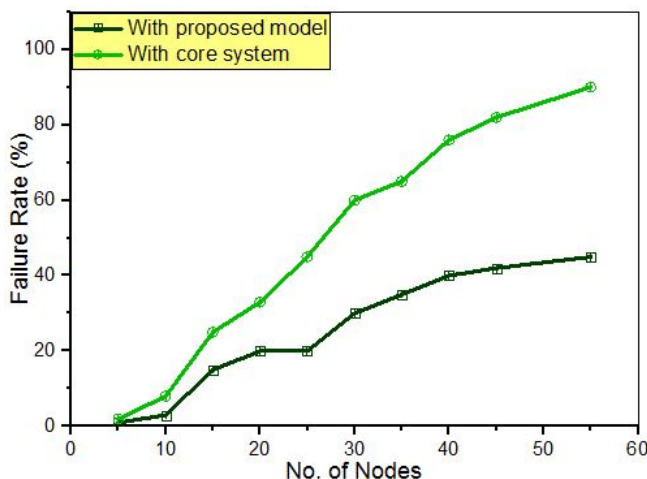


Fig. 7. Nodes Failure Comparison

## V. CONCLUSION

This paper presents a model “DistB-SDoIndustry” based on distributed Blockchain technology with SDN-IoT architecture for Industry 4.0 applications. Basically, we have highlighted two technologies such as SDN and Blockchain, in order to provide robust privacy and reliable security in the Industry 4.0 environment efficiently. Moreover, we have considered the SDN-IoT model for dividing our whole architecture into different layers securely. After that, we have also addressed Blockchain for improving data security and confidentiality. There is no interference with the third party in the presented system; furthermore, this paper has implemented the SDN-IoT model in different parameters like the security rate. This depends on no. of data packets and packets failure rate based on no. of nodes in the proposed networking model. Still, the implementation of Blockchain is in a developing stage. In the future, this study will be added to the complete implementation of Blockchain competently. Moreover, the authors will analyze the different types of attacks, like Denial of Service (DoS) attacks, flooding attacks, etc. In addition, the authors will also evaluate more parameters such as throughput, packet arrival rate, the response time of data, etc. as well as assess the performances in numerous ways of the presented architecture.

## REFERENCES

- [1] Z. Shi, Y. Xie, W. Xue, Y. Chen, L. Fu, and X. Xu, “Smart factory in industry 4.0,” *Systems Research and Behavioral Science*, vol. 37, no. 4, pp. 607–617, 2020.
- [2] F. S. T. da Silva, C. A. da Costa, C. D. P. Crovato, and R. da Rosa Righi, “Looking at energy through the lens of industry 4.0: A systematic literature review of concerns and challenges,” *Computers & Industrial Engineering*, p. 106426, 2020.
- [3] M. J. Islam, M. Mahin, S. Roy, B. C. Debnath, and A. Khatun, “Distblacknet: A distributed secure black sdn-iot architecture with nfv implementation for smart cities,” in *2019 International Conference on Electrical, Computer and Communication Engineering (ECCE)*. IEEE, 2019, pp. 1–6.
- [4] B. K. Mukherjee, S. I. Pappu, M. J. Islam, and U. K. Acharjee, “An sdn based distributed iot network with nfv implementation for smart cities,” in *International Conference on Cyber Security and Computer Science*. Springer, 2020, pp. 539–552.
- [5] A. Tsuchiya, F. Fraile, I. Koshijima, A. Ortiz, and R. Poler, “Software defined networking firewall for industry 4.0 manufacturing systems,” *Journal of Industrial Engineering and Management (JIEM)*, vol. 11, no. 2, pp. 318–333, 2018.
- [6] D. V. Medhane, A. K. Sangaiah, M. S. Hossain, G. Muhammad, and J. Wang, “Blockchain-enabled distributed security framework for next-generation iot: An edge cloud and software-defined network-integrated approach,” *IEEE Internet of Things Journal*, vol. 7, no. 7, pp. 6143–6149, 2020.
- [7] M. Farnaghi and A. Mansourian, “Blockchain, an enabling technology for transparent and accountable decentralized public participatory gis,” *Cities*, vol. 105, p. 102850, 2020.
- [8] S. K. Singh, Y.-S. Jeong, and J. H. Park, “A deep learning-based iot-oriented infrastructure for secure smart city,” *Sustainable Cities and Society*, p. 102252, 2020.
- [9] U. Bodkhe, S. Tanwar, K. Parekh, P. Khanpara, S. Tyagi, N. Kumar, and M. Alazab, “Blockchain for industry 4.0: A comprehensive review,” *IEEE Access*, vol. 8, pp. 79 764–79 800, 2020.
- [10] R. F. Al-Mutawa and F. A. Eassa, “A smart home system based on internet of things,” *International Journal of Advanced Computer Science and Applications*, vol. 11, no. 2, 2020. [Online]. Available: <http://dx.doi.org/10.14569/IJACSA.2020.0110234>

- [11] T. Han, S. R. U. Jan, Z. Tan, M. Usman, M. A. Jan, R. Khan, and Y. Xu, "A comprehensive survey of security threats and their mitigation techniques for next-generation sdn controllers," *Concurrency and Computation: Practice and Experience*, vol. 32, no. 16, p. e5300, 2020.
- [12] T. Adbeb, W. Di, and M. Ibrar, "Software-defined networking (sdn) based vanet architecture: Mitigation of traffic congestion," *International Journal of Advanced Computer Science and Applications*, vol. 11, no. 3, 2020. [Online]. Available: <http://dx.doi.org/10.14569/IJACSA.2020.0110388>
- [13] C.-S. Yang, "Maritime shipping digitalization: Blockchain-based technology applications, future improvements, and intention to use," *Transportation Research Part E: Logistics and Transportation Review*, vol. 131, pp. 108–117, 2019.
- [14] R. Zhang, R. Xue, and L. Liu, "Security and privacy on blockchain," *ACM Computing Surveys (CSUR)*, vol. 52, no. 3, pp. 1–34, 2019.
- [15] G. Rathee, M. Balasaraswathi, K. P. Chandran, S. D. Gupta, and C. Boopathi, "A secure iot sensors communication in industry 4.0 using blockchain technology," *JOURNAL OF AMBIENT INTELLIGENCE AND HUMANIZED COMPUTING*, 2020.
- [16] Q. Wang, X. Zhu, Y. Ni, L. Gu, and H. Zhu, "Blockchain for the iot and industrial iot: A review," *Internet of Things*, vol. 10, p. 100081, 2020.
- [17] T. Alladi, V. Chamola, R. M. Parizi, and K.-K. R. Choo, "Blockchain applications for industry 4.0 and industrial iot: A review," *IEEE Access*, vol. 7, pp. 176 935–176 951, 2019.
- [18] A. Bagdadee, L. Zhang, and M. Remus, *A Brief Review of the IoT-Based Energy Management System in the Smart Industry*, 01 2020, pp. 443–459.
- [19] K. K. Karmakar, V. Varadharajan, S. Nepal, and U. Tupakula, "Sdn enabled secure iot architecture," in *2019 IFIP/IEEE Symposium on Integrated Network and Service Management (IM)*, 2019, pp. 581–585.
- [20] A. Yazdinejad, R. M. Parizi, A. Dehghantanha, Q. Zhang, and K. R. Choo, "An energy-efficient sdn controller architecture for iot networks with blockchain-based security," *IEEE Transactions on Services Computing*, vol. 13, no. 4, pp. 625–638, 2020.
- [21] A. Rahman, M. K. Nasir, Z. Rahman, A. Mosavi, S. Shahab, and B. Minaei-Bidgoli, "Distblockbuilding: A distributed blockchain-based sdn-iot network for smart building management," *IEEE Access*, 2020.
- [22] I. H. Abdulqadder, S. Zhou, D. Zou, I. T. Aziz, and S. M. A. Akber, "Multi-layered intrusion detection and prevention in the sdn/nfv enabled cloud of 5g networks using ai-based defense mechanisms," *Computer Networks*, p. 107364, 2020.
- [23] F. H. Pohrmen, R. K. Das, and G. Saha, "Blockchain-based security aspects in heterogeneous internet-of-things networks: A survey," *Transactions on Emerging Telecommunications Technologies*, vol. 30, no. 10, p. e3741, 2019.
- [24] S. B. Rane and Y. A. M. Narvel, "Re-designing the business organization using disruptive innovations based on blockchain-iot integrated architecture for improving agility in future industry 4.0," *Benchmarking: An International Journal*, 2019.
- [25] J. Al-Jaroodi, N. Mohamed, and I. Jawhar, "A service-oriented middleware framework for manufacturing industry 4.0," *ACM SIGBED Review*, vol. 15, no. 5, pp. 29–36, 2018.
- [26] I. Jamai, L. B. Azzouz, and L. A. Saïdane, "Security issues in industry 4.0," in *2020 International Wireless Communications and Mobile Computing (IWCMC)*. IEEE, 2020, pp. 481–488.
- [27] E. Oztemel and S. Gursev, "Literature review of industry 4.0 and related technologies," *Journal of Intelligent Manufacturing*, vol. 31, no. 1, pp. 127–182, 2020.
- [28] J. Prinsloo, S. Sinha, and B. von Solms, "A review of industry 4.0 manufacturing process security risks," *Applied Sciences*, vol. 9, no. 23, p. 5105, 2019.
- [29] F. Naeem, M. Tariq, and H. V. Poor, "Sdn-enabled energy-efficient routing optimization framework for industrial internet of things," *IEEE Transactions on Industrial Informatics*, 2020.
- [30] A. Rahman, M. J. Islam, F. A. Sunny, and M. K. Nasir, "DistBlockSDN: A Distributed Secure Blockchain based SDN-IoT Architecture with NFV Implementation for Smart Cities," *In Press: International Conference on Innovation in Engineering and Technology (ICIET)*, vol. 23, p. 24, IEEE, 2019.
- [31] R. Brozzi, D. Forti, E. Rauch, and D. T. Matt, "The advantages of industry 4.0 applications for sustainability: Results from a sample of manufacturing companies," *Sustainability*, vol. 12, no. 9, p. 3647, 2020.

# Small-LRU: A Hardware Efficient Hybrid Replacement Policy

Purnendu Das<sup>1</sup>, Bishwa Ranjan Roy<sup>2\*</sup>  
Department of Computer Science  
Assam University Silchar  
INDIA

**Abstract**—Replacement policy plays a major role in improving the performance of the modern highly associative cache memories. As the demand of data intensive application is increasing it is highly required that the size of the Last Level Cache (LLC) must be increased. Increasing the size of the LLC also increases the associativity of the cache. Modern LLCs are divided into multiple banks where each bank is a set-associative cache. The replacement policy implemented on such highly associative banks consume significant hardware (storage and area) overhead. Also the Least Recently Used (LRU) based replacement policy has an issue of dead blocks. A block in the cache is called dead, if the block is not used in the future before its eviction from the cache. In LRU policy, a dead block can not be remove early until it become LRU-block. So, we have proposed a replacement technique which is capable of removing dead block early with reduced hardware cost between 77% to 91% in comparison to baseline techniques. In this policy random replacement is used for 70% ways and LRU is applied for rest of the ways. The early eviction of dead blocks also improves the performance of the system by 5%.

**Keywords**—Replacement policies; cache memories; last level cache; hardware overheads; dead block

## I. INTRODUCTION

Replacement policy plays the most significant role in the performance of highly set-associative cache architecture. In multi-level cache, the first level cache (L1) is allotted as private cache to individual core whereas the large last level cache (LLC) is shared by all the cores. To reduce the access latency the LLC is also divided into multiple banks where each bank is a set-associative cache. The data distribution among the banks is based on different data mapping policies [1]. In this work we consider Static Non Uniform Cache Access (SNUCA) where each block has a fixed bank to be mapped and the bank is called the *home-bank* of the block [1]. Fig. 1 shows a multicore processor having 4 cores. Each core has a private L1 cache and a part of shared L2 cache. To make the design simple we have not divided the L2 into more than 4 banks. The rest of the paper follows the same multicore processor as shown in Fig. 1. Each bank is a set-associative cache as shown in Fig. 2.

Today's data intensive applications demand larger and higher associative cache (specially LLC). These highly associative cache reduces the conflict misses and hence improves the performance of the system. But these highly associative cache has some overheads in terms of hardware. One such

overhead is because of maintaining replacement policy for such highly associative banks. In set-associative cache (or bank), each set maintains its own replacement policy. This policy is required to replace an existing block from the cache.

As mentioned above, each set in the set-associative cache has its separate replacement hardware. For an  $N$ -way set associative cache, each set maintains separate hardware for its replacement policy. To insert a newly incoming block in the set, one of the existing block need to be evicted first known as the *victim block*. The purpose of replacement policy is to select a victim block to replace it with the recently requested block. One of the most popular and well known replacement policy is called Least Recently Used (LRU) policy. The concept of this technique is well known and not necessary to discuss here. To maintain LRU policy in each set having  $N$  ways, each way must uniquely represent its age relative to the other ways. For example, if there are only 2 ways in a cache then each set needs only 1 bit to maintain the relative age. Bit-0 means old and bit-1 means new. Similarly 2 bits are required to uniquely maintain the relative age in case of a 4-way set associative cache. Hence for a  $N$ -way set associative cache,  $\log_2 N$  bits are required to maintain the relative age of each ways in the set. The details about the hardware overheads required to maintain replacement policy is discussed in Section II.

The three important operations of any replacement policy are:

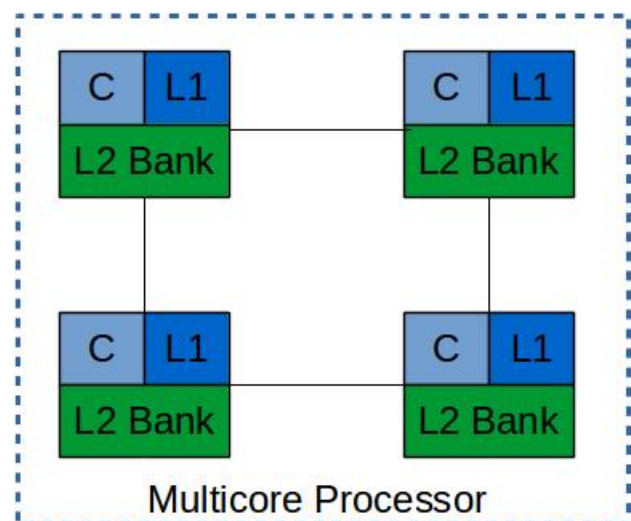


Fig. 1. An Example of Multicore Processor having 4 cores and 4 LLC banks.

\*corresponding author

- **Eviction:** During the replacement process which block should be evicted from the cache. In case of LRU replacement policy, the eviction mechanism always selects the least recently used block as a victim block.
- **Insertion:** The newly incoming block replaces the victim block in the cache. But its position w.r.t. the replacement policy is determined by its insertion mechanism. In case of LRU replacement policy the newly incoming block is always placed in the MRU position.
- **Promotion:** This operation handles about what to do when a block is being accessed from the cache. In case of LRU replacement policy, the promotion mechanism makes such blocks as MRU. It means that if a block  $B$  present in the bank and a request has been made to access the block. In this case the block will be accessed from the bank and hence it considered as hit. But after the access, the block will be moved to the MRU position.

The main reason to use such highly associative caches even in presence of hardware overheads is the performance. A highly associative cache reduces the conflict misses in the system and hence improves the system performance. These motivates researchers to reduce the hardware overhead of such highly associative caches. In this work we are mainly targeting the hardware overhead of the replacement policy of such cache.

The LRU based replacement policy are simple but faces a problem of dead blocks [2], [3], [4], [5], [6]. Since the insertion mechanism of LRU inserts a block at the MRU position and the eviction mechanism selects the victim from the LRU position, it takes long time in a highly associative cache to make a block LRU from MRU. It has been found that there are some blocks which are accessed only once in the cache. Such blocks though never used again but not possible to remove until they become LRU. Such block are called *dead-blocks*. Alternatively, a block is called dead-block at a particular instance, if the block will never used in the future before evicting it from the cache. The LRU based policy faces the issue of dead block and many technology has already been proposed to reduce

the presence of such blocks [7]. Most of these dead block prediction policies are costly in terms of storage capacity as they require to maintain additional bits for prediction. Our proposed work is capable of early eviction of dead block with minimized hardware cost. Though the proposed policy is not as smart as [7] in terms of dead block prediction but managed to reduce hardware cost significantly. The proposed policy attempts to reduce hardware cost of LRU based techniques with high dead block prediction ability so that hit rate of the memory can be improved

The organization of the paper is as follows. The next section discuss about the background and related works. The proposed Small-LRU is discussed in Section III. Section IV gives the experimental analysis and finally Section V concludes the paper.

## II. BACKGROUND

The simplest replacement policy is known as FIFO (First In First Out) which uses a straight forward strategy to replace victim block while the most widely used traditional replacement policy is LRU (Least Recently Used) which selects a victim block based on reference history. Some similar policies are MRU (Most Recently Used) and Random. In case of  $N$ -ways set-associative cache, LRU policy require  $N \times \log_2 N$  bits to represent a set. For example, a 4-ways set-associative cache require  $\log_2 4 = 2$ bits to represent a single way and  $4 \times 2 = 8$  bits require to represent a complete set. most of these traditional policies require the same hardware cost to implement. Random policy does not require any additional bit to select victim block but not suitable for general purpose.

Replacement policy is major area of research from the last two decades. The work in replacement policy can be divided into the following two categories: (a) Performance oriented like [7], [8], [9], [10] (b) Overhead reduction oriented like [11], [12]. The most efficient replacement policy was proposed in 1965 [13], which states that the evicted block must be the block which will reuse in the furthest future. The policy is considered as the optimal replacement policy. Unfortunately it is not possible to implement this optimal policy in any physical computer as it needs the knowledge of future. Hence all the practical replacement policy proposed are trying to come closer to this policy in terms of performance. Note that, in case of replacement policy, the performance means reduction in the number of cache misses. There is a huge gap still exist between all the implementable replacement policies and the Optimal replacement policy [7], [11].

Many recent replacement policies have attempted to mimic the functionality of the optimal replacement policy using the modern techniques like Machine Learning and Artificial Intelligence [7], [10]. Even after all such attempts the gap is still exists. In [10] the authors have introduced a replacement policy which can predict future based on its past experiences. Another similar technique has been proposed in [14].

Removing dead blocks from the cache is also a major responsibility of the replacement policies. The LRU based replacement policies fail to remove dead blocks early [7]. Since detecting dead block also needs the knowledge of future it can only be predicted with some efficient prediction mechanism. Some well known dead block prediction based replacement

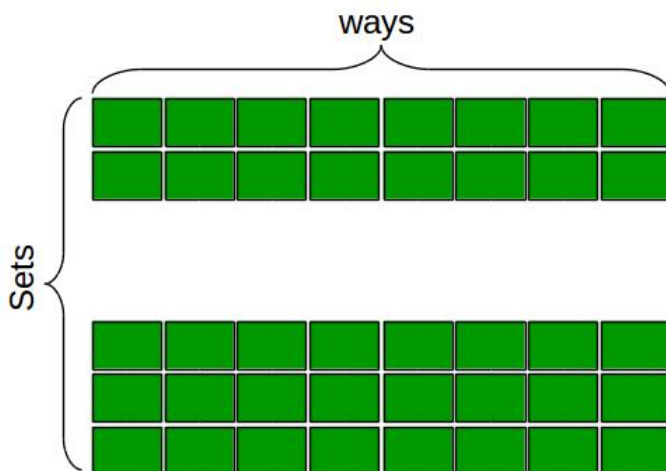


Fig. 2. An Example of 8-way Set Associative L2 Bank as shown in Fig. 1.

policies are [7], [8], [15], [16]. Most of these techniques has high hardware overhead as they need to maintain some prediction tables. Some low-overhead based dead block predictors are [11], [12]. The technique proposed in this paper is also a low-overhead based replacement policy.

### III. SMALL-LRU

In this work a highly associative set is divided into two parts: LRU-Part and Random-Part. The LRU-Part maintains LRU replacement policy while the Random-Part maintains random replacement policy. LRU-Part is relatively small and have only 30% of the total ways. Thus the technique is called Small-LRU. During the eviction of a block, the Random-Part selects a victim block randomly and place it to the LRU-part. To accommodate the block coming from Random-Part, LRU-Part removes its least recently used block. An example of Small-LRU is shown in Fig. 3.

The main advantage of Small-LRU can be divided into two parts: (a) Reducing hardware overhead and (b) Reducing the presence of dead blocks in the cache. The hardware overhead is largely reduced as the random replacement policy needs almost zero overhead. The dead blocks are present in any large sized cache memories. Since 70% of the cache is random the dead blocks will be evicted early from the set. On the other hand, if there is a highly used block being randomly picked by the random-part then block will be given chances by placing it in the MRU position of the LRU-Part. Another hit to the block will again move the block into the Random-Part.

The highly associative sets are used to remove the conflict misses in the cache. The proposed design is still equally capable of removing the conflict misses but with significantly reduced hardware overheads. Experimental analysis as discussed in Section IV shows that the proposed mechanism is good for most of the benchmarks.

#### A. Replacement Operations

**Insertion Policy:** Every time a newly inserted block will be placed in the Random-Part. The insertion policy needs to move a block from Random-Part to the LRU-Part for making room for the incoming block. To place the migrated block from Random-Part to the LRU-Part, the LRU-Part removes its LRU block.

**Promotion Policy:** If block from LRU-Part is need to be promoted then it will be placed on the random-part. Promoting a block from LRU-Part to Random-Part also needs some adjustments. Before doing such promotion a random block

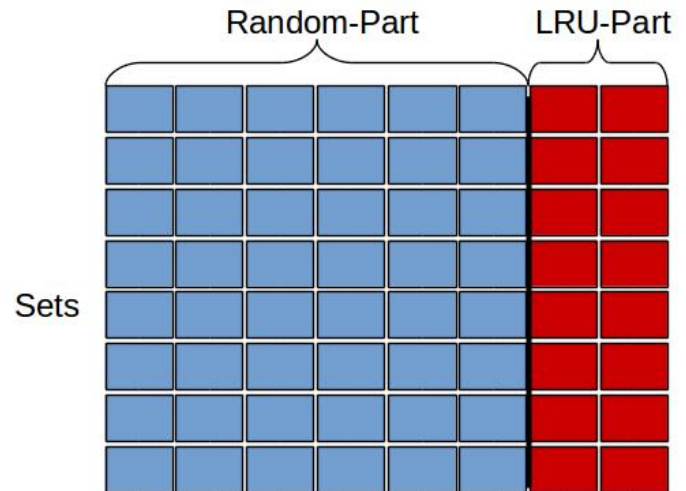


Fig. 3. An Example of Small-LRU, Implemented in a 8-way Set-Associative Cache.

from the random-part will be moved to the LRU position of the LRU-Part.

**Eviction Policy:** Every time the evicted block is the LRU block of the LRU-Part.

#### B. Advantage of Small-LRU

1) **Hardware Cost:** As mentined in Section II, LRU policy requires to maintain  $N \times \log_2 N$  bits to represent each set of an  $N$ -ways set-associative cache. Since Small-LRU use random replacement policy on 70% ways and LRU only on 30% ways, the hardware overhead is reduces by more than 70%. In Small-LRU the number of additional bits required is  $M \times \log_2 M$ , where  $M = 0.3 \times N$ . Table I represents the improvement achived in terms of hardware cost. It is observed that Small-LRU policy have reduced storage cost of LRU policy by 77% to 91% depending on the degree of associativity.

2) **Dead Block Prediction:** Early prediction of dead block is always a challenging task as discussed in Section I. Removing a dead block early from the cache is always a better choice. In Small-LRU, R-Part randomly moves a block to the MRU position of LRU-Part. The LRU-Part has only 30% ways from the total ways available in the cache. Hence even the original cache is highly associative a dead block can be removed early from the cache. A smarter replacement policy like DIP [9] or RRIP [15] in the LRU-Part may help to remove dead blocks even better.

TABLE I. COMPARISON OF THE BITS REQUIRED TO IMPLEMENT ORIGINAL LRU POLICY AND THE PROPOSED SMALL-LRU.

| Associativity | Bits Required in Original LRU | Bits Required in Small LRU | Reduction in Small-LRU |
|---------------|-------------------------------|----------------------------|------------------------|
| 8             | 24                            | 2                          | 91%                    |
| 16            | 64                            | 8                          | 87%                    |
| 32            | 160                           | 29                         | 82%                    |
| 64            | 384                           | 81                         | 78%                    |
| 128           | 896                           | 199                        | 77%                    |

TABLE II. SPECIFICATIONS USED FOR DESIGNING SMALL-LRU.

| Specification        | Values                                  |
|----------------------|---|
| Cores used           | 4                                       |
| Levels used in cache | 2                                       |
| Private cache        | L1                                      |
| Shared cache         | L2 (total 4 banks)                      |
| L2 cache             | 512KB (per bank), 8-way set associative |
| L1 cache             | 64KB, 2-way set associative             |
| Size of cache-block  | 64B                                     |

#### IV. EXPERIMENTAL ANALYSIS

System Architecture is implemented in gem5, a full-system simulator [17]. We have simulated a multicore processor (4 cores) with two level of cache memory. The upper level cache L1 is used as private cache to each core and last level cache LLC is shared among the cores. The baseline replacement policy as well as the proposed policies are implemented using Ruby module. PARSEC benchmark [18] applications are simulated in ALPHA system architecture.

To analyze the performance, all the benchmark applications are executed on the target machine designed using proposed replacement policy as well as the baseline replacement policies 200 million cycles. the specification of target machine used to implement Small-LRU is shown in Table II.

##### A. Result Analysis with Baseline-1

We considered LRU policy as baseline-1 to compare the result of our proposed policy. Statistical comparison in terms of MPKI (Miss Per Kilo Instructions) is shown in Fig. 4. It is observed from the figure that that Small-LRU have reduced MPKI by removing dead block early from cache. Fig. 5 shows

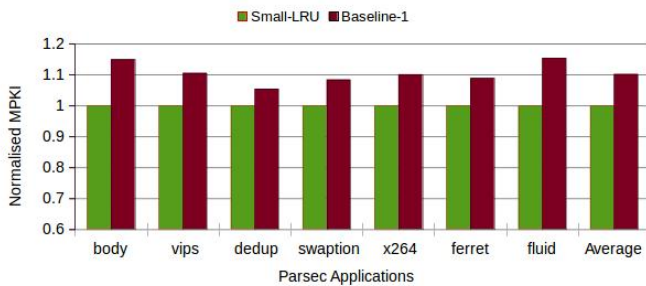


Fig. 4. Normalized Comparison of Small-LRU with Baseline-1 over MPKI.

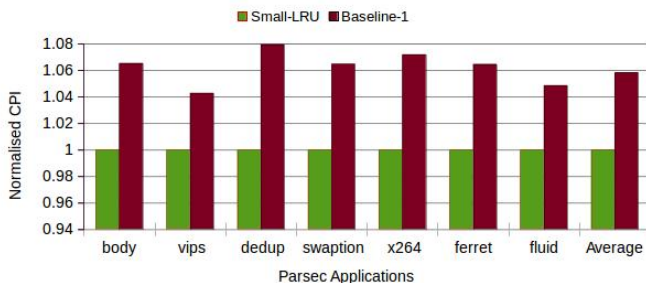


Fig. 5. Normalized Comparison of Small-LRU with Baseline-1 over CPI.

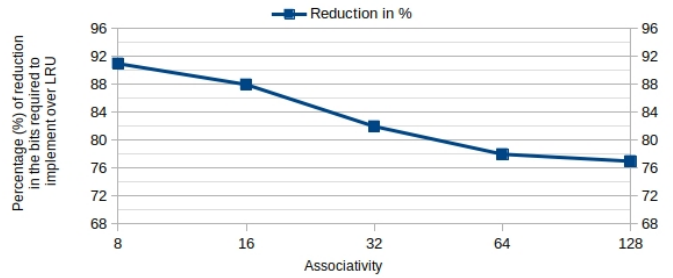


Fig. 6. The Percentage of Bits Reduction to Implement Small-LRU over Baseline-1.

the improvement in CPI (Cycle Per Instructions) due to the reduction in the number of cache miss. The proposed policy reduces MPKI by 10% and CPI by 5% on average. Though the improvement in the performance of the system is not significant with Small-LRU policy but the major advantage of this policy is the reduction of hardware cost without suffering the performance of the system. Fig. 6 depicts the percentage of bits reduction to implement Small-LRU compared to baseline-1 which is between 77% to 91%.

##### B. Result Analysis with Baseline-2

We have also compared Small-LRU with a multicore processor having 16-way associative banks. We call this design as Baseline-2. Fig. 7 shows that Small-LRU reduces the MPKI by 10% in comparison to the MPKI of baseline-2. By comparing the CPI of both the techniques it is observed that Small-LRU improves the performance of the system by 5.5% in comparison to baseline-2. Fig. 8 shows the statistical

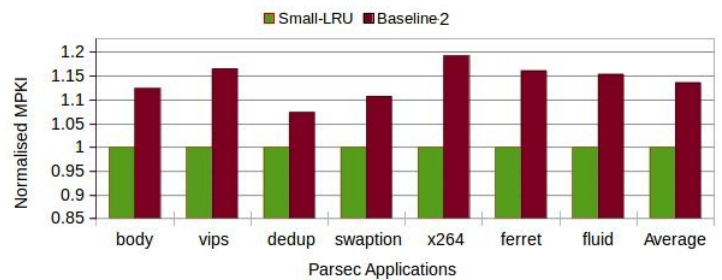


Fig. 7. Normalized Comparison of Small-LRU with Baseline-2 over MPKI.

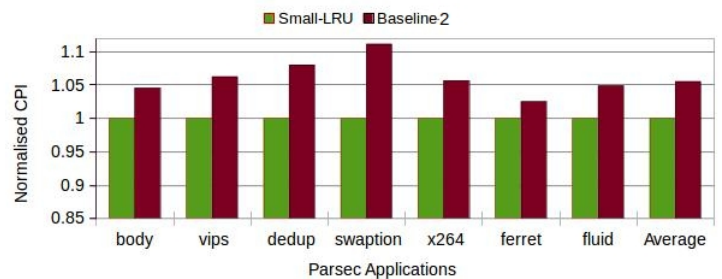


Fig. 8. Normalized Comparison of Small-LRU with Baseline-2 over CPI.



comparison of CPI between the Small-LRU and Baseline-2. Same as Baseline-1, the CPI improvement of Small-LRU over Baseline-2 is not significantly higher but enough to prove that Small-LRU reduces hardware overhead without degrading the performance.

## V. CONCLUSION

Replacement policy plays a major role in improving the performance of the modern highly associative cache memories. As the demand of data intensive application is increasing it is highly required that the size of the Last Level Cache (LLC) must be increased. Increasing the size of the LLC also increases the associativity of the cache. Modern LLCs are divided into multiple banks where each bank is a set-associative cache. The replacement policy implemented on such highly associative banks consume significant hardware (storage and area) overhead. Also the Least Recently Used (LRU) based replacement policy has an issue of dead blocks. A block in the cache is called dead, if the block is not used in the future before its eviction from the cache. Removing such dead block early from the cache is not possible in LRU policy.

In this paper we have proposed Small-LRU policy to reduce the hardware cost by more than 70% and also improves the performance by removing early dead blocks. In this policy random replacement is used for 70% ways and LRU is applied for rest of the ways. A block is always inserted in the Random-Part and evicted from the LRU-Part. Early eviction of dead blocks improves the MPKI and CPI of system using Small-LRU by 10% and 5.5%, respectively.

## REFERENCES

- [1] C. Kim, D. Burger, and S. W. Keckler, "An adaptive, non-uniform cache structure for wire-delay dominated on-chip caches," in *Proceedings of the 10th International Conference on Architectural Support for Programming Languages and Operating Systems*, ser. ASPLOS X. New York, NY, USA: Association for Computing Machinery, 2002, p. 211–222.
- [2] F. Juan and L. Chengyan, "An improved multi-core shared cache replacement algorithm," in *2012 11th International Symposium on Distributed Computing and Applications to Business, Engineering Science*, Oct 2012, pp. 13–17.
- [3] K. Morales and B. K. Lee, "Fixed segmented lru cache replacement scheme with selective caching," in *2012 IEEE 31st International Performance Computing and Communications Conference (IPCCC)*, Dec 2012, pp. 199–200.
- [4] A. Wierzbicki, N. Leibowitz, M. Ripeanu, and R. Wozniak, "Cache replacement policies revisited: the case of p2p traffic," in *IEEE International Symposium on Cluster Computing and the Grid, 2004. CCGrid 2004.*, April 2004, pp. 182–189.
- [5] W. A. Wong and J. . Baer, "Modified lru policies for improving second-level cache behavior," in *Proceedings Sixth International Symposium on High-Performance Computer Architecture. HPCA-6 (Cat. No.PR00550)*, Jan 2000, pp. 49–60.
- [6] Smith and Goodman, "Instruction cache replacement policies and organizations," *IEEE Transactions on Computers*, vol. C-34, no. 3, pp. 234–241, March 1985.
- [7] M. Kharbutli and Y. Solihin, "Counter-based cache replacement and bypassing algorithms," *IEEE Transactions on Computers*, vol. 57, no. 4, pp. 433–447, April 2008.
- [8] Y. Xie and G. H. Loh, "Pipp: promotion/insertion pseudo-partitioning of multi-core shared caches," in *ISCA*, 2009.
- [9] M. K. Qureshi, A. Jaleel, Y. N. Patt, S. C. Steely, and J. S. Emer, "Adaptive insertion policies for high performance caching," in *ISCA*, 2007.
- [10] A. Jain and C. Lin, "Back to the future: Leveraging belady's algorithm for improved cache replacement," in *2016 ACM/IEEE 43rd Annual International Symposium on Computer Architecture (ISCA)*, June 2016, pp. 78–89.
- [11] P. Das and B. R. Roy, "Splitways: An Efficient Replacement Policy for Larger Sized Cache Memory," *International Journal of Engineering and Advanced Technology (IJEAT)*, vol. 9, no. 1, 2019.
- [12] S. Das, N. Polavarapu, P. D. Halwe, and H. K. Kapoor, "Random-lru: A replacement policy for chip multiprocessors," in *LSI Design and Test*, 2013, pp. 204–2013.
- [13] L. A. Belady, "A study of replacement algorithms for a virtual-storage computer," *IBM Systems Journal*, vol. 5, no. 2, pp. 78–101, 1966.
- [14] A. Jain and C. Lin, "Rethinking belady's algorithm to accommodate prefetching," in *2018 ACM/IEEE 45th Annual International Symposium on Computer Architecture (ISCA)*, June 2018, pp. 110–123.
- [15] A. Jaleel, K. B. Theobald, S. C. Steely, and J. Emer, "High performance cache replacement using re-reference interval prediction (rrip)," in *Proceedings of the 37th Annual International Symposium on Computer Architecture*, ser. ISCA '10. New York, NY, USA: Association for Computing Machinery, 2010, p. 60–71. [Online]. Available: <https://doi.org/10.1145/1815961.1815971>
- [16] K. J. Deris and A. Baniyadi, "Analysis of non-optimal lru decisions in high-performance processors," in *2008 International Conference on Microelectronics*, Dec 2008, pp. 458–461.
- [17] N. Binkert, B. Beckmann, G. Black, S. K. Reinhardt, A. Saidi, A. Basu, J. Hestness, D. R. Hower, T. Krishna, S. Sardashti *et al.*, "The gem5 simulator," *ACM SIGARCH Computer Architecture News*, vol. 39, no. 2, pp. 1–7, 2011.
- [18] C. Bienia, "Benchmarking modern multiprocessors," Ph.D. dissertation, Princeton University, January 2011. [Online]. Available: <http://parsec.cs.princeton.edu/>

# Parameter Estimation of the ALBA Autonomous Surface Craft

Melanie M. Valdivia-Fernandez<sup>1</sup>, Brayán A. Monroy-Ochoa<sup>2</sup>, Daniel D. Yanyachi<sup>3</sup>, Juan C. Cutipa-Luque<sup>4</sup>

Electronic Engineering Professional School<sup>1,2,3</sup>

Universidad Nacional de San Agustín de Arequipa, Av. Venezuela s/n, Cercado, Arequipa, Peru

Pedro Paulet's Astronomic and Aerospace Institute<sup>4</sup>

Universidad Nacional de San Agustín de Arequipa, Cerro San Francisco s/n, Characato, Arequipa, Peru

**Abstract**—Arequipa region holds the largest extension of the Peruvian littoral at the Pacific sea, has also fresh water resources composed of rivers and lagoons from the coast to the Andes highland. The ALBA vehicle is a low-cost autonomous surface vessel with open source architecture that is being developed to support water monitoring tasks in the region. This article deals with the nonlinear identification problem for an autonomous surface craft and the maximum likelihood estimation approach is used to estimate its parameters. The parametric nonlinear model is considered with simulated and experimental data. The results shows good fitting values when two, three and a maximum four parameters are estimated.

**Keywords**—Autonomous surface craft; parameter estimation; modeling; maximum likelihood; nonlinear; zigzag

## I. INTRODUCTION

Water, the most precious resource for human being, is being vulnerable to contamination at present since there is enough evidence [1]. The *Mar de Grau* is the Peruvian sea with abundant marine species, presents a littoral of 3079.50 km and a breadth extension of 370.4 km (200 nmi). Arequipa is the south region of Peru with the largest littoral, more than 500 km. There is an open area for a sustainable exploration and monitoring these resources and the fresh waters that feed them.

The monitoring of sea conditions is commonly carried out using manned vessels, following standard international procedures and agreements. However, these large vessels cannot work in coastal areas and estuary locals due to the risk of crashing with rocks, irregularities in seabed, and currents. The use of autonomous surface crafts (ASCs) is an alternative and has advantages, such as low dimensions, zero human risk, able to explore shallow waters.

The main contribution of this paper is to estimate a greater number of parameters for the ALBA ASC, using the maximum likelihood approach with simulated and actual data, the method used in this paper is called maximum likelihood estimation (MLE) and allows us to identify many parameters at a time, therefore MLE is used for large samples and is very versatile and accurate because it works estimating not only from the values obtained of the inertial sensor and of the position sensor so it is reliable and even being able to have initial conditions. The organization of this document is as follows: Section 1 provides topics on the importance of water and its exploration using ASCs. For this, the identification of parameters and the importance of using the MLE method is done; Section 2

presents the work done on parameter estimation and ASCs; Section 3 presents the mathematical model of the ALBA ASC in nonlinear representation; Section 4 presents the maximum likelihood parameter estimation approach to identify the ALBA ASC; Section 5 presents the experimental tests and their achieved estimated parameters; Section 6 provides the conclusion.

## II. BACKGROUND

The ASCs are executing different missions around the world and their developments involve multidisciplinary areas, such as modeling, identification, navigation, control, guidance, path planning, etc. The Charlie ASC, for instance, carries out surface micro layer sampling with its real-time platform composed of navigation, guidance and control [2]. Another ASC, powered with solar energy, has navigation, guidance and collision avoidance systems to accomplish missions of water quality and greenhouse gas emissions measurements in lakes. [3].

In [4], the authors describe the modeling and identification of an ASC in a wide range of speeds and glide conditions, obtaining good estimated parameters that have been used in the proportional-derivative (PD) controller synthesis. There are other approaches to estimate nonlinear model parameter, such as the symbolic regression using genetic programming [5]. The model parameter can also be identified using experimental towing tank and open water self-propelled tests, as in [6]. The nonlinear parameters estimation is presented in [7] using the maximum likelihood approach and applied to autonomous underwater vehicle in [8]. This method maximizes the likelihood function of innovation variables, which is the difference between the output measured variable and the output estimated variable. In [9], the authors present the recursive least squares optimization approach to determine the linear and nonlinear parameters of an autonomous underwater vehicle (AUV). In [10], the authors present the parametric identification model of a ship based on the least squares approach, validating by means of high precision of identified hydrodynamic derivatives. Identifiability property can be verified before to apply a parameter estimation approach and there are some linear algebra tools that can solve this problem indeed when the model to be identified is nonlinear. These tools can be found in [11], [12], [13], [14].

### III. ALBA ASC MODELING

ALBA is a low-cost ASC, developed on an inflatable boat, used for water quality monitoring and scientific study. This vehicle has a trolling motor and a servomechanism that changes the force direction for maneuvering. The vehicle has a control architecture composed of inertial navigation sensors, wireless communication, and microcontrollers to execute navigation, guidance and control algorithms. Fig. 1 shows the picture of the cited vehicle under Lake tests, a detailed description of its development can be found in [15]. Fig. 2 presents a diagram of the servomechanism structure developed for turning the trolling motor, with 0.2462 m height and 0.511 m long. The servomotor is located on the right side and covers a space of 0.062×0.094×0.029 m. Fig. 2(b) shows the clamps where the trolling helm is installed.

The ALBA model is expressed according to standard notation used in maritime vehicles [16]. The dynamics of 6 degrees-of-freedom (DOF) are represented with two coordinate systems, one named earth-fixed  $\eta = (x, y, z, \phi, \theta, \psi)$  and another named body-fixed  $\nu = (u, v, w, p, q, r)$  (Fig. 3). The dynamic equation is given by (1) and the kinematic transformation between earth-fixed and body-fixed frames is expressed by (2), respectively:

$$M\dot{\nu} + C(\nu)\nu + D(\nu)\nu + g(\eta) = \tau, \quad (1)$$

$$\dot{\eta} = J(\eta)\nu, \quad (2)$$

where  $J(\eta)$  is the coordinate transformations;  $M$  is the mass matrix composed of rigid body mass and added mass;  $C(\nu)$  includes terms of centripetal, Coriolis and rigid body;  $D(\nu)$  is the damping matrix and  $\tau$  is the control effort vector. Table I presents the main features of the ALBA ASC.



Fig. 1. ALBA ASC Operating in a Lake [15].

For the ALBA ASC, it is reduced in dynamics to 2-DOF that correspond to the body  $\eta = (v, r)^T$  and inertial  $\nu = (y, \psi)^T$  frames, respectively. The hydrostatic forces in these two directions are null and the strip theory is used to obtain theoretical vehicle parameters [17]. The nonlinear model is given for very small values of the propeller angle  $\delta_p$  and surge velocity constant  $u_0 = \text{cte}$ .

$$(m - Y_{\dot{v}})\dot{v} + (mx_G - Y_{\dot{r}})\dot{r} = Y_v u_0 v + Y_r u_0 r + Y_{v|v}|v| + Y_{r|r}|r| - m u_0 r + b + k_{u_0} \sin \delta_p \quad (3)$$

$$(mx_G - N_{\dot{v}})\dot{v} + (I_z - N_{\dot{r}})\dot{r} = N_v u_0 v + N_r u_0 r + N_{v|v}|v| + N_{r|r}|r| - m x_G u_0 r - l_{x_G} k_{u_0} \sin \delta_p \quad (4)$$

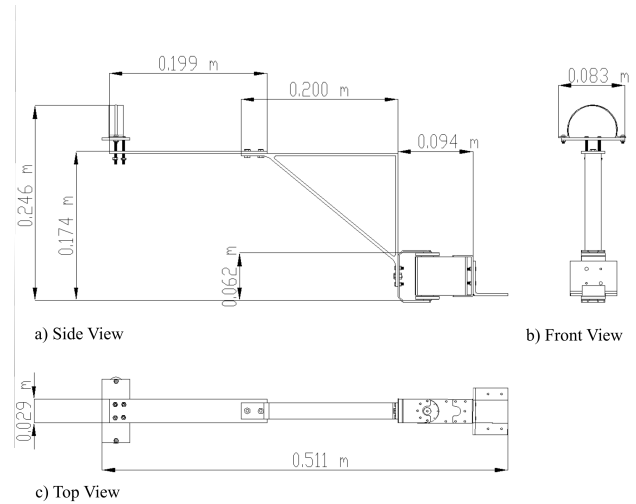


Fig. 2. Turning Propeller Mechanism.

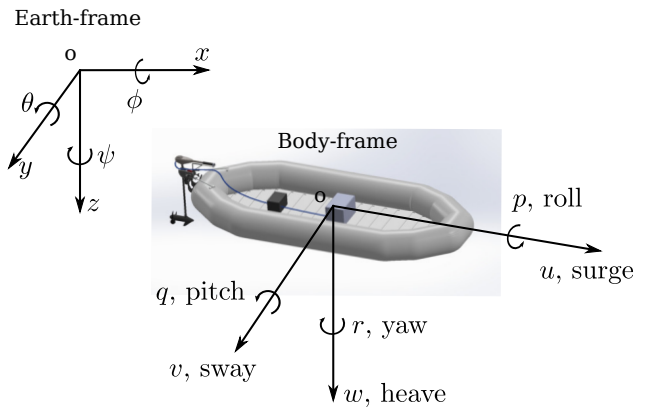


Fig. 3. ALBA ASC Coordinate System.

$$\dot{y} = u_0 \sin(\psi) + v \cos(\psi), \quad (5)$$

$$\dot{\psi} = r, \quad (6)$$

where  $k_{u_0} = 99$  is the propeller coefficient,  $l_{x_G} = 1.79$  m is the distance from the propeller location to the gravity center, and the rest hydrodynamic derivatives are described in Table II. The propeller angle in this vehicle is limited by  $\pm 25^\circ$  degrees considered sufficient to generate effort to maneuver the whole craft.

Using straightforward transformation, the nonlinear model of the ALBA ASC (3)-(6) can be expressed in the standard

TABLE I. ALBA ASC MAIN FEATURES

| Features                         | Value   |
|----------------------------------|---------|
| Length ( $L$ )                   | 3.14 m  |
| Width                            | 1.45 m  |
| Gravity center ( $x_G$ )         | 0.15 m  |
| Propeller location ( $l_{x_G}$ ) | 1.79 m  |
| Mass ( $m$ )                     | 125 kg  |
| Inertia around z axis ( $I_z$ )  | 123.84  |
| Cruise speed ( $u_0$ )           | 1 m/s   |
| Maximum speed                    | 4 m/s   |
| Autonomy                         | 3 hours |

TABLE II. ALBA HYDRODYNAMIC PARAMETERS

| Parameters       | Value     | Unit                                  | Description          |
|------------------|-----------|---------------------------------------|----------------------|
| $Y_{\dot{\psi}}$ | -571.0619 | kg                                    | Added mass           |
| $Y_{\dot{r}}$    | -101.9671 | kg.m/rad                              | Added mass           |
| $N_{\dot{\psi}}$ | -101.9671 | k.m                                   | Added mass           |
| $N_{\dot{r}}$    | -232.9491 | kg.m <sup>2</sup> /rad                | Added mass           |
| $Y_v$            | -139      | kg/s                                  | Linear drag          |
| $Y_{v v }$       | -854.7687 | kg/m                                  | Nonlinear drag       |
| $Y_r$            | -43.8331  | kg.m <sup>2</sup> /rad.s <sup>2</sup> | Linear cross drag    |
| $Y_{r r }$       | -228.8418 | kg.m/rad <sup>2</sup>                 | Nonlinear cross drag |
| $N_v$            | 70.5207   | kg.m/s                                | Linear cross drag    |
| $N_{v v }$       | -143.9553 | kg                                    | Nonlinear cross drag |
| $N_r$            | -101.9671 | kg.m <sup>2</sup> /rad.s <sup>2</sup> | Linear drag          |
| $N_{r r }$       | -884.0436 | kg.m <sup>2</sup> /rad <sup>2</sup>   | Nonlinear drag       |
| $b$              | -10       | kg                                    | Offset in sway       |

form:

$$\begin{aligned}\dot{\mathbf{x}} &= \mathbf{f}(\mathbf{x}, \mathbf{u}), \\ \mathbf{y} &= \mathbf{g}(\mathbf{x}, \mathbf{u}),\end{aligned}\quad (7)$$

where  $\mathbf{x} = [v \ r \ y \ \psi]^T$  is the state space vector,  $\mathbf{y} = [r \ y \ \psi]^T$  is the output vector,  $\mathbf{u} = \delta_p$  is the control input, and  $\mathbf{f}$  and  $\mathbf{g}$  are nonlinear functions. More details of nonlinear and linearizing models can be found in [15].

#### IV. PARAMETER ESTIMATION APPROACH

The parameters estimation approach is presented here based on an optimization problem. The goal is to maximize a likelihood cost function, which means an expression of the output error or difference between measurement and output variables (Fig. 4). This is a nonlinear approach because the model to be identified presents nonlinearities (7). There should be included the parameter vector  $\Theta$ , the discrete time measurement variable  $\mathbf{z}(k)$  with  $N$  samples,  $\mathbf{w}$  state noise variable with its distribution matrix  $G$ , and  $\mathbf{v}$  measurement noise variable with distribution matrix  $F$ . Then, the nonlinear model to be used in MLE approach is rewritten as:

$$\begin{aligned}\dot{\mathbf{x}} &= \mathbf{f}(\mathbf{x}, \mathbf{u}, \Theta) + G(\mathbf{w}) \\ \mathbf{y} &= \mathbf{g}(\mathbf{x}, \mathbf{u}, \Theta) \\ \mathbf{z}(k) &= \mathbf{y}(k) + F\mathbf{v}(k)\end{aligned}\quad (8)$$

The prediction error is given by:

$$\mathbf{q}(k) = [\hat{\mathbf{y}}(k) - \mathbf{y}(k)]\quad (9)$$

Then, a likelihood function is expressed in function of  $\mathbf{q}$  and its respective covariance matrix  $\mathcal{B}$ :

$$p(\mathbf{y} | \Theta) = (2\pi)^{-m/2} |\mathcal{B}|^{-n/2} \cdot \exp \left[ -\frac{1}{2} \sum_{k=1}^n [\mathbf{q}(k, \Theta)]^T \mathcal{B}^{-1} [\mathbf{q}(k, \Theta)] \right]\quad (10)$$

where  $n$  is the dimension of state space vector  $\mathbf{x}$  and  $m$  is the dimension of the measurement vector  $\mathbf{y}$ . The value of  $\Theta$  is estimated through the maximization of this likelihood function, as follows:

$$\hat{\Theta} = \arg \max_{\Theta} p(\mathbf{y} | \Theta)\quad (11)$$

The likelihood expression (10) can be transformed using the relation  $-\ln(p(\mathbf{y} | \Theta))$ , and neglected the constant term [18], [7]:

$$J(\Theta) = \frac{1}{2} \sum_{i=1}^N \{ [\mathbf{q}(k, \Theta)]^T \mathcal{B}^{-1} [\mathbf{q}(k, \Theta)] + \ln |\mathcal{B}| \}\quad (12)$$

Therefore, minimize the functional  $J(\Theta)$  is equivalent to maximize the likelihood function with a great advantageous for computational purposes. There are many types of algorithms that solve the problem of the optimization, such as Gauss-Newton (GN) and Levenberg-Marquardt (LM).

To achieve an optimal  $\Theta$ , the cost function  $J(\Theta)$  should be approximated to a parabolic function using the well known Taylor series [7]:

$$\begin{aligned}J(\Theta_0 + \Delta\Theta) &\cong J(\Theta_0) + \Delta\Theta^T \left. \frac{\partial J}{\partial \Theta} \right|_{\Theta=\Theta_0} \\ &+ \frac{1}{2} \Delta\Theta^T \left. \frac{\partial^2 J}{\partial \Theta \partial \Theta^T} \right|_{\Theta=\Theta_0} \Delta\Theta\end{aligned}\quad (13)$$

where  $\Theta_0$  is the nominal vector parameter. The optimization is obtained under the constraint:

$$\frac{\partial}{\partial \Theta} [J(\Theta_0 + \Delta\Theta)] = 0\quad (14)$$

Solving expression (13) with the constraint given in (14), the variation of the estimated parameter vector  $\Delta\hat{\Theta}$  is:

$$\Delta\hat{\Theta} = - \left[ \left. \frac{\partial^2 J}{\partial \Theta \partial \Theta^T} \right|_{\Theta=\Theta_0} \right]^{-1} \left. \frac{\partial J}{\partial \Theta} \right|_{\Theta=\Theta_0}\quad (15)$$

Let the Hessian matrix be a non-singular matrix, the estimated parameter vector can be expressed as follows:

$$\hat{\Theta} = \Theta_0 + \Delta\hat{\Theta}\quad (16)$$

Considering the approximation (13) for the cost function, in the next iteration, the process will be repeated assuming the estimated vector as a nominal parameter  $\Theta_0 = \hat{\Theta}$ . Therefore, the generalized recursive equation is expressed as:

$$\Theta_{j+1} = \Theta_j - [\nabla_{\Theta}^2 J(\Theta_j)]^{-1} \nabla_{\Theta}^T J(\Theta_j)\quad (17)$$

where  $\nabla$  is the gradient of  $J$  whose Hessian matrix is  $\nabla_{\Theta}^2 J(\Theta_j)$ . Compute of Hessian matrix demands a huge computational effort which can be avoided by the GN algorithm:

$$\nabla_{\Theta}^2 J(\Theta) = \sum_{k=1}^n [\nabla_{\Theta} \hat{\mathbf{y}}_k(\Theta)]^T [\mathcal{B}]^{-1} [\nabla_{\Theta} \hat{\mathbf{y}}_k(\Theta)]\quad (18)$$

where the terms of the second order derivatives are removed. The gradient of the estimated output,  $\nabla_{\Theta} \hat{\mathbf{y}}_k(\Theta)$ , is named sensitivity function.

The LM algorithm is an extension of the GN, whose principal idea consists in modify  $\nabla_{\Theta}^2 J(\Theta)$  by the expression  $\nabla_{\Theta}^2 J(\Theta) + \lambda I$  in the Eq. (17). The inversion of the matrix is not yielded in explicit manner and, now, it will be solved

by singular value decomposition (SVD) according to the expression:

$$[\nabla_{\Theta}^2 J(\Theta) + \lambda \mathbf{I}] \Delta \hat{\Theta} = \nabla_{\Theta}^T J(\Theta_j) \quad (19)$$

The above LM algorithm [19] solves the problem of singularity in the Hessian matrix. Additionally, the LM algorithm works like a GN algorithm for small values of  $\lambda$ . These optimization algorithms are already present in libraries of non-commercial and commercial softwares, such as Gnu-Octave and Matlab.

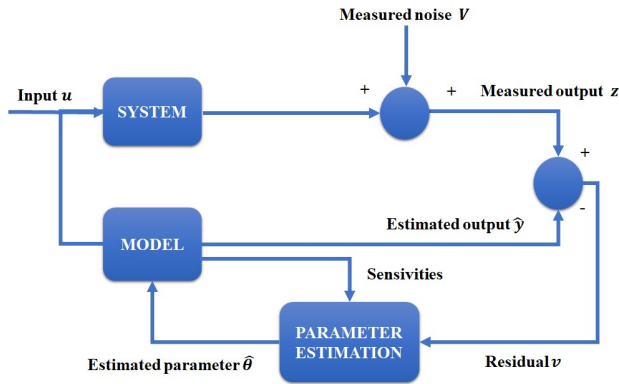


Fig. 4. Parameter Estimation via Output Error.

## V. RESULTS

This section gives the results of the approach applied to the ALBA ASC. For data generated in simulator and data obtained experimentally, the parameter vector defined in (8) is expressed relative to the nonlinear model as:

$$\Theta = [ \Theta_1 \quad \Theta_2 \quad \Theta_3 \quad \Theta_4 ]^T, \quad (20)$$

where  $\Theta_1 = Y_{\dot{v}}$ ,  $\Theta_2 = I_z$ ,  $\Theta_3 = N_v$ , and  $\Theta_4 = b$ . Three cases are analyzed in order to estimate a greater number of parameters.

### A. Simulated

The nonlinear model of the ALBA ASC expressed by (3)-(6) is implemented in Matlab/Simulink as shown in Fig. 5. The upper block named ALBA USV (unmanned surface vehicle) presents the nonlinear dynamics of the vehicle, and the lower block named Zig-zag maneuvering presents the zigzag maneuver generated numerically.

Fig. 6 presents the plot of data generated using this software resource, where  $\psi$  is the yaw angle in zigzag course due to the switching control of the propeller angle  $\delta_p$  between  $\pm 20^\circ = \pm 0.3490$  rad. This switching control signal is in closed loop and is activated by the yaw angle limits given also between  $\pm 20^\circ = \pm 0.3490$  rad. The yaw rate angle  $r$  follows an oscillatory behavior indicating the necessary angular rate of the vehicle to approach this zigzag maneuver, the initial condition for this numerical test is  $v = 0.001m/s$ ,  $r = 0.001rad/s$ ,  $y = 0.001m$  and  $\psi = -30 \times \pi/180rad$ . The simulated data were obtained with different hydrodynamic

parameters moved purposefully to 25% respect to the theoretical values (nominal values) given in Table II.

Table III presents the case where two parameters are estimated using the MLE algorithm. The fitting between identified model response and simulated data is 94.56% for the yaw rate  $r$ , 57.71% for the  $y$  position, and 58.88% for the yaw angle  $\psi$  (Table IX). Table IV presents the case where three parameters are estimated using the MLE algorithm. The fitting between identified model response and simulated data is 94.55% for the yaw rate  $r$ , 56.77% for the  $y$  position, and 58.98% for the yaw angle  $\psi$ . (Table IX). Table V presents the case where four parameters are estimated using the MLE algorithm. The fitting between identified model response and simulated data is 90.49% for the yaw rate  $r$ , 71.38% for the  $y$  position, and 83.16% for the yaw angle  $\psi$  (Table IX). Fig. 7 presents this last comparison validating the approach used here.

TABLE III. ESTIMATION WITH 2 PARAMETERS.

| Parameter  | Estimated value | Nominal value | Maneuver              |
|------------|-----------------|---------------|-----------------------|
| $\Theta_3$ | 102.551309      | 70.520700     | $\pm 20(\pi/180)$ rad |
| $\Theta_4$ | -8.139283       | -10.00000     | $\pm 20(\pi/180)$ rad |

TABLE IV. ESTIMATION WITH 3 PARAMETERS.

| Parameter  | Estimated value | Nominal value | Maneuver              |
|------------|-----------------|---------------|-----------------------|
| $\Theta_2$ | 172.83619       | 123.844500    | $\pm 20(\pi/180)$ rad |
| $\Theta_3$ | 102.722541      | 70.520700     | $\pm 20(\pi/180)$ rad |
| $\Theta_4$ | -8.3401932      | -10.00000     | $\pm 20(\pi/180)$ rad |

TABLE V. ESTIMATION WITH 4 PARAMETERS.

| Parameter  | Estimated value | Nominal value | Maneuver              |
|------------|-----------------|---------------|-----------------------|
| $\Theta_1$ | -897.157150     | -571.061900   | $\pm 20(\pi/180)$ rad |
| $\Theta_2$ | 188.315315      | 123.844500    | $\pm 20(\pi/180)$ rad |
| $\Theta_3$ | 103.501259      | 70.520700     | $\pm 20(\pi/180)$ rad |
| $\Theta_4$ | -8.923712       | -10.00000     | $\pm 20(\pi/180)$ rad |

### B. Experimental

The experimental tests were carried out in the Tingo lagoon in Arequipa. The inertial navigation system provided the yaw rate  $r$  and yaw angle  $\psi$  at sampling time of 0.1 s, a global positioning system (GPS) provided the position  $y$  at sampling time of 1 s. There was a need to use oversampling technique from 1 s to 0.1 s for the  $y$  data in order to feed the estimation algorithm used here. The control  $\delta_p$  signal is also provided at 0.1 s sampling time from the embedded electronic of ALBA ASC, described in [15].

Fig. 8 presents the plot of experimental data, where  $\psi$  is the yaw angle in zigzag course due to the switching control of the propeller angle  $\delta_p$  between  $\pm 20^\circ = \pm 0.3490$  rad. This switching control signal is in closed loop and is activated by the yaw angle limits given also between  $\pm 20^\circ = \pm 0.3490$  rad, the initial condition for this numerical test is  $v = 0m/s$ ,  $r = 0rad/s$ ,  $y = 0.01m$  and  $\psi = -30 \times \pi/180rad$

. The yaw rate angle  $r$  follows an oscillatory behavior indicating the necessary angular rate of the vehicle to approach this zigzag maneuver.

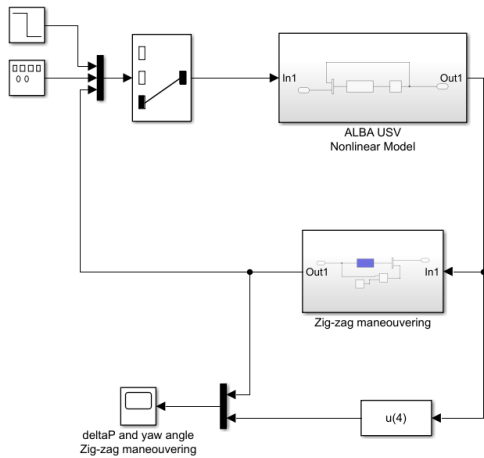


Fig. 5. Nonlinear Model Simulink.

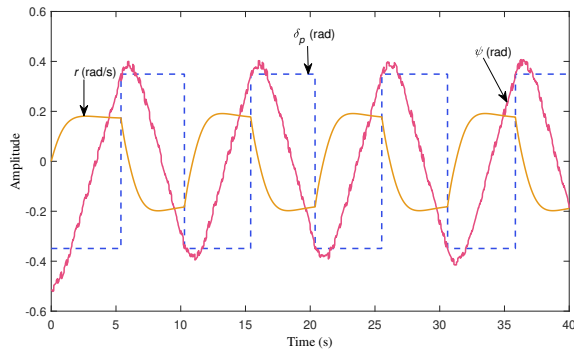


Fig. 6. Simulated Zigzag Maneuver.

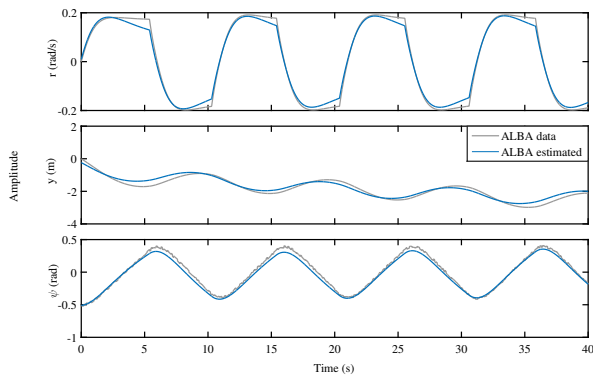


Fig. 7. Comparison Responses between the Identified Model and the Simulated Data.

Table VI presents the case where two parameters are estimated using the MLE algorithm. The fitting between identified model response and experimental data is 76.05% for the yaw rate  $r$ , 16.27% for the  $y$  position, and 63.35% for the yaw angle  $\psi$  (Table IX). Table VII presents the case where three parameters are estimated using the MLE algorithm. The fitting between identified model response and experimental data is 75.78% for the yaw rate  $r$ , 15.15% for the  $y$  position, and 63.45% for the yaw angle  $\psi$  (Table IX). Table VIII presents the case where four parameters are estimated using the MLE algorithm. The fitting between identified model response and simulated data is 75.71% for the yaw rate  $r$ , 15.31% for the  $y$  position, and 64.46% for the yaw angle  $\psi$  (Table IX). Fig. 9 presents this comparison validating the approach used here.

Table IX summarizes the numerical and experimental results carried out to estimate the main parameters of the ALBA. The fitting between the identified model and experimental data are good and above to 50%, except for the  $y$  position. As noted in model structure (3)-(6), its differential equation does not exert significant contribution in the whole model. Moreover, in autonomous vehicles [16], the  $y$  kinematic commonly compromises the observability and controllability linear properties. Here, there is an unsolved and open area for autonomous surface craft consisting in to examine identifiability properties and advances recently developed for biologic systems [13], [14].

## VI. CONCLUSIONS

A nonlinear model for the ALBA autonomous surface vehicle was identified using the maximum likelihood estimation approach. This approach was initially tested numerically with data obtained through the vehicle dynamic simulator. The approach was then applied to the vehicle data obtained in experimental test maneuvers. The four estimated parameters compose the identified system for the ALBA, a low-cost vehicle destined to monitor water conditions of lagoons and shallow water of the Pacific sea. The fitting between identified model responses and data is quite good and above 50%, guarantying the proposed approach and its application to this class of maritime vehicles. The fitting of  $y$  position was not good and analysis using identifiability properties should be conducted further.

TABLE VI. ESTIMATION WITH 2 PARAMETERS.

| Parameter  | Estimated value | Nominal value | Maneuver              |
|------------|-----------------|---------------|-----------------------|
| $\Theta_3$ | 68.478772       | 70.520700     | $\pm 20(\pi/180)$ rad |
| $\Theta_4$ | -7.552532       | -10.000000    | $\pm 20(\pi/180)$ rad |

TABLE VII. ESTIMATION WITH 3 PARAMETERS.

| Parameter  | Estimated value | Nominal value | Maneuver              |
|------------|-----------------|---------------|-----------------------|
| $\Theta_2$ | 140.339921      | 123.844500    | $\pm 20(\pi/180)$ rad |
| $\Theta_3$ | 70.648283       | 70.520700     | $\pm 20(\pi/180)$ rad |
| $\Theta_4$ | -7.874004       | -10.000000    | $\pm 20(\pi/180)$ rad |

TABLE VIII. ESTIMATION WITH 4 PARAMETERS.

| Parameter  | Estimated value | Nominal value | Maneuver              |
|------------|-----------------|---------------|-----------------------|
| $\Theta_1$ | -819.501860     | -571.061900   | $\pm 20(\pi/180)$ rad |
| $\Theta_2$ | 124.193961      | 123.844500    | $\pm 20(\pi/180)$ rad |
| $\Theta_3$ | 73.561674       | 70.520700     | $\pm 20(\pi/180)$ rad |
| $\Theta_4$ | -7.160477       | -10.000000    | $\pm 20(\pi/180)$ rad |

TABLE IX. FITTING BETWEEN IDENTIFIED MODEL RESPONSE AND DATA (SIMULATED AND EXPERIMENTAL).

| Fitting  | Number of estimated parameters |        |        |              |        |        |
|----------|--------------------------------|--------|--------|--------------|--------|--------|
|          | Simulated                      |        |        | Experimental |        |        |
| Response | two                            | three  | four   | two          | three  | four   |
| $r$      | 94.56%                         | 94.55% | 90.49% | 76.05%       | 75.78% | 75.71% |
| $y$      | 57.71%                         | 56.77% | 71.38% | 16.27%       | 15.15% | 15.31% |
| $\psi$   | 58.88%                         | 58.98% | 83.16% | 63.35%       | 63.45% | 64.46% |

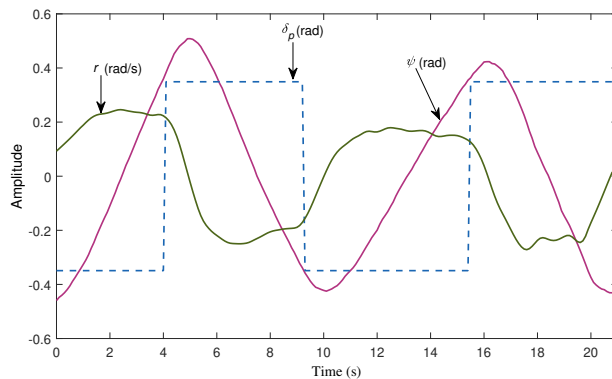


Fig. 8. Experimental Zigzag Maneuver.

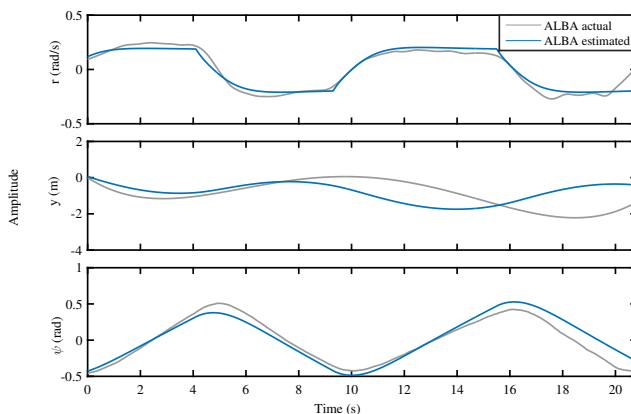


Fig. 9. Comparison Responses between the Identified Model and the Experimental Data.

ACKNOWLEDGMENT

The authors thank the *Universidad Nacional de San Agustín de Arequipa* for the incentives and financial support given to the construction of the ALBA vehicle, under grant number TP3-2018-UNSA.

REFERENCES

- [1] V. Rajendran, S. Nirmaladevi D, B. Srinivasan, C. Rengaraj, and S. Mariyaselvam, "Quality assessment of pollution indicators in marine water at critical locations of the gulf of mannar biosphere reserve, tuticorin," *Marine Pollution Bulletin*, vol. 126, pp. 236 – 240, 2018. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0025326X17309372>
- [2] M. Caccia, M. Bibuli, G. Bruzzone, R. Bono, and E. Spirandelli, "Charlie, a testbed for usv research," *IFAC Proceedings Volumes*, vol. 42, no. 18, pp. 97 – 102, 2009, 8th IFAC Conference on Manoeuvring and Control of Marine Craft. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S1474667016318778>
- [3] M. Dunbabin and A. Grinham, "Experimental evaluation of an autonomous surface vehicle for water quality and greenhouse gas emission monitoring," 06 2010, pp. 5268 – 5274.
- [4] C. R. Sonnenburg and C. A. Woolsey, "Modeling, identification, and control of an unmanned surface vehicle," *Journal of Field Robotics*, vol. 30, no. 3, pp. 371–398, 2013. [Online]. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1002/rob.21452>
- [5] D. Moreno, E. Besada, J. López, D. Chaos, J. Aranda, and J. Cruz, "Identificación de un modelo no lineal de un vehículo marino de superficie usando regresión simbólica," *Actas de las Jornadas de Automática*, pp. 850–855, 2015.
- [6] S. Lack, E. Rentzow, and T. Jeansch, "Experimental parameter identification for an open-frame rov: Comparison of towing tank tests and open water self-propelled tests," *IFAC-PapersOnLine*, vol. 52, no. 21, pp. 271 – 276, 2019, 12th IFAC Conference on Control Applications in Marine Systems, Robotics, and Vehicles CAMS 2019. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S2405896319322049>
- [7] E. Morelli and V. Klein, *Aircraft System Identification: Theory and Practice*. Sunflyte Enterprises, 2016. [Online]. Available: <https://books.google.com.pe/books?id=HsW-AQAACA AJ>
- [8] J. C. Cutipa-Luque and D. C. Donha, "Auv identification and robust control," *IFAC Proceedings Volumes*, vol. 44, no. 1, pp. 14735 – 14741, 2011, 18th IFAC World Congress. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S1474667016459963>
- [9] S. Randeni, A. Forrest, R. Cossu, Z. Leong, D. Ranmuthugala, and V. Schmidt, "Parameter identification of a nonlinear model: replicating the motion response of an autonomous underwater vehicle for dynamic environments," *Nonlinear Dynamics*, 11 2017.
- [10] C. Jian, Z. Jiayuan, X. Feng, Y. Jianchuan, Z. Zaojian, Y. Hao, X. Tao, and Y. Luchun, "Parametric estimation of ship maneuvering motion with integral sample structure for identification," *Applied Ocean Research*, vol. 52, pp. 212 – 221, 2015. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0141118715000814>
- [11] M. Saccomani and L. D'Angiò, "Examples of testing global identifiability with the daisy software," *IFAC Proceedings Volumes*, vol. 42, no. 10, pp. 48 – 53, 2009, 15th IFAC Symposium on System Identification. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S1474667016386219>
- [12] N. Meshkat, C. E.-z. Kuo, and J. DiStefano III, "On finding and using identifiable parameter combinations in nonlinear dynamic systems biology models and combos: A novel web implementation," *PLoS One*, vol. 9, no. 10, p. e110261, 2014.
- [13] A. F. Villaverde, A. Barreiro, and A. Papachristodoulou, "Strike-gold user manual structural identifiability taken as extended-generalized observability with lie derivatives and decomposition," 2016.
- [14] T. S. Ligon, F. Fröhlich, O. T. Chiş, J. R. Banga, E. Balsa-Canto, and J. Hasenauer, "GenSSI 2.0: multi-experiment structural identifiability analysis of SBML models," *Bioinformatics*, vol. 34, no. 8, pp. 1421–1423, 11 2017. [Online]. Available: <https://doi.org/10.1093/bioinformatics/btx735>
- [15] B. A. Monroy-Ochoa and J. C. Cutipa-Luque, "Development of a low-cost unmanned surface vehicle for water quality monitoring," *International Journal of Control and Automation*, vol. 13, no. 4, pp. 1197–1207, Jul. 2020.

- [16] T. I. Fossen, *Handbook of marine craft hydrodynamics and motion control*. John Wiley & Sons, 2011.
- [17] J. N. Newman, *Marine hydrodynamics*. The MIT press (40th anniversary edition), 2018.
- [18] L. Ljung, *System Identification: Theory for the user*, 2nd ed. Prentice Hall, 1999.
- [19] B. Maciel, L. Goes, E. Hemerly, and N. Neto, "Flight path reconstruction and parameter estimation using output-error method," *Journal of Shock and Vibration*, vol. 13, pp. 379–392, 2006.



# Mobility-Aware Container Migration in Cloudlet-Enabled IoT Systems using Integrated Muticriteria Decision Making

Mutaz A. B. Al-Tarawneh

Computer Engineering Department, Mutah University  
Karak, JORDAN 61710

**Abstract**—Service migration plays a vital role in continuous service delivery in Internet of Things (IoT) systems. This paper presents a mobility-aware container migration algorithm for Cloudlet-enabled IoT systems. The proposed algorithm is based on an integrated multicriteria decision making (MCDM) approach. It has been implemented using a specialized simulation tool and compared to other existing migration algorithms. Simulation results demonstrate the ability of the proposed algorithm to achieve up to 48%, 48%, 20% and 36% improvement in migration time, service downtime, migration reliability and service loss rate, respectively as compared to other migration algorithms. The proposed algorithm is capable of perceiving the run-time dynamics of IoT systems and appropriately manage the process of container migration.

**Keywords**—Internet of Things (IoT); container; migration; Cloudlet; criteria; decision making

## I. INTRODUCTION

Nowadays, the tremendous growth of commodity-available smart objects has resulted in the ubiquity of Internet of Things (IoT) applications. These objects can be either fixed or mobile and generate streams of big data with immense processing and memory requirements that surpass the capabilities of user devices. Therefore, IoT applications urged the use of Cloud Computing as a processing back-end that fulfills the processing requirements of such applications [1].

Following the expansion in IoT applications and the elevation of their processing requirements and the urgency of their latency constraints, the inseparable relation between the IoT applications and the Cloud has been hindered due to the quasi-central nature of the Cloud-based services and resources [2]. In order to cope with the processing and timing requirements of IoT applications and overcome the limitations of the Cloud-based IoT platforms, several edge computing models [3] such as Multi-access Edge Computing [4], Fog Computing [5] and Cloudlet Computing [6] have been proposed. For instance, Cloudlet Computing is considered as a middle-way between Fog Computing and Cloud Computing models. It places computing clusters with sufficient network bandwidth and processing capabilities in proximity to the IoT devices. A Cloudlet is basically a small data center, which is usually instantiated at few hops away from IoT devices and employ virtualization at either virtual machine or container levels [7].

According to the Cloudlet Computing paradigm, IoT application modules and services are placed on close tiny data centers aiming at satisfying the stringent timing requirements

of the requests generated from IoT devices. Computing requests - generated from IoT devices are offloaded to such nearby tiny data centers which in turn perform the required functionality as dictated by the received requests. Once the received request is executed, execution results are sent back to the source IoT device [7]. In order to efficiently utilize the Fog/Cloudlet resources and minimize any potential performance overheads, a lightweight virtualization technique is employed, in which application modules and user data are encapsulated in containers [8, 9].

While Cloudlet Computing provides Cloud-like services in proximity to IoT devices and offers low-latency services to these devices, its advantages could be diminished in case of mobile or non-stationary IoT devices, where the access point that connects the IoT device to the Cloudlet may change [10]. In such scenario, user or IoT device mobility will increase the number of hops between the IoT device the Cloudlet hosting the associated container and, in turn, negatively impact the latency requirements of IoT devices. In other words, as the number of hops between the IoT device and its associated Cloudlet increases, the latency of IoT service requests will increase in a manner that resembles that of Cloud-based platforms, which shrinks the expected performance of Cloudlet-based platforms. Therefore, as the mobile IoT devices move away from their associated Cloudlets, both processing and control of their application modules should be handed over to the Cloudlet in proximity to the access point to which the IoT device is currently connected. This process can be performed by migrating the container, associated with the mobile IoT device, to a new Cloudlet that is placed near to the IoT device and respects that device's functional and non-functional requirements [11].

Apparently, IoT devices requesting services from nearby Cloudlets may be mobile and require continuity of service delivery across different locations. However, maintaining service delivery for mobile users and IoT devices is a challenging process. The study in [12] has proposed a framework for a Fog Computing architecture in which virtual machine migration is supported. They have assumed that each user has a virtual machine running on a particular Cloudlet. This Cloudlet is considered as an endpoint that provides services to the associated user. A virtual machine could be migrated to another Cloudlet based on user location, direction of movement, running applications, computing capacity of the destination Cloudlet and the network capacity. However, the proposed framework has not been evaluated and its suitability to existing Fog

environments has not been verified. On the other hand, [13] has proposed a genetic algorithm-based model for virtual machine migration in Mobile Cloud Computing (MCC) environments considering both user mobility and the current workload of the potential destination Cloudlet in order to reduce the number of virtual machine migrations. In addition, [14] has proposed a simulation framework based on the iFogSim tool [15] to support user mobility by migrating virtual machines across Fog nodes. Furthermore, [16] has proposed a Fog computing architecture that supports user mobility by implementing a route optimization algorithm that improves the performance of the handover mechanisms. Similarly, [17] has proposed a framework for Mobile Edge Computing (MEC) environments in which service migration is implemented. Their goal was to ensure service continuity for mobile users and reduce both service downtime and overall migration time. Although previous research efforts have tackled service migration in various edge computing environments, their efforts have either supported migration at the virtual machine level or ignored important factors such as user location and mobility patterns.

More recently, the work in [18] has proposed an autonomic container migration approach for Fog computing environments. Their approach seeks to ensure continuous service delivery by migrating application modules to Fog nodes that are deemed to respect service delivery deadlines. The efficacy of their approach has been quantified in terms of network usage, execution cost and service execution delays. However, they did not consider other important factors such as service downtime, migration time, service loss and the reliability of the migration process. On the other hand, the work in [19] is the most recent effort that models container migration in Fog/Cloudlet computing environments. It models container-based migration and supports realistic user mobility patterns based on the SUMO urban simulation tool [20, 21]. In addition, they have proposed several baseline container migration algorithms that are categorized as: lowest latency (LL) in which containers are migrated to the Cloudlet that is expected to deliver the lowest service latency, closest access point (CAP) where containers are migrated to the Cloudlet connected to the closest access point to current user location and closest server Cloudlet (CSC) in which containers are migrated to the closest server Cloudlet to current user location. However, none of these migration algorithms has considered migration reliability, service downtime and overall migration time when making the migration decision.

This work proposes a container migration algorithm for Cloudlet computing environments based on an integrated multicriteria decision making (MCDM) approach that integrates the Entropy [22] and Technique of Order Preference Similarity to the Ideal Solution (TOPSIS) [23] methods. The proposed approach seeks to migrate a container associated with a particular mobile IoT device to a nearby server Cloudlet considering user location, mobility pattern, migration time and migration reliability. While MCDM has been used in related disciplines such as cloud service selection [24, 25] and task scheduling in Cloud Computing [26, 27], no prior attempts have been made to apply MCDM for container migration in Cloudlet computing environments. The main contributions of this work can be summarized as follows:

- Proposing an MCDM-based algorithm for container mi-

gration in Cloudlet-enabled IoT systems.

- Implementing the proposed algorithm using a specialized simulation tool with realistic migration models and user mobility patterns.
- Performing a series of experiments to assess the performance of the proposed migration algorithm and compare it against that of existing algorithms.

The rest of this paper is organized as follows. Section II explains the research methodology. Section III presents and discusses the obtained results and Section IV summarizes and concludes this paper.

## II. RESEARCH METHODOLOGY

This section explains the research methodology followed in order to implement the proposed migration algorithm including the system model, mathematical formulation of the MCDM techniques and the algorithmic design of the proposed algorithm.

### A. System Model

In this paper, a hierarchical Cloudlet-enabled IoT system is assumed as shown in Fig. 1. Service requests generated from IoT devices are sent to the nearby Cloudlet which in turn performs the required task and communicate the results back to the IoT devices. If a task required by a particular IoT device can not be performed by the Cloudlet or requires further processing, the Cloudlet will forward the received request to the Cloud. Each Cloudlet supports a lightweight virtualization technique in which user's data and application modules are encapsulated in containers.

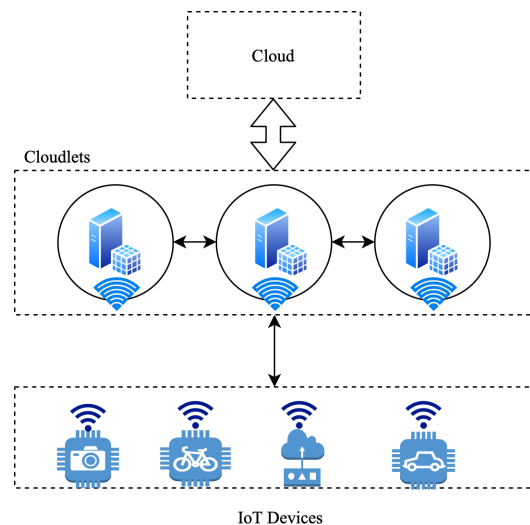


Fig. 1. IoT System Model.

On the other hand, IoT devices are assumed to be mobile i.e. change their locations over time. At any time instant, each IoT device will be connected to the Cloudlet hosting its associated container and can access that Cloudlet's services through a particular access point. In order to prevent service interruption and guarantee low-latency service delivery, container migration (from a source to a destination Cloudlet) should be performed

in a timely manner. Determining the time instant at which the migration process should be initiated and choosing the most appropriate Cloudlet to host the migrating container are two major issues in container migration. In this work, each Cloudlet is considered as the decision maker responsible for addressing the two aforementioned issues. In other words, each Cloudlet is assumed to monitor the current location of its associated IoT devices and decide if a migration should be performed and what is the most suitable Cloudlet to host the migrating container considering both migration reliability and expected migration time.

The Cloudlets layer considered in this work consists of a set of Cloudlet nodes  $C = \{SC_1, SC_2, \dots, SC_N\}$ . The total number of cloudlets in the environment is denoted as  $N$ . In addition, each Cloudlet node is characterized by a triplet  $\langle Cn_i(cpu)^{cap}, Cn_i(mem)^{cap}, Cn_i(bw)^{cap} \rangle$  where  $Cn_i(cpu)^{cap}$  is the number of CPU cores available on Cloudlet  $i$ ,  $Cn_i(mem)^{cap}$  is the amount of memory available on node  $i$  and  $Cn_i(bw)^{cap}$  is the available bandwidth. The mobile IoT devices receive services from their associated containers which are deployed on various Cloudlet nodes available in the IoT environment. Each container ( $con_j$ ) is allocated to a suitable Cloudlet that satisfies its resource requirements. The resource requirements of a particular container can be denoted as a triplet  $\langle Con_j(cpu)^{req}, Con_j(mem)^{req}, Con_j(bw)^{req} \rangle$  where  $Con_j(cpu)^{req}$  is the number of CPU cores required by container  $j$ ,  $Con_j(mem)^{req}$  is the required memory and  $Con_j(bw)^{req}$  is the required bandwidth.

On the other hand, Fig. 2 shows the migration model assumed in this work. This migration model is based on [14, 19]. This model defines both the migration zone and the migration point on the map. Whereas the migration zone defines the area in which migration decisions are always computed and evaluated, the migration point represents the point on the map at which container migration should be initiated. The decision on whether to migrate a container associated with a particular user (or IoT device) should be made by the Cloudlet to which that IoT device is currently connected. For instance, Fig. 2 shows two users (user-1 and user-2) moving on the map with different speeds, directions and geographical positions. Their attributes should be continuously monitored by their current Cloudlet that is responsible for making the right decision for each user. As shown, user-1 is moving towards the access point (AP) that connects it to its current Cloudlet while user-2 is moving away from the AP. Hence, no migration should be performed for the container associated with user-1 while a migration for the container associated with user-2 should be initiated once it reaches the migration point. The migration point can be set either statically based on the coverage area of the AP or dynamically based on some dynamic attributes such as the user speed. For instance, Fig. 2 shows an example of a fixed migration point that is set at a distance that is 70% of the radius of the coverage area of the AP.

Once a Cloudlet decides to migrate the container associated with a particular IoT device, it must consult its migration algorithm to select the most suitable destination Cloudlet from a set of potential Cloudlets confined within the migration cone. The migration cone is defined in terms of the two adjacent directions to the current direction of the user and an angle that defines the relative region between the user and the AP. For

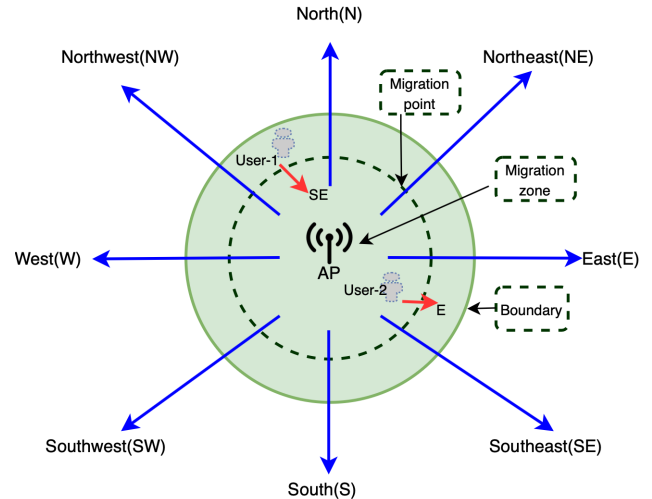


Fig. 2. Migration Model [19].

example, the cone associated with user-2 is confined between the Northeast and the Southeast directions. The value of the cone's angle could be a fixed value or a dynamic value that is changed at run-time subject to some optimization constraints.

### B. MCDM-Based Migration Algorithm

This section explains the proposed MCDM-based migration algorithm. Fig. 3 highlights the main steps followed by the proposed migration algorithm in order to select the most appropriate destination Cloudlet for a migrating container. As mentioned in Section II-A, each Cloudlet continuously monitors the movement of its associated IoT devices. Once a Cloudlet observes that an associated mobile IoT device is moving away from its associated AP, it must prepare and initiate the migration process by performing a number of steps based on an integrated Entropy-TOPSIS multicriteria decision making approach. First, the migration cone defined by the migration model must be constructed. Second, if the migration cone is not empty (i.e. it contains at least two possible destination Cloudlets that can satisfy all processing, memory and bandwidth requirements of the migrating container), the source Cloudlet will construct a decision matrix  $D_{ij}$  based on two criteria, namely, migration time ( $mt_{ij}$ ) and migration reliability ( $mr_{ij}$ ).

Apparently, migration time represents the expected time required to complete transferring the migrating container from source Cloudlet  $i$  to destination Cloudlet  $j$ . The total time required to migrate a container from Cloudlet  $i$  to Cloudlet  $j$  is computed as shown in equation 1.

$$mt_{ij} = \frac{M_{con_i}}{bw_{ij}} + T_{U_i} + T_{D_{ij}} \quad (1)$$

Where  $mt_{ij}$  is the total migration time,  $M_{con_i}$  is the memory size of the container migrating from Cloudlet  $i$ ,  $bw_{ij}$  is the available network bandwidth between Cloudlets  $i$  and  $j$ ,  $T_{U_i}$  is the uplink latency of Cloudlet  $i$  and  $T_{D_{ij}}$  is the latency by distance between Cloudlets  $i$  and  $j$ .

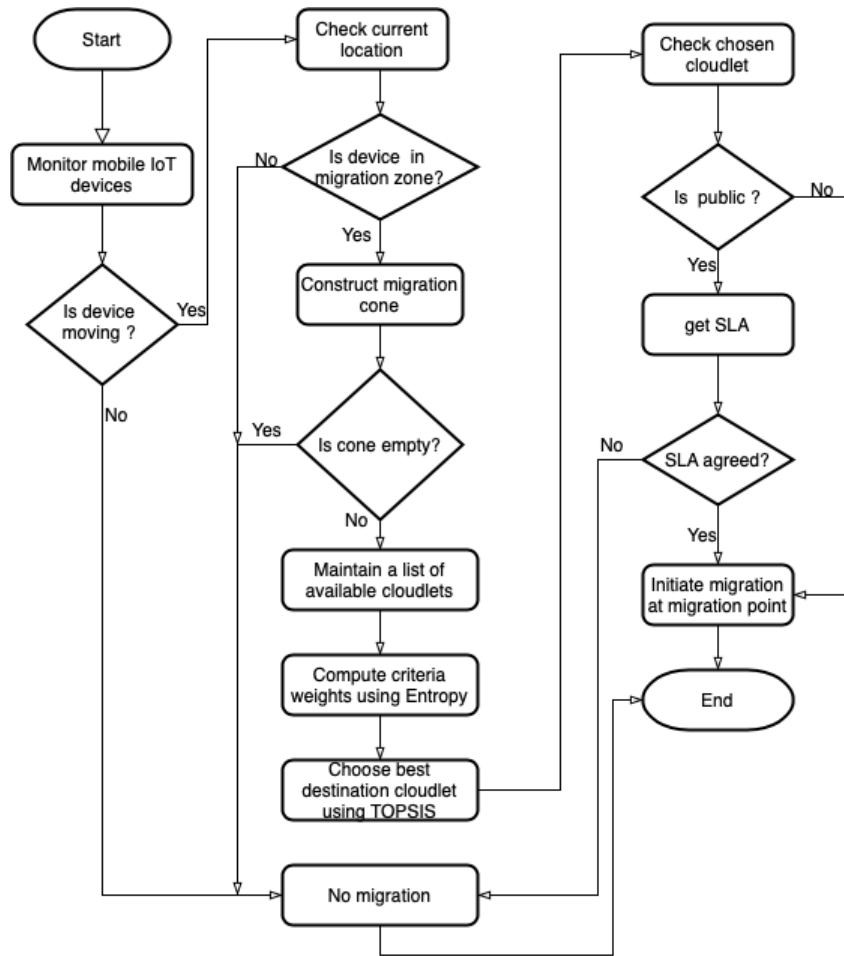


Fig. 3. Migration Algorithm Flow-chart.

On the other hand, migration reliability ( $mr_{ij}$ ) quantifies the probability that the migration process (from Cloudlet  $i$  to Cloudlet  $j$ ) will be completed successfully during the anticipated migration time. In this work, migration reliability is calculated based on the reliability of the communication link between Cloudlets  $i$  and  $j$  and the reliability of the destination Cloudlet  $j$ . In general, the reliability of a particular subsystem  $s$  can be computed as shown in equation 2 [28].

$$R_s(t) = e^{-\int_0^t h(\tau) d\tau} \quad (2)$$

Where  $h(t)$  is the hazard rate function associated with subsystem  $s$ . Assuming that both the communication link, between Cloudlets  $i$  and  $j$ , and destination Cloudlet  $j$  have constant hazard functions, their reliability values can be computed as shown in equation 3.

$$\begin{aligned} R_{link(i,j)}(t) &= e^{-\lambda_{ij}t} \\ R_{C_j}(t) &= e^{-\lambda_{C_j}t} \end{aligned} \quad (3)$$

Where  $R_{link(i,j)}$  is the reliability of the communication link between Cloudlets  $i$  and  $j$ ,  $R_{C_j}$  is the reliability of

Cloudlet  $j$ ,  $\lambda_{ij}$  is the constant hazard rate of the communication link and  $\lambda_{C_j}$  is the constant hazard rate of Cloudlet  $j$ . Since the communication link and the destination Cloudlet constitute a serial reliability model, the reliability of the migration process (from Cloudlet  $i$  to Cloudlet  $j$ ) can be computed by multiplying these individual reliability values [29], as shown in equation 4.

$$mr_{ij}(t) = R_{link(i,j)}(t) * R_{C_j}(t) = e^{-(\lambda_{ij} + \lambda_{C_j})t} \quad (4)$$

Where  $mr_{ij}$  is the reliability of migrating a container from Cloudlet  $i$  to Cloudlet  $j$ . Intuitively, by setting the value of  $t$  to  $mt_{ij}$ , the reliability of a migration process whose expected finish time is  $mt_{ij}$  can be computed. Having obtained the migration time and migration reliability of each possible destination Cloudlet  $j$ , the source Cloudlet will construct the decision matrix as shown in equation 5.

$$D_{ij} = \begin{bmatrix} mt_{i1} & mr_{i1} \\ mt_{i2} & mr_{i2} \\ mt_{i3} & mr_{i3} \\ \vdots & \vdots \\ mt_{in_c} & mr_{in_c} \end{bmatrix} \quad (5)$$

Where  $n_c$  is the number of Cloudlets available within the constructed migration cone. Third, once the decision matrix is obtained, the proposed migration algorithm will employ an Entropy-based method to assign a weight to each decision criterion as weight calculation is an essential step towards MCDM-based alternative (Cloudlet) selection approaches. The computed weight signifies the importance given to each criterion when making the migration decision. The computation of criteria weights based on the Entropy method can be summarized as follows:

**Step 1:** Computation of normalized feature weight ( $P_{ij}$ ) for the  $i^{th}$  alternative and the  $j^{th}$  criterion.

$$P_{ij} = \frac{d_{ij}}{\sum_{i=1}^{n_c} (d_{ij})}, (1 \leq i \leq n_c, 1 \leq j \leq m) \quad (6)$$

Where  $d_{ij}$  is the value of  $j^{th}$  criterion under the  $i^{th}$  alternative,  $n_c$  is the number of possible alternatives (i.e. Cloudlets) and  $m$  is the number of criteria.

**Step 2:** Calculation of the output entropy ( $e_j$ ) for each criterion  $j$ .

$$e_j = -K * \sum_{i=1}^{n_c} (P_{ij} * \ln(P_{ij})), (1 \leq j \leq m) \quad (7)$$

$$K = \frac{1}{\ln(n_c)}$$

**Step 3:** Calculation of the degree of diversification ( $g_j$ ) for each criterion  $j$ .

$$g_j = |1 - e_j|, 1 \leq j \leq m \quad (8)$$

**Step 4:** Computation of the weight ( $w_j$ ) of each criterion.

$$w_j = \frac{g_j}{\sum_{j=1}^m (g_j)}, (1 \leq j \leq m) \quad (9)$$

Fourth, having obtained the criteria weights using the Entropy method, the migration algorithm proceeds to select the most suitable destination Cloudlet to receive the migrating container using the TOPSIS method. The rationale behind the TOPSIS method is to choose the alternative with the shortest euclidean distance from the ideal solution and the farthest distance from the negative-ideal solution. The TOPSIS method employed to select the ideal alternative (Cloudlet) proceeds as follows:

**Step1:** The decision matrix ( $D_{ij}$ ) is normalized using vector normalization technique and then weighted using the criteria weights - computed using the Entropy method to obtain the weighted normalized performance matrix.

$$\bar{d}_{ij} = \frac{d_{ij}}{\sqrt{\sum_{i=1}^{n_c} (d_{ij})^2}}, (1 \leq i \leq n_c, 1 \leq j \leq m) \quad (10)$$

$$\bar{d}_{ij} = w_j * \bar{d}_{ij}, (1 \leq j \leq m)$$

**Step 2:** Determine the best condition ( $A_j^+$ ) and the worst condition ( $A_j^-$ ) with respect to each criterion ( $j$ ) and construct

the best conditions vector ( $A_b$ ) and the worst conditions vector ( $A_w$ ).

$$A_j^+ = \begin{cases} \max_i(\bar{d}_{ij}) & j \in J^+, (1 \leq i \leq n_c) \\ \min_i(\bar{d}_{ij}) & j \in J^-, (1 \leq i \leq n_c) \end{cases} \quad (11)$$

$$A_j^- = \begin{cases} \min_i(\bar{d}_{ij}) & j \in J^+, (1 \leq i \leq n_c) \\ \max_i(\bar{d}_{ij}) & j \in J^-, (1 \leq i \leq n_c) \end{cases} \quad (12)$$

Where  $J^+$  indicates beneficiary criteria i.e. criteria with positive impact (e.g. migration reliability) while  $J^-$  represents non-beneficiary criteria i.e. criteria with negative impact on the migration process (e.g. migration time). The best and worst conditions vectors are then constructed as:

$$A_b = \{A_j^+ \mid (1 \leq j \leq m)\} \quad (13)$$

$$A_w = \{A_j^- \mid (1 \leq j \leq m)\}$$

**Step 3** Calculate the Euclidean distance between each alternative  $i$  and each of the best condition vector ( $A_b$ ) and the worst condition vector ( $A_w$ ).

$$L_{ib} = \sqrt{\sum_{j=1}^m (\bar{d}_{ij} - A_b[j])^2}, (1 \leq i \leq n_c) \quad (14)$$

$$L_{iw} = \sqrt{\sum_{j=1}^m (\bar{d}_{ij} - A_w[j])^2}, (1 \leq i \leq n_c)$$

Where  $L_{ib}$  and  $L_{iw}$  are the distances from the target alternative (Cloudlet)  $i$  to the best and worst condition vectors, respectively.

**Step 4** Calculate the performance score ( $P_i$ ) assigned to each alternative.

$$P_i = \frac{L_{iw}}{L_{iw} + L_{ib}}, (1 \leq i \leq n_c) \quad (15)$$

**Step 5** Rank alternatives (Cloudlets) according to their performance scores ( $P_i$ ).

After performing the steps dictated by the TOPSIS method, the Cloudlet with the highest performance score is chosen as the most suitable destination Cloudlet to receive the migrating container and the migration process will be initiated upon device's arrival at the migration point. If the chosen Cloudlet is a public Cloudlet, its service level agreement (SLA) parameters will be compared against the requirements of the IoT device (user) whose container is involved in the migration process. If the chosen Cloudlet satisfies the requirements of the associated IoT device, the migration process will be initiated once that device reaches the predefined migration point. On the other hand, if the chosen Cloudlet is not guaranteed to fulfil the requirements of the associated IoT device, the migration process will be aborted and the associated container will remain on its original Cloudlet.

### III. RESULTS AND DISCUSSION

In order to evaluate the performance of the proposed migration algorithm, it has been implemented using the MobFogSim simulation tool [19]. MobFogSim is a recently developed tool that allows development and evaluation of container migration in Fog and Cloudlet computing environments with concrete support for realistic user mobility patterns. The proposed algorithm has been validated and compared against other migration algorithms that are originally implemented within the MobFogSim tool. Table I summarizes the most important input simulation parameters. Whereas some simulation parameters were given fixed values, other parameters such as network bandwidth and the constant hazard rates were randomly picked from the designated ranges. Hence, the results presented in this section represent the arithmetic mean of the results of 20 different simulation runs.

TABLE I. SIMULATION SETTINGS.

| Parameter                              | Value                         |
|--|-------------------------------|
| Size of container's execution state    | 12.8 MB                       |
| Number of Cloudlets                    | 144                           |
| Number of Cloudlets per access point   | 1                             |
| Radius of access point coverage( $r$ ) | 500 m                         |
| Migration point policy                 | fixed ( $r=40$ m)             |
| IoT device speed                       | 20 Km/h                       |
| Network bandwidth (between Cloudlets)  | [2,8] MB/s                    |
| Hazard rate                            | $[1 * 10^{-7}, 15 * 10^{-7}]$ |

Fig. 4 depicts the total migration time achieved under the proposed policy along with the migration time achieved under existing algorithms, namely, lowest latency (LL), closest access point (CAP) and closest server Cloudlet (CSC) [19]. The results are annotated with their 95% confidence interval. In order to reasonably assess the performance of different migration algorithms, they were simulated under different migration techniques: cold migration (Cold) and live migration (Live). In cold migration: first, the container to be migrated is frozen/stopped to make sure that it no longer modifies its state. Second, the whole execution state of that container is checkpointed and then transferred while being stopped. Third, the container is resumed at the destination once all its execution state is available there. On the other hand, live migration first halts container on the source Cloudlet and copies its minimal and kernel execution state to the destination Cloudlet so that the container can resume its execution there. Only after that and while the container is running, it transfers all the remaining state including the memory pages, which represent the major portion of the whole state. As shown in Fig. 4, the proposed algorithm has achieved the lowest migration time as compared to other migration algorithms under both the Cold and Live migration techniques; the employed Entropy-TOPSIS decision making approach has allowed the proposed algorithm to observe the degree of diversification available among the considered potential destination Cloudlets (i.e. those available within the migration cone) and choose the Cloudlet that would yield the lowest migration time as compared to other algorithms. In addition, the Cold migration technique has lower migration time when compared to the Live migration technique under all possible migration algorithms. On average, the proposed algorithm has achieved 48%, 42% and 43% improvement in the migration time as compared to the LL, CAP and CSC algorithms, respectively.

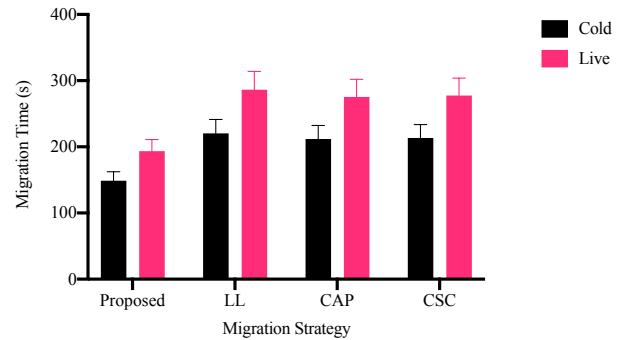


Fig. 4. Migration Time.

On the other hand, Fig. 5 illustrates the total downtime observed under the proposed and existing migration algorithms. The downtime is the time interval during which the service required by the mobile IoT devices will not be available due to the migration process. As shown, the proposed algorithm has resulted in lower downtime as compared to other algorithms under both migration techniques. It can be observed that the downtime is equal to the migration time when the Cold migration technique is employed. On the other hand, the downtime is much lower than the migration time when the Live migration technique is used; under the Live migration technique, the migrating container keeps on running while the majority of its state is being moved to the chosen destination Cloudlet. The container is stopped only for the transfer of a minimal amount of its overall state, after which the container runs at the destination Cloudlet. Nevertheless, the proposed algorithm can still achieve lower downtime as compared to other migration algorithms due to its ability to select the most appropriate destination Cloudlet with sufficient network bandwidth to handle the migration process. The amount of improvement the proposed algorithm has achieved in downtime - as compared to other algorithms is equivalent to that observed in migration time.

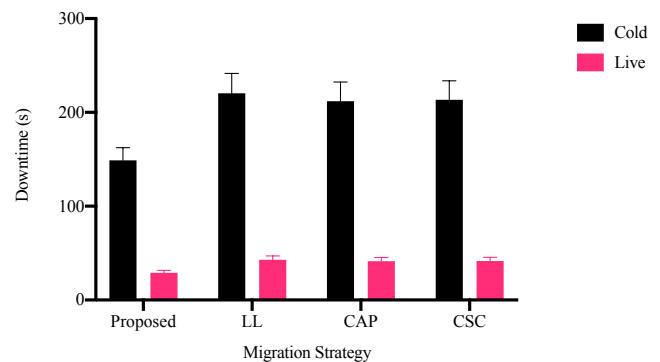


Fig. 5. Service Downtime.

Fig. 6 shows migration reliability obtained by the proposed and the existing migration algorithms under Cold and Live migration techniques. As shown, the proposed algorithm was able to surpass other migration algorithms in terms of mi-

gration reliability. Apparently, migration reliability decreases when moving from Cold to Live migration technique due to the increase in migration time - under the same hazard rate. This can be explained by the exponential relationship shown in equation 4. However, the proposed algorithm can still outperform other algorithms and pick the most suitable destination Cloudlet considering the reliability factor. Overall, the proposed algorithm has achieved 17%, 20% and 20% average improvement in migration reliability as compared to LL, CAP and CSC algorithms, respectively. On the other hand, Fig. 7 portrays the service loss rate under the considered migration algorithms. Service loss rate is the percent of service requests that could not be handled due to the migration process. This rate is clearly dependent on the downtime. As shown, the proposed algorithm was able to yield the lowest service loss rate as compared to other algorithms under both Cold and Live migration techniques. The proposed algorithm has achieved 36%, 27% and 28% reduction in service loss rate when compared to LL, CAP and CSC, respectively.

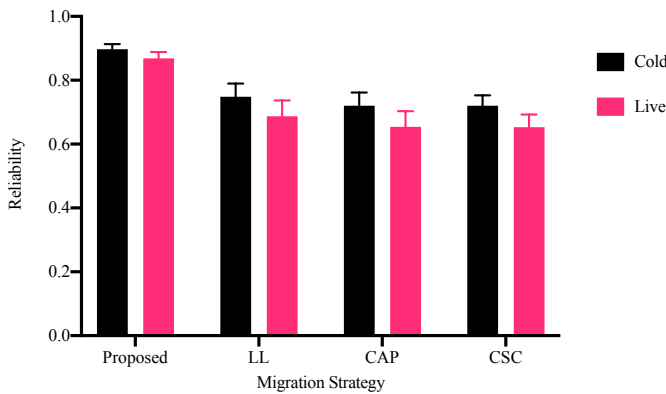


Fig. 6. Migration Reliability.

The presented results prove the ability of the proposed migration algorithm to perceive the run-time dynamics of the Cloudlet-enabled IoT environment and choose the most appropriate destination Cloudlet that would optimize the considered performance metrics. This is due to the fact that the proposed MCDM-based migration algorithm is able to observe the degree of variation - with respect to each migration criterion within the considered destination Cloudlets, reasonably assign criteria weights (using the Entropy method) and appropriately assign performance scores to the considered Cloudlets (using TOPSIS) when making a migration decision. This fact is confirmed by the results shown in Tables II and III which represent two sample migration scenarios. As shown, migration time has been assigned higher weight than migration reliability in the first scenario (Table II) while the opposite is observed in the second scenario (Table III).

On the other hand, Fig. 8 illustrates how criteria weights were assigned during the simulation process for 12 consecutive migration events. These results prove the ability of the proposed migration algorithm to sense the dynamic changes in the decision variables (criteria) which, in turn, steers the decisions making process i.e. destination Cloudlet selection

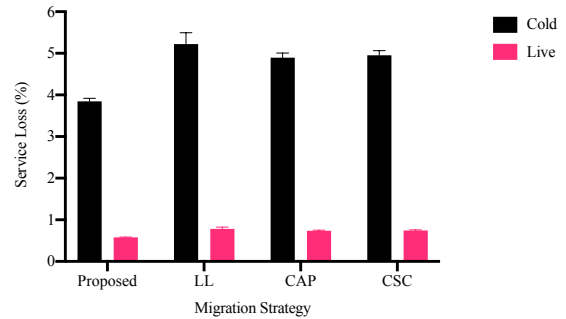


Fig. 7. Service Loss Rate.

TABLE II. MIGRATION SCENARIO 1.

| Alternative ID  | Migration Time (s) | Migration Reliability | Performance Score |
|-----------------|--------------------|-----------------------|-------------------|
| A               | 128.6              | 0.85                  | 0.97              |
| B               | 128.7              | 0.82                  | 0.95              |
| C               | 128.7              | 0.73                  | 0.91              |
| D               | 144.4              | 0.80                  | 0.82              |
| E               | 163.6              | 0.91                  | 0.62              |
| F               | 167.7              | 0.72                  | 0.56              |
| G               | 219.5              | 0.89                  | 0.13              |
| H               | 220.8              | 0.85                  | 0.11              |
| I               | 219.2              | 0.62                  | 0.02              |
| Criteria Weight | 0.81               | 0.19                  |                   |

TABLE III. MIGRATION SCENARIO 2.

| Alternative ID  | Migration Time (s) | Migration Reliability | Performance Score |
|-----------------|--------------------|-----------------------|-------------------|
| A               | 177.1              | 0.86                  | 0.99              |
| B               | 177.2              | 0.79                  | 0.76              |
| C               | 186.9              | 0.80                  | 0.73              |
| D               | 177.1              | 0.78                  | 0.71              |
| E               | 177.1              | 0.76                  | 0.63              |
| F               | 177.1              | 0.75                  | 0.61              |
| G               | 177.0              | 0.69                  | 0.45              |
| H               | 177.2              | 0.63                  | 0.35              |
| I               | 223.9              | 0.64                  | 0.02              |
| Criteria Weight | 0.39               | 0.61                  |                   |

towards the most appropriate alternative taking into account the current status of the IoT environment.

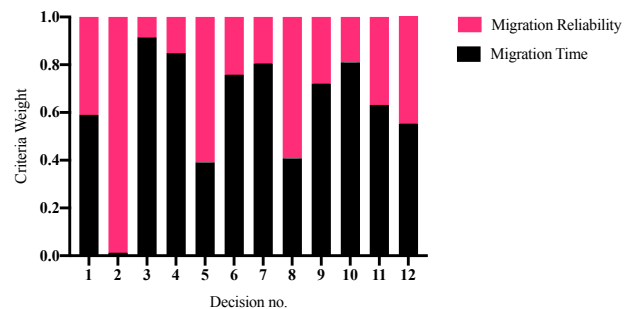


Fig. 8. Criteria Weights for Different Migration Decisions.

#### IV. CONCLUSION

In this paper, a mobility-aware container migration algorithm for Cloudlet-enabled IoT systems has been presented and evaluated. The proposed algorithm has employed an Entropy-TOPSIS integrated multicriteria decision making approach to

select the most appropriate destination Cloudlet for a migrating container. It has been implemented using a specialized simulation tool and compared against exiting migration algorithms. Simulation results have proved the ability of the proposed algorithm to outperform other algorithms in terms of migration time, service downtime, migration reliability and service loss rate. They have also confirmed the ability of the proposed algorithm to perceive the run-time dynamics of the IoT environment and accurately steer the destination Cloudlet selection process.

## REFERENCES

- [1] M. De Donno, K. Tange, and N. Dragoni, "Foundations and evolution of modern computing paradigms: Cloud, iot, edge, and fog," *IEEE Access*, vol. 7, pp. 150936–150948, 2019.
- [2] S. Parikh, D. Dave, R. Patel, and N. Doshi, "Security and privacy issues in cloud, fog and edge computing," *Procedia Computer Science*, vol. 160, pp. 734 – 739, 2019.
- [3] Y. Ai, M. Peng, and K. Zhang, "Edge computing technologies for internet of things: a primer," *Digital Communications and Networks*, vol. 4, no. 2, pp. 77 – 86, 2018.
- [4] S. Shahzadi, M. Iqbal, T. Dagiuklas, and Z. Qayyum, "Multi-access edge computing: open issues, challenges and future perspectives," *Journal of Cloud Computing*, vol. 6, 12 2017.
- [5] S. P. Singh, A. Nayyar, R. Kumar, and A. Sharma, "Fog computing: from architecture to edge computing and big data processing," *The Journal of Supercomputing*, vol. 75, no. 4, pp. 2070–2105, 2019.
- [6] H. Yao, C. Bai, M. Xiong, D. Zeng, and Z. Fu, "Heterogeneous cloudlet deployment and user-cloudlet association toward cost effective fog computing," *Concurrency and Computation: Practice and Experience*, vol. 29, no. 16, p. e3975, 2017, e3975 epe.3975.
- [7] A. Yousefpour, C. Fung, T. Nguyen, K. Kadiyala, F. Jalali, A. Niakanlahiji, J. Kong, and J. P. Jue, "All one needs to know about fog computing and related edge computing paradigms: A complete survey," *Journal of Systems Architecture*, vol. 98, pp. 289 – 330, 2019.
- [8] J. P. Martin, A. Kandasamy, and K. Chandrasekaran, "Exploring the support for high performance applications in the container runtime environment," *Human-centric Computing and Information Sciences*, vol. 8, no. 1, p. 1, 2018.
- [9] A. Celesti, D. Muldari, A. Galletta, M. Fazio, L. Carnevale, and M. Villari, "A study on container virtualization for guarantee quality of service in cloud-of-things," *Future Generation Computer Systems*, vol. 99, pp. 356 – 364, 2019.
- [10] L. F. Bittencourt, J. Diaz-Montes, R. Buyya, O. F. Rana, and M. Parashar, "Mobility-aware application scheduling in fog computing," *IEEE Cloud Computing*, vol. 4, no. 2, pp. 26–35, 2017.
- [11] C. Puliafito, C. Vallati, E. Mingozzi, G. Merlino, F. Longo, and A. Puliafito, "Container migration in the fog: A performance evaluation," *mdpi sensors*, vol. 19, 2019.
- [12] L. F. Bittencourt, M. M. Lopes, I. Petri, and O. F. Rana, "Towards virtual machine migration in fog computing," in *2015 10th International Conference on P2P, Parallel, Grid, Cloud and Internet Computing (3PGCIC)*, 2015, pp. 1–8.
- [13] M. Islam, A. Razzaque, and J. Islam, "A genetic algorithm for virtual machine migration in heterogeneous mobile cloud computing," in *2016 International Conference on Networking Systems and Security (NSysS)*, 2016, pp. 1–6.
- [14] M. M. Lopes, W. A. Higashino, M. A. Capretz, and L. F. Bittencourt, "Myifogsim: A simulator for virtual machine migration in fog computing," in *Proceedings of The 10th International Conference on Utility and Cloud Computing*, New York, NY, USA, 2017, p. 47–52.
- [15] H. Gupta, A. Vahid Dastjerdi, S. K. Ghosh, and R. Buyya, "ifogsim: A toolkit for modeling and simulation of resource management techniques in the internet of things, edge and fog computing environments," *Software: Practice and Experience*, vol. 47, no. 9, pp. 1275–1296, 2017.
- [16] Y. Bi, G. Han, C. Lin, Q. Deng, L. Guo, and F. Li, "Mobility support for fog computing: An sdn approach," *IEEE Communications Magazine*, vol. 56, pp. 53–59, 05 2018.
- [17] A. Machen, S. Wang, K. K. Leung, B. J. Ko, and T. Salonidis, "Live service migration in mobile edge clouds," *IEEE Wireless Communications*, vol. 25, no. 1, pp. 140–147, 2018.
- [18] J. Martin and A. Kandasamy, "Mobility aware autonomic approach for the migration of application modules in fog computing environment," *Journal of Ambient Intelligence and Humanized Computing*, pp. 1–20, 2020.
- [19] C. Puliafito, D. M. Gonçalves, M. M. Lopes, L. L. Martins, E. Madeira, E. Mingozzi, O. Rana, and L. F. Bittencourt, "Mobfogsim: Simulation of mobility and migration for fog computing," *Simulation Modelling Practice and Theory*, vol. 101, p. 102062, 2020.
- [20] M. Behrisch, L. Bieker, J. Erdmann, and D. Krajzewicz, "Sumo - simulation of urban mobility: An overview," in *SIMUL 2011, The Third International Conference on Advances in System Simulation*, 2011, pp. 63–68.
- [21] P. A. Lopez, M. Behrisch, L. Bieker-Walz, J. Erdmann, Y. Flötteröd, R. Hilbrich, L. Lücken, J. Rummel, P. Wagner, and E. Wiessner, "Microscopic traffic simulation using sumo," in *2018 21st International Conference on Intelligent Transportation Systems (ITSC)*, 2018, pp. 2575–2582.
- [22] S. R. Komaragiri and D. N. Kumar, *Multicriterion Analysis in Engineering and Management*. PHI Learning Pvt. Ltd., New Delhi, India, 2010.
- [23] J. Papathanasiou and N. Ploskas, *Multiple Criteria Decision Aid : Methods, Examples and Python Implementations*. Cham: Springer International Publishing, 2018.
- [24] J. Araujo, P. Maciel, E. Andrade, G. Callou, V. Alves, and P. Cunha, "Decision making in cloud environments: an approach based on multiple-criteria decision analysis and stochastic models," *Journal of Cloud Computing*, vol. 7, no. 1, p. 7, 2018.
- [25] C. Jatoth, G. R. Gangadharan, U. Fiore, and R. Buyya, "Selcloud: a hybrid multi-criteria decision-making model for selection of cloud services," *Soft Computing*, vol. 23, no. 13, pp. 4701–4715, 2019.
- [26] R. Khorsand and M. Ramezani, "An energy-efficient task-scheduling algorithm based on a multi-criteria decision-making method in cloud computing," *International Journal of Communication Systems*, vol. 33, no. 9, p. e4379, 2020.
- [27] S. C. Nayak, S. Parida, C. Tripathy, B. Pati, and C. R. Panigrahi, "Multi-criteria decision-making techniques for avoiding similar task scheduling conflict in cloud computing," *International Journal of Communication Systems*, vol. 33, no. 13, p. e4126, 2020.
- [28] E. A. Elsayed, "Overview of reliability testing," *IEEE Transactions on Reliability*, vol. 61, no. 2, pp. 282–291, 2012.
- [29] C. E. Ebeling, *An Introduction to Reliability and Maintainability Engineering*. Waveland Press, 2019.



# Disaster Recovery in Cloud Computing Systems: An Overview

Abedallah Zaid Abualkishik<sup>1</sup>

College of Computer Information  
Technology  
American University in the Emirates  
Dubai, United Arab Emirates

Ali A. Alwan<sup>2</sup>

Department of Computer Science  
Kulliyah of Information and  
Communication Technology  
International Islamic University  
Malaysia, Selangor, Malaysia

Yonis Gulzar<sup>3</sup>

Department of Management  
Information Systems  
College of Business Administration  
King Faisal University  
Al-Ahsa, Saudi Arabia

**Abstract**—With the rapid growth of internet technologies, large-scale online services, such as data backup and data recovery are increasingly available. Since these large-scale online services require substantial networking, processing, and storage capacities, it has become a considerable challenge to design equally large-scale computing infrastructures that support these services cost-effectively. In response to this rising demand, cloud computing has been refined during the past decade and turned into a lucrative business for organizations that own large datacenters and offer their computing resources. Undoubtedly cloud computing provides tremendous benefits for data storage backup and data accessibility at a reasonable cost. This paper aims at surveying and analyzing the previous works proposed for disaster recovery in cloud computing. The discussion concentrates on investigating the positive aspects and the limitations of each proposal. Also examined are discussed the current challenges in handling data recovery in the cloud context and the impact of data backup plan on maintaining the data in the event of natural disasters. A summary of the leading research work is provided outlining their weaknesses and limitations in the area of disaster recovery in the cloud computing environment. An in-depth discussion of the current and future trends research in the area of disaster recovery in cloud computing is also offered. Several work research directions that ought to be explored are pointed out as well, which may help researchers to discover and further investigate those problems related to disaster recovery in the cloud environment that have remained unresolved.

**Keywords**—Cloud computing; data backup; disaster recovery; multi-cloud

## I. INTRODUCTION

Since its introduction in the commercial sector, cloud computing has undergone a significant change in storing and securing information. With cloud computing, data are run in a collection of nodes including servers and remote computers, which enables users to remotely access the data at any time and from any location. The cloud service providers wish to ensure the delivery of flexible services offered in such a way that keeps users separated from the underlying infrastructure. Cloud computing is important when applied to data recovery due to its flexibility, cost-effectiveness, reliability, and scalability. However, since the internet constitutes an open network for sharing information and conducting transactions, it possesses many security and privacy risks as well as availability issues, particularly for businesses [1,2]. This

problem has been addressed using many different approaches including distributed computing, server clustering, and wide-area networking [3].

Small and Medium Business (SMB) corporations are progressively coming to terms with the fact that the cloud service offers many benefits in terms of managing and facilitating their business. They can acquire immediate access to effective business applications and significantly expand their infrastructure resources, all at a minimal expense[4]. Cloud computing is understood as a strategy to enhance existing capabilities and to dynamically introduce new functionalities without investments in different infrastructures, offer training to new employees, and ensure the accreditation of new software packages to expand IT abilities [5]. In today's business environment, the data services operated by CPs encounter many challenges in ensuring a high level of reliability of data services before and after disasters [6]. Data services must ensure reliability and flexibility through an effective and practical DR plan. The data reliability and flexibility are essential requirements for any firm to maintain financial success and sustain the future growth of the organization [6]. The main issue concerning disaster recovery in the cloud computing context is how to provide an effective plan for data backup and recovery that guarantees high data reliability at a reasonable cost prior to a disaster. Thus, a number of solutions have been offered focusing on disaster recovery and data backup in a single-cloud paradigm [6-9].

This paper attempts to highlight and discuss the existing research done on disaster recovery in cloud computing including single and multi-cloud environments. The surveyed studies are thoroughly evaluated to identify the strengths and limitations of each work. Besides, the current and future trends related to disaster recovery in the cloud environment are discussed. In the focus of the discussion are the major issues concerning data backup and recovery in the cloud paradigm.

The remainder of this paper is structured as follows: Section II explains the four categories of disasters that are most likely to occur. In Section III, an overview of disaster recovery is outlined while its different types are explained in Section IV. The discussion covers the three current recovery techniques of cold site recovery, warm site recovery, and hot site recovery. Section V examines the concept of DR in the

---

This work is fully supported by the Kulliyah of Information and Communication Technology, International Islamic University Malaysia, Malaysia

cloud computing context. The research challenges linked to DR in cloud computing are discussed in Section VI. The most noteworthy research done on data DR in cloud computing is reported and discussed in Section VII. Furthermore, a descriptive summary of the related works covered in this survey is given outlining the merits and demerits of each work. Subsequently discussed are the current challenges such as the number of replications, storage cost, and data reliability in a multi-cloud DR process in Section VIII. Moreover, several future work directions deserving to be explored further are included in observations throughout the paper. The research conclusions are presented in the final Section IX.

## II. TYPES OF DISASTER

Disasters, whether man-made or natural, can result in costly service interruption. For many organizations, adopting cloud computing constitutes the most reliable way of

obtaining a dedicated and shared model that can serve DR at low cost and sustain a high speed of access [10-12]. Disaster is defined as any kind of event that leads to critical or devastating damage to a system and results in compromising the availability and the continuity of the system's operations and services for an unknown period. Thus, due to the huge negative impact of any kind of such disaster on the essential services of the system, many businesses, and public services strive to install effective disaster recovery mechanisms that can preserve the sensitive data and decrease the downtime to the minimum level (service disruption). Disasters can be classified into four main classes based on their nature and type, namely climate disaster deliberated and/or intended disruption, damage or loss of utilities and services, and system equipment malfunction. These four types of disasters are further elaborated in Table I.

TABLE I. DISASTER EVENTS CATEGORIES

| Category                                      | Incident Type   | Description  |
|---|---|--|
| <b>Climate disasters</b>                      | Flood   | <ul style="list-style-type: none"><li>• Rapid and uncontrolled increment in water level in a stream, natural or artificial lake, dam, or coastal area.</li><li>• Fires that cause severe and serious damage in properties can be ignited by inadvertent acts such as lightning, arsonists, smokers, or burning wood or any other inflammable materials.</li><li>• Natural disasters may occur in certain areas on earth. This includes landslides due to heavy rain or heavy objects falling from high places such as rocks.</li><li>• Strong winds with high speed that might strike some regions especially in low atmospheric pressure areas like deserts.</li><li>• The source of this type of natural disaster includes any substances such as chemical, airborne radioactive particles that compromise the surrounding environment and threaten the population, particularly in urban areas. This type of disaster also includes pollution of the air due to the emission of some toxic substances in the event of earthquakes and hurricanes.</li></ul> |
|   | Fire  |  |
|   | Subsidence and landslides                               |  |
|   | Windstorm   |  |
|   | Contamination and environmental hazards                 |  |
| <b>Deliberated and/or intended disruption</b> | Arson   | <ul style="list-style-type: none"><li>• Arson is a deliberate act of setting fire with the intention of vandalism that causes damage to property such as buildings, bridges, vehicles, and private homes.</li><li>• When a group of workers are dissatisfied with their work conditions and want to show a form of refusal to perform work through collective action such as demonstration or strike.</li><li>• This type of disaster comes in a form of threatening others using violence to create fear to accomplish personal, political, or ideological goals. Acts of terrorism do not distinguish between civilians and/or government officials and may target anyone in society.</li><li>• Deliberate destruction or damage of equipment to hinder a particular group.</li></ul>  |
|   | Labor dispute/ Industrial Action                        |  |
|   | Act of terrorism  |  |
|   | Act of sabotage   |  |
| <b>Loss of utilities and services</b>         | Electrical power failure and Network services breakdown | <ul style="list-style-type: none"><li>• This type of disaster encompasses several harmful activities that lead to the interruption of normal electrical power services. Furthermore, a disruption in network services is not considered as an immediate disaster; however, network breakdown is problematic if the outrage negatively affects the ability of the company to provide services to its clients, vendors, and business partners.</li></ul>   |
| <b>Equipment or system failures</b>           | Cooling plant failure<br>A/C failure                    | <ul style="list-style-type: none"><li>• Interruption of the cooling plant that can cause the unavailability of services and facilities.</li><li>• Interruption of the air conditioning system that can cause the unavailability of services and facilities.</li><li>• Interruption of fire suppression that can cause the unavailability of services and facilities.</li><li>• Internal power outage.</li><li>• Piece of equipment that physically fails in such a way as to impair its availability and performance.</li></ul>  |
|   | Fire suppression failure                                |  |
|   | Internet failure  |  |
|   | Equipment failure                                       |  |

#### IV. AN OVERVIEW OF DISASTER RECOVERY

DR refers to planning the minimization of data loss and recovery when such losses occur in terms of the expected legal, regulatory, financial, and reputational effects. Regardless of the type of industry, unforeseen events can bring business operations to a standstill and incur extensive financial loss and/or reputational damage to an organization [3, 12-16]. Therefore, a data recovery plan is critical to maintaining continuity by providing all the solutions and steps needed to restore normal operations as soon as possible. Most business organizations throughout the world rely on data to derive competitive advantage and thrive in the marketplace yet give little thought to potential data losses and the consequences thereof. DR is usually the responsibility of the IT department as it is principally concerned with the recovery of computing systems and data after a breach. A breach may be caused by a natural disaster such as a fire, storm, or flood, yet it can also have man-made causes, such as a power outage, malware, data theft, or other malevolent practices. DR preparedness usually requires the implementation of a Disaster Recovery Plan (DRP) so that the steps and procedures to be followed after an incident and can be codified beforehand [8,12, 17-18].

Thus, a DRP constitutes an essential and necessary aspect of any functional enterprise [3, 8, 17-21]. It consists of a set of procedures and predefined policies that attempt to ensure the continuity of the critical business services and sustain the organization's mission by providing the usual services to the target clients during and after the disaster. One of the essential tasks of any effective DRP is to help firms rebuild and restore their system after the failure of their software and hardware components. Unlike fault tolerance that ensures the continuity of the operations due to a failure occurring in one of the system components, DR is more concerned with serious damage and long-term disruption of the business services [8,17, 20]. A DRP intends to manage and maintain the system that is affected by events that have an immediate impact on the availability and the continuity of the services. This includes but is not limited to recovery against cyber-attacks that threaten security, natural disasters, and server outages. A typical disaster recovery plan includes certain steps that ensure the rapid implementation of the DRP to restore the system to its normal state. Many critical parameters should be considered when designing a DRP, which encompasses Critical Business Functions (CBFs), Maximum Acceptable Outage (MAO), Recovery Time Objective (RTO), and Business Impact Analysis (BIA). The most critical parameter in the case of a disaster is CBFs, which include a set of functions very critical in sustaining the business continuity of the services by the organization. Any long-term interruption of these services means that the organization fails to execute its critical operations. There is also a strong relationship between the service disruption and the maximum time that a function can be unavailable without affecting the main mission of the organization, which is called (MAO). Also, to ensure smooth continuity in the organization's service, the maximum time before recovery should be computed accurately. It should be noted that for any DRP the RTOs must be either greater than or equal to the MAO since the RTO represents the timeframe

for the recovery process to be completed. Similarly, the BIA represents the risk analysis that examines the CBFs and the MAO in order to determine the impact the function failure has on a business. The BIA can also be used to specify the priority of recovery attempts that need to be accomplished [8,17].

#### V. TYPES OF DISASTER RECOVERY

The various types of DR upon which others are built include cold site recovery, warm site recovery, and hot site recovery. As shown in Table II, the current technology standards for platform recovery can be implemented using one of the following techniques [3, 7, 17, 22]:

**Hot site:** Computers are configured and equipped with a list of software and data to accept the production load when the primary server is down. The fail-over is typically (if required) obtained through cluster configuration. The standby cluster configuration is separate and distinguished from the master database configuration.

**Warm site:** Computer hardware is pre-configured and supplied with a list of software. Once a disaster occurs, the Domain Name System (DNS) is switched and redirected to the backup site, and the server accepts the production load. The services have to be restarted manually.

**Cold site:** In cold site, the hardware elements of the computer need a set of software associated with a set of data to be generated or restored before promoting the system into a productive state.

Generally, if a disaster occurs at one of the sites, the business is successfully switched to other sites. DR for large-scale hazards usually requires shutting off the power to all utilities and evacuating the facility if required, with the exact tasks to be performed to protect personnel and save lives as identified in the DRP. Many natural disasters, such as flooding or major fires, can cause extensive damage to storage media, in which case specialized and professional data recovery techniques must be used. The physical recovery of data is conducted through different means depending on the extent of the damage, and it may require the use of custom hardware and software recovery systems such as spin-stand data recovery from physically damaged media and data carving [7, 22-23].

TABLE II. STANDARDS PLATFORM RECOVERY

| Option    | RTO Coverage            | Description   | Cost Indication |
|-----------|-------------------------|---|-----------------|
| Hot Site  | Minutes (5 min – 4 hrs) | The hot site option needs a high attention level from the administrative staff of the organization. The age of data is dependent on the data recovery strategy. | High            |
| Warm Site | Hours (4 – 24 hrs)      | The warm site option denotes that the organization has sufficient resources to recover the system. Nevertheless, some extra work is needed to make it live.     | Medium          |
| Cold Site | Days (1 – 7 days)       | The cold site needs to reconstruct the system in a way the recovered data is transferred to another location.   | Low             |

Another type of DR concerns the backing up of critical business data into one or more geographically distributed DCs so that there is a very low probability of all the sites being affected at once. An important aspect of DR is information assurance, which is implemented through multiple Network Attached Storage (NAS) and Storage Area Networks (SANs). Information assurance and recoverability can also be ensured using grid computing and cloud storage. The DRP developer considers the cost involved (including an evaluation of the cost of planning against the cause of failure), the optimal facility location, the optimal data allocation units and the method to be followed for data replication. In consequence, autonomous and semi-autonomous remote data backup and recovery processes have proven to be more popular as storage costs have decreased and bandwidth increased.

## VI. AN OVERVIEW OF DISASTER RECOVERY IN CLOUD COMPUTING

A proactive disaster recovery plan constitutes an essential requirement to sustain long-term success for organizations. A set of well-planned measures that system recovery in the case of a disaster is necessary to ensure the continuity of the services and ensure the availability of daily business activities. A well prepared DRP is very beneficial and can be considered as a long-term investment for many organizations. It is disputable, however, if we acknowledge the fact that the immediate impact of the DRP is unclear and its potential benefits may be rejected. However, cloud-based data backup and recovery has become predominant and proven to be a cost-effective strategy compared to other non-cloud-based approaches [7]. In the cloud environment, the idea of virtualization is no longer relevant to the specifications of the hardware on which it runs. This independency between virtualization technology and hardware often means that organizations are able to safely migrate their data, OS, and software tools to the cloud taking into consideration the financial advantages. The performance of the recovery process is considerably influenced by the network bandwidth and the scalability of services. In other words, high network bandwidth with sustainable scalability of services ensures the rapid commencement of the recovery process. After a disaster, all operations can be re-executed again within a few hours according to the compatibility of the IT structure and the cloud-based DR. It is worth noting that most of the data backup and recovery processes are fully automated and requires either minimal or no human intervention [7,24,25, 37 - 38].

The primary importance of utilizing cloud architecture for implementing a DR strategy is the consequent increase in the overall resilience of the organization's processes and applications. Most CSPs use the geographically distributed model of data backup and redundancy so that companies experiencing widespread outages in their networks due to a disaster can recover within a few hours and with minimal disruption [10]. For example, the Amazon cloud stores mission-critical customer applications at multiple geographically dispersed DCs and uses the "fail gracefully" design philosophy. If there is a momentary outage in an application at one location, the customer is notified immediately. The application is then automatically switched to

another location while downstream circuit breakers prevent any failure of processes and interfaces that rely on that application [26]. An important concern in DR service delivery is continuity of service to enable applications to come back online very soon after a disaster [3, 15, 18, 21, 25, 27- 28, 37].

There are numerous benefits of adopting disaster recovery in the cloud. Nevertheless, several weaknesses may prevent people from exploiting disaster recovery in the cloud. Table III summarizes the advantages and disadvantages of DR in the cloud [17, 29 - 30,38].

TABLE III. ADVANTAGES AND DISADVANTAGES OF DISASTER RECOVERY IN THE CLOUD

| Advantages  | Disadvantages  |
|---|--|
| <ul style="list-style-type: none"><li>• The arrival and maturity of cloud computing represent a paradigm shift where many of the same functionalities can be shifted to the cloud.</li><li>• DRPs using cloud architecture are attractive for small and medium-sized enterprises.</li><li>• Companies can "outsource" their computing requirements and continuity planning to CSPs.</li><li>• The cloud has a rapid turnaround time with outages lasting no more than a few hours.</li><li>• In-house personnel can work with the CSP to redirect customers to the cloud during a disaster.</li></ul> <p>The entire process remains transparent to customers worldwide who do not experience the effects from the disruption.</p> | <ul style="list-style-type: none"><li>• Customers may be concerned about security and data confidentiality since company data are transferred to a third party.</li><li>• A company has no control over where its data will be stored.</li><li>• There have been many incidents of company insiders engaging in malpractice.</li><li>• Companies become dependent on CSP.</li></ul> <p>The long-term viability of the CSP becomes a source of concern for the company.</p> |

## VII. ISSUES AND CHALLENGES OF DISASTER RECOVERY IN THE CLOUD

Since its adoption by a large number of corporations in the world, cloud computing has become an indispensable element in running the essential business operations for large, medium, and small-scale organizations. This is due to its unique ability to ensure the availability of the services and provide resources that are efficient and reliable while maintaining a reasonable cost. The cloud model relies on the concept of pay-as-you-go, which means the user can request the needed resources from the cloud service provider and be billed according to the used resources. Many service models have been incorporated in cloud computing. This includes but is not limited to Infrastructure as a Service (IaaS), Software as a Service (SaaS), and Platform as a Service (PaaS). Other service models are also provided by the cloud providers, for instance, Database as a Service (DBaaS). Despite the tremendous benefits of cloud computing in running the essential business operations for organizations, some are reluctant to fully adopt the cloud paradigm due to security and privacy issues. Thus, cloud computing has not been fully exploited by many organizations.

A variety of reasons and challenges may prevent many organizations from moving towards disaster recovery planning in cloud computing. In the following, we outline the most critical factors that contribute to rejecting cloud computing [2, 13, 17, 22, 24- 28, 31, 36 - 38]:

**Lack of Full Control of Data:** Sharing data with cloud providers can result in losing the full control of data. Since the data backup is executed by the cloud service provider, clients may feel concerned about their data dependency with the CPs and the risk of data loss. Hence, it is crucial for these organizations that they select the most reliable service provider who can guarantee the integrity and the privacy of their data. Given these concerns, many organizations may be reluctant to move their businesses on the cloud.

**Operation Cost:** Operation cost to run the organization's business on the cloud constitutes a critical factor that influences the decision to adopt it. However, the actual cost of running user business on the cloud after switching to a data recovery service is reduced. This reduction in the operating cost may attract many users to adopt the cloud as their preferred platform to run their businesses. The goal of any cloud service provider is to always propose an effective data recovery plan with the least cost. The operating cost of disaster recovery comprises of the following components:

1) Setting up and implementation costs, which denotes the cost of migrating and implementing the organization's business on the cloud. This cost will reduce in the long-term run for the business.

2) Operation cost represents the estimated cost for the daily activities to operate with the data. This includes the operating costs of data storage, transfer, and processing.

3) Disaster costs indicate the total cost of data recovery in the event of a disaster as well as the estimated cost of the damage for unrecoverable disasters. The potential cost of the disaster has a significant impact on the total cost of the services on the cloud.

**Speed of Response in Failure Detection:** The duration of the time to detect and report to the system failure is very crucial to sustain a high level of availability and reliability. The speed of response to system failure reflects the period in which the system is down and all services are inoperable. Therefore, it is an essential objective for any cloud provider to ensure a fast reaction to the service disruption of the system. However, in certain cases, multiple backup sites are engaged, which makes it difficult to immediately distinguish between service disruption and network failure and take the necessary action for detecting and reporting the problem.

**Security:** A cyber-terrorist attack is a typical example of a man-made disaster whereby the system resources are attacked for a variety of reasons. Such attacks may cause data corruption and destroy the system. Hence, any form of data protection must ensure a high level of security and rapid data recovery. They constitute the key elements that influence any decision to adopt disaster recovery services.

**Replication Latency:** The concept of a disaster recovery plan relies on performing data backup through replication. There are two different strategies of data replication that can be utilized, namely synchronous and asynchronous replication strategies. Synchronous replication strategy aims at ensuring a high probability of fulfilling the requirements of the Recovery Point Objective (RPO) and Recovery Time Objective (RTO). Nevertheless, the synchronous replication strategy incurs

higher cost than the asynchronous replication strategy, which in turn may negatively affect the system performance. A larger number of tiers in the web application leads to a significant increment in the Round Trip Time (RTT) between the primary site and the backup site. Although the asynchronous replication strategy is cheaper, it does not deliver the same level of quality service for disaster recovery. Thus, organizations should strike a balance between cost and desired performance taking into consideration the requirements of their particular situation. Furthermore, the latency in data replication constitutes a major concern when deciding whether to adopt the cloud as their preferred platform to run the business.

**Security of Data Storage:** One of the essential benefits of cloud services is that it offers an adequate solution to the issue of data storage. It allows organizations to store their data by providing unlimited space at a reasonable cost. The extensive usage of cloud services leads to a steady increment in the amount of data required for storage. Cloud storage services offer greater flexibility and thus save the budget. Using cloud storage requires less investment than purchasing conventional data storage devices. The architecture of a cloud storage system comprises of four layers, namely physical storage, infrastructure management, application interface, and access. The smooth and reliable running of the applications requires a distributed computing environment that ensures availability, reliability of services, and balancing the workload among all servers. However, data security requires centralized storage in which data are placed in one single storage point. This means that the security of the stored data is at high risk if any failure occurs on the cloud service provider.

**Lack of Redundancy:** When a disaster occurs on the primary site running the services, the cloud service provider immediately activates the secondary site and redirects the incoming requests and services toward the secondary site to ensure the continuity of the business. Running services on the secondary site will have negative implications on the future data backup process as no replication technique (synchronous or asynchronous) can be performed. Failure in running future data backup due to the outage of the primary site thus increases the risk of data loss since one single local storage (secondary site) is available. However, this issue can be easily resolved once the primary site is restored. Overall, any disaster recovery strategy should provide the best solution possible to ensure the precise assessment of all the potential types of risk and examine their negative implications.

## VIII. SOLUTIONS OF DISASTER RECOVERY IN CLOUD COMPUTING

This section examines several of the proposed solutions that are relevant to disaster recovery in the cloud computing environment. We attempt to evaluate these solutions by highlighting the merits and limitations of each approach. A summary table given at the end describes some characteristics of the works considered.

The study completed by Pokharel et al. [32] introduces the Geographical Redundancy Approach (GRA) to disaster recovery in the cloud system. GRA is analyzed using the Markov model, and the experiment result shows that it

accomplishes a high availability and survivability while sustaining a low downtime and low cost. However, the proposed approach is not evaluated in terms of measuring the RTO and the RPO that are considered an important measure in evaluating any DR solution. Most importantly, the proposed solution fits only single-cloud systems and may not apply to multi-cloud systems where multiple remote independent clouds are interconnected.

A comprehensive survey is offered by Wood et al. [27] who list the current disaster recovery solutions and practices concentrating on the most critical factors that affect the disaster recovery process. The three categories of disaster recovery mechanisms that are defined are the hot backup site, warm backup site, and cold backup site. The study also discusses the issue of failover and failback that may occur in the event of a disaster, emphasizing on how to restore the control to the primary site and ensure the continuity of the business-critical services.

The study completed by Jian-hua and Nan [33] describes the typical cloud storage architecture that consists of a storage layer, an infrastructure management layer, an application interface layer, and an access layer. It also explains the typical architecture of disaster recovery deployment in the cloud system that manages the cloud storage in the inter-private cloud model. It stores the application data in the server, remotely connected to another set of backup servers distributed over different areas. Each backup server has another two backup servers, the local backup server (LBS) and the remote backup server (RBS). An incremental data backup approach is used to progressively update the data in order to decrease the usage of network bandwidth and accelerate the data backup process. Several enhancements in the service experience lead to reduced data traffic and transmission cost, which includes carrying out data compression and encryption before the data backup process. The model is designed to work in a single-cloud environment that replicates the original data. Creating one single replica is very crucial and increases the risk of data loss particularly in the event of a disaster.

The work introduced by Javaraiah [34] highlights the issue of online data backup in cloud computing systems. The approach concentrates on managing the data backup process on the consumer's premises to reduce cost. The approach is designed to handle complicated issues associated with the online data backup process in the cloud along with DR. Among the critical issues considered is eliminating the dependency on other cloud providers when performing the data backup operation. Various experiments are conducted, and the results have shown that the proposed solution achieves low costs data backup and simplifies the migration process of data from one CP to another. Nevertheless, this work is limited as it focuses exclusively on the issue of data DR in a single-cloud environment and does not address the maintenance of business services during and after a disaster.

Sengupta and Annervaz [31] address the issue of disaster recovery in multi-sites architecture where the data backup resides in multiple distributed locations. Data Distribution Plan for multi-site Disaster Recovery (DDP-DR) is proposed that offers different plans for data distribution based on

Protection Level (PL) and Placement Constraint (PC). PL denotes the degree of reliability required by the client against the simultaneous datacenter failures while PC denotes the constraint on some DC locations either to be included or excluded from the list of potential locations for the data backup. DDP-DR derives the optimal plan based on the most critical business and operational factors such as cost of data storage and replication, Recovery Time Objective (RTO), and Recovery Point Objective (RPO). Several experiments are conducted to evaluate the efficiency of the proposed solution in different scenarios. However, the proposed solution does not include computing the network cost for data transmission during the backup process and is limited to one client. In some real-life scenarios, there may be more than one client within the same DR architecture.

Grolinger et al. [35] discuss the problem of disaster data management. They emphasize that most of the current data management solutions designed for disaster recovery lack the integration capabilities in order to minimize the negative impact on user data. The proposed framework called Knowledge as a Service (KaaS) handles the cloud data management process during a disaster. It stores as much as possible from the disaster-related data, thus sustaining the interoperability and the integration of the data. Facilitating data integration relies on using knowledge acquisition and knowledge delivery. Knowledge acquisition includes information extraction and retrieval to develop a sound structure for the disaster data while knowledge delivery is used to integrate information from different data sources and forward it to the target clients. However, the proposed framework is not tested and evaluated empirically in order to determine its efficiency and effectiveness. Moreover, not discussed is the issue of disaster recovery in multi-sites where backup data need to be distributed among several remote locations. Lastly, the proposed solution does not incorporate the issue of deriving the optimal plan for data backup during the disaster.

Saquib et al. [6] proposed a new model named Disaster Recovery as a Service for database applications in cloud computing systems. The proposed model provides a solution for disaster recovery with zero data loss and fast recovery. The proposed model exploits the synchronous technique for data replication to ensure minimum RPO and RTO. However, the study lacks the empirical comparison with other cloud-based disaster recovery solutions that would determine its effectiveness. Moreover, the solution is limited to single cloud systems and may not fit multi-cloud systems.

Satoshi Togawa and Kazuhide Kanenishi [14] introduce a new framework of disaster recovery for e-learning systems that sustain business operations during natural disasters such as earthquakes and tsunami. A prototype that works in a private cloud model is developed based on IaaS architecture. The proposed framework incorporates a distributed storage system to ensure that the framework continues sustaining e-learning services even after the disaster. Several experiments are conducted that prove its effectiveness. However, the work fails to examine the framework in terms of the critical business operational metrics of cost, RTO, and RPO. In addition, it is tailored to work in a single cloud environment

where only one single data backup is performed. Any failure in the backup site may thus result in data loss and long-term service disruption.

Lenk [8] focuses on the issue of data deployment for distributed systems in the event of a disaster. The proposed deployment method utilizes the Cloud Standby Disaster Recovery for warm standby in the cloud and runs on different clouds with many cloud providers. The method enables independent and automated data deployment. The method is tested in several experiments, and the results show that the recovery time is reduced significantly. Nevertheless, the fault-tolerance of the deployment method is not investigated.

Jena and Mohanty [9] investigate the issue of disaster recovery in intercloud systems exploiting the genetic algorithm for resource allocation. The main aim is to provide fast track and balanced mapping procedures for impatient tasks in the cloud system. The proposed approach utilizes the genetic algorithm and Pareto optimal mapping to manage resource allocation while sustaining a high utilization rate of the processors, high throughput, and producing a low carbon footprint. A variety of experiments are conducted to evaluate the performance of the proposed approach. The proposed

solution is tested by producing the optimal plan for resource allocation for impatient tasks. Nevertheless, the proposed solution is limited to a single cloud with multiple data centers distributed over many remote locations. Besides, not considered are managing the replication plan to generate a minimum number of data backup without compromising the reliability requirements for the user. Also, the algorithm is not evaluated in terms of Recovery Time Objective and Recovery Point Objective. These two parameters are very essential in the investigation of disaster recovery solutions in cloud computing systems.

Sabbaghi et al. [3] propose a framework formed by integrating five essential types of proven redundancy techniques that have a major impact on the uptime of services in cloud DCs. This work focuses on how disasters can be controlled in a cloud computing DC and how to keep the organization's business running in the event of a disaster. The proposed framework is evaluated through a survey of networking professionals and experts. The results are provided for evaluation but do not include the performance metrics RTO and RPO. Table IV summarizes the previous approaches of DR in the cloud computing environment.

TABLE IV. SUMMARY OF PREVIOUS APPROACHES OF DISASTER RECOVERY IN THE CLOUD

| Author and Year         | Type of DR Cloud | Scope     | CP No. | Parameters                               | Limitations  |
|-------------------------|------------------|-----------|--------|--|--|
| Pokharel et al. [32]    | Single Cloud     | DR        | 1      | Infrastructure cost, Downtime            | Did not discuss RTO and RPO  |
| Wood et al.[27]         | Single Cloud     | DR        | 1      | Cost, RTO, RPO, Performance              | Did not provide RTO and RPO analysis to ensure continuity  |
| Jian-hua & Nan [33]     | Single Cloud     | DR        | 1      | Storage cost                             | Did not present RTO and RPO analysis; no experimental result   |
| Javaraiah[34]           | Single Cloud     | DR        | 1      | Infrastructure cost                      | Did not discuss the parameters RTO and RPO; did not ensure continuity  |
| Sengupta & Annervaz[31] | Single Cloud     | DR        | 1      | Storage cost, Protection level, RTO, RPO | The proposed model only considered the case of one customer, single-cloud multiple DCs.  |
| Grolinger et al.[35]    | Single Cloud     | DR        | 1      | Storage space                            | Did not discuss data recovery; did not provide the full framework; did not use the performance metrics RTO and RPT to test the framework |
| Saqib et al.[6]         | Single Cloud     | DR and BC | 1      | Infrastructure cost, RTO, RPO            | Did not provide performance analysis; did not ensure BC  |
| Togawa & Kanenishi[14]  | Single Cloud     | DR and BC | 1      | Migration Time                           | Did not discuss the parameters RTO and RPO; did not ensure BC  |
| A. Lenk[8]              | Single Cloud     | DR        | 1      | Cost, Time, RPO                          | Did not discuss the parameters RTO and RPO   |
| Jena & Mohanty[9]       | Single Cloud     | DR        | 1      | Cost, Time                               | Did not discuss the parameters RTO and RPO   |
| Sabbaghi et al.[3]      | Single Cloud     | DR        | 1      | Cost, Time                               | Did not discuss the parameters RTO and RPO   |

## IX. DISCUSSION AND FUTURE WORK RECOMMENDATION

This section highlights and discusses the issues and challenges relevant to DR examined in this paper. Also, this section presents the future directions towards DR in cloud computing. Most of today's company services rely on IT systems, some of them being of critical importance to society such as financial services and health care services. Even a very short period of downtime or a very small amount of data loss may result in huge economic losses or social problems. Therefore, most important business and public services use

DR mechanism in order to protect their critical data and minimize the downtime caused by catastrophic system faults. Among the types of technologies adopted in DR, systems are asynchronous backup or continuous synchronization of data and preparing standby systems in geographically separated places. During the past decade, cloud computing has emerged as the new service paradigm and is gaining in popularity. A vast number of services are now being built on the cloud platform. These services utilize the resources of a cloud platform with a pay-as-you-go pricing model. The on-demand nature of cloud computing vastly reduces the cost and RTO of

DR whose peak resource demands are much higher than average demands. However, data DR represents a kind of service that possesses the highest data reliability requirements. How to perform data DR service using the cloud computing paradigm to maximize data reliability while reducing cost and RTO still constitutes a challenge. Similar to other computer systems, cloud computing systems also risk dependency, failure detection, security, human-caused damage, natural disasters, and the like. All of these risks may lead to cloud service interruption or even loss of data. To ensure high data reliability, CSPs deploy several data protection strategies. For example, popular distributed storage systems currently used in cloud platforms such as Amazon S3, Google GFS, and Apache HDFS have adopted 3-replicas data redundant mechanism by default. However, in the case of an entire data center failure, data may still be lost. In order to avoid this problem, some CSPs use geographical data dispersion to protect the most critical data, while data centers in distinct locations owned by one CSP use similar software stack, infrastructures purchased in bulk, operation mechanism, and management team. There are still risks of multiple data center failures due to common causes shared across data centers. Also, the number of data centers owned by one CSP is limited. In case some of them become unreachable, the surviving data centers may not apply to customers due to geographical distance, especially in the event of emergency data restoration. Thus, no matter how many preventive measures are being taken, the possibility of data reliability disruption in a cloud cannot be ignored. According to public reports, even the most advanced cloud services have encountered several instances of wide-area outages and the shutting down of public services. Therefore, the best solution for DR service is to utilize multiple data centers from different CSPs. Some researchers focused on how to backup data in a cloud computing environment. Javaraiah [34], for example, introduces online backup and DR and eliminates the dependency on CPs. Sengupta and Annervaz [31] proposed a plan for multi-site DR where backup data can reside in multiple data centers, including the public cloud.

DR in cloud computing has the potential to become a frontrunner in promoting a secure, virtual, and economically viable IT solution in the future. One of the challenges for data management in a cloud environment is how to design a model that tests data storage at low cost, and RTO with high data reliability. Below are summarized the most critical issues relevant to DR in cloud computing that can be observed:

**Cloud Data Storage:** DR in the cloud possesses potential side effects that affect data availability and data access performance. Moreover, it inevitably reduces the replication level of cloud data, and the location of replicas becomes more important which needs further research focusing on data access performance.

**Cost-effective:** The cost-effective cloud data storage solution is still at its validation stage, where the approaches provided are based on experimental environments. Therefore, effective solutions are needed to focus on implementing a prototype of the solution in the cloud.

**Privacy and Confidentiality:** A significant and critical issue is that cloud data storage must guarantee privacy and confidentiality of the data used for DR. Therefore, an effective approach that addresses the issue of privacy and data confidentiality in the cloud data storage is required.

## X. CONCLUSION

This paper has discussed and examined the issue of disaster recovery in the cloud computing environment. An in-depth analysis of the state of the art for DR in cloud computing has been given, together with an overview of the process of disaster recovery for computer systems. The elements of DR in cloud computing have been reported, which includes overview, definition, and types of DR. Also discussed were the details of cloud-based DR analyzed using traditional approaches. In addition, we also identified the main issues and challenges of DR mechanisms that need to be resolved. Several disaster recovery platforms have been described. A comprehensive review of the previous studied of DR in the cloud in both public cloud and privately-owned resources has been conducted. The paper concludes that data DR services must ensure reliability and flexibility through an effective and practical DR plan that constitute vital initiatives for any organization to prosper and sustain growth. Finally, the paper has examined the current trends in the area of disaster recovery in cloud computing and has pointed out future work directions in the field of cloud-based DR to identify the most recent issues and challenges that need to be explored further.

## REFERENCES

- [1] Alzain MA, Soh B, Pardede E (2011). MCDB: Using Multi-clouds to Ensure Security in Cloud Computing. 2011 IEEE Ninth International Conference on Dependable, Autonomic and Secure Computing, Sydney, NSW, Australia.
- [2] Tebaa M, Hajji SEL (2014). From Single to Multi-clouds Computing Privacy and Fault Tolerance. IERI Procedia, 10, 112-118.
- [3] Sabbaghi F, Mahboubi A, Othman SH (2017). Hybrid Service for Business Contingency Plan and Recovery Service as a Disaster Recovery Framework for Cloud Computing. Journal of Soft Computing and Decision Support Systems, 4(4), 1-10.
- [4] Chen D, Zhao H (2012). Data Security and Privacy Protection Issues in Cloud Computing. 2012 International Conference on Computer Science and Electronics Engineering, Hangzhou, China.
- [5] Marston S, Li Z, Bandyopadhyay S, Zhang J, Ghalsasi A (2011). Cloud computing — The business perspective. Decision Support Systems, 51(1), 176-189.
- [6] Saquib Z, Tyagi V, Bokare S, Dongawe S, Dwivedi M, Dwivedi J (2013). A new approach to disaster recovery as a service over cloud for database system. 2013 15th International Conference on Advanced Computing Technologies (ICACT), Rajampet, India.
- [7] Suguna S, Suhasini A (2014). Overview of data backup and disaster recovery in cloud. International Conference on Information Communication and Embedded Systems (ICICES2014), Chennai, India.
- [8] Lenk A (2015). Cloud Standby Deployment: A Model-Driven Deployment Method for Disaster Recovery in the Cloud. IEEE 8th International Conference on Cloud Computing, New York, USA.
- [9] Jena T, Mohanty J (2016). Disaster recovery services in intercloud using genetic algorithm load balancer. International Journal of Electrical and Computer Engineering (IJECE), 6(4), 1828-1838.
- [10] Prazeres A, Lopes E (2013). Disaster Recovery – A Project Planning Case Study in Portugal. Procedia Technology, 9, 795-805.
- [11] Matos R, Andrade EC, Maciel P (2014). Evaluation of a disaster recovery solution through fault injection experiments. 2014 IEEE



- International Conference on Systems, Man, and Cybernetics (SMC), San Diego, CA, USA.
- [12] Andrade E, Nogueira B (2018). Performability Evaluation of a Cloud-Based Disaster Recovery Solution for IT Environments. *Journal of Grid Computing*, 16(2), 1-19.
- [13] Yang P, Kong B, Li J, Lu M (2010). Remote disaster recovery system architecture based on database replication technology. 2010 International Conference on Computer and Communication Technologies in Agriculture Engineering, Chengdu, China.
- [14] Togawa S, Kanenishi K (2013). Private Cloud Cooperation Framework of E-Learning Environment for Disaster Recovery. 2013 IEEE International Conference on Systems, Man, and Cybernetics, Manchester, UK.
- [15] Chang V (2015). Towards a Big Data system disaster recovery in a Private Cloud. *Ad Hoc Networks*, 35, 65-82.
- [16] Alshammari MM, Alwan AA, Nordin A, Al-Shaikhli IF (2017). Disaster recovery in single-cloud and multi-cloud environments: Issues and challenges. 4th IEEE International Conference on Engineering Technologies and Applied Sciences (ICETAS), Bahrain.
- [17] Alhazmi OH (2016). A Cloud-Based Adaptive Disaster Recovery Optimization Model. *Computer and Information Science*, 9(2), 58.
- [18] Alshammari MM, Alwan AA, Nordin A, Abualkishik AZ (2018). Disaster Recovery with Minimum Replica Plan for Reliability Checking in Multi-Cloud. *Procedia computer science*, 130(C), 247-254.
- [19] Lenk A, Tai S (2014). *Cloud Standby: Disaster Recovery of Distributed Systems in the Cloud*. New York, USA.
- [20] Osama E-T, Munir M, Lela P (2016). Assessing IT disaster recovery plans: The case of publicly listed firms on Abu Dhabi/UAE security exchange. *Information and Computer Security*, 24(5), 514-533.
- [21] Alshammari MM, Alwan AA (2018). Disaster Recovery and Business Continuity of Database Services in Multi-Cloud. International Conference on Computer Applications & Information Security, ICCAIS, Riyadh, Saudi Arabia.
- [22] Khoshkholghi MA, Abdullah A, Latip R, Subramaniam S, Othman M (2014). Disaster recovery in cloud computing: A survey. *Computer and Information Science*, 7(4), 39-54.
- [23] Ameigeiras P, Ramos-Muñoz JJ, Schumacher L, Prados-Garzon J, Navarro-Ortiz J, López-Soler JM (2015). Link-level access cloud architecture design based on SDN for 5G networks. *IEEE network*, 29(2), 24-31.
- [24] Chintureena SV (2014). Ensured Availability of resources in a highly reliable mode through Enhanced approaches for Effective Disaster Management in Cloud. International Conference on Electronics and Communication System (ICECS), Coimbatore, India.
- [25] Aobing S, Tongkai J, Qiang Y, Song Y (2013). Virtual machine scheduling, motion and disaster recovery model for IaaS cloud computing platform. IEEE Conference Anthology, China.
- [26] Jaiswal V, Sen A, Verma A (2014). Integrated Resiliency Planning in Storage Clouds. *IEEE Transactions on Network and Service Management*, 11(1), 3-14.
- [27] Wood T, Cecchet E, Ramakrishnan KK, Shenoy PJ, van der Merwe JE, Venkataramani A (2010). Disaster Recovery as a Cloud Service: Economic Benefits & Deployment Challenges. *HotCloud*, 10, 8-15.
- [28] Liu G, Shen H (2017). Minimum-Cost Cloud Storage Service Across Multiple Cloud Providers. *IEEE/ACM Transactions on Networking*, 25(4), 2498-2513.
- [29] Shi X, Guo K, Lu Y, Chen X (2014). Survey on Data Recovery for Cloud Storage. International Conference on Trustworthy Computing and Services, Beijing, China.
- [30] Attiya I, Zhang X (2017). Cloud Computing Technology: Promises and Concerns. *International Journal of Computer Applications*, 159(9), 32-37.
- [31] Sengupta S, Annervaz KM (2012). Planning for Optimal Multi-site Data Distribution for Disaster Recovery. International Workshop on Grid Economics and Business Models, Paphos, Cyprus.
- [32] Pokharel M, Lee S, Park JS (2010). Disaster Recovery for System Architecture Using Cloud Computing. 2010 10th IEEE/IPSJ International Symposium on Applications and the Internet, Seoul, South Korea.
- [33] Jian-hua Z, Nan Z (2011). Cloud Computing-based Data Storage and Disaster Recovery. 2011 International Conference on Future Computer Science and Education, Xi'an, China.
- [34] Javaraiah V (2011). Backup for cloud and disaster recovery for consumers and SMBs. 5th IEEE International Conference on Advanced Telecommunication Systems and Networks (ANTS), Bangalore, India.
- [35] Grolinger K, Capretz MAM, Mezghani E, Exposito E (2013). Knowledge as a Service Framework for Disaster Data Management. 2013 Workshops on Enabling Technologies: Infrastructure for Collaborative Enterprises, Hammamet, Tunisia.
- [36] Sengupta S, Annervaz KM (2014). Multi-site data distribution for disaster recovery—A planning framework. *Future Generation Computer Systems*, 41, 53-64.
- [37] Mohammad M. Alshammari, Ali A. Alwan, Azlin Nordin, Abedallah Zaid Abualkishik (2020). Data backup and recovery with minimum replica plan in multi-cloud environment. *International Journal of Grid and High Performance Computing*, 12(2), 201-120.
- [38] Mohammad Matar Al-Shammari and Ali A. Alwan. Disaster Recovery and Business Continuity for Database Services in Multi-Cloud. Proceedings of the 1st International Conference on Computer Applications & Information Security (ICCAIS' 2018), 4 – 6 April 2018, Riyadh, Saudi Arabia.

# An IoT based Urban Areas Air Quality Monitoring Prototype

Martin M. Soto-Cordova<sup>1</sup>, Martha Medina-De-La-Cruz<sup>2</sup>, Anderson Mujajico-Mariano<sup>3</sup>  
Universidad de Ciencias y Humanidades  
Lima, Peru

**Abstract**—According to the World Health Organization, the most affected places with the presence of polluting gases and particles in suspension are urban areas due to the emissions corresponding to human activities, they have also caused diseases and deaths in millions of people in the world. This paper describes the process of design and implementation of an electronic prototype applying the Internet of Things concept with a cloud storage and processing service. This device has the purpose of monitoring in real-time the air quality through the presence of pollutant gases and PM10 and PM2.5 suspended particles to carry out later studies that contribute to prevention measures in the health care of the population.

**Keywords**—Air pollution; air quality; Arduino; Internet of Things (IoT); cloud service; MQTT; Air Quality Index (AQI); sensors

## I. INTRODUCTION

Air is a vital resource for all living beings depends mainly on the process of respiration humans. The respiratory frequency of an adult person is 12 to 20 breaths per minute and children 40 to 60 [1]. In the process of respiration, the people are exposed to inhale polluting gases and suspended particles in the air.

According to the World Health Organization (WHO), outdoor air pollution is a major environmental health problem affecting everyone in low-, middle-, and high-income countries. Ambient (outdoor) air pollution in both cities and rural areas was estimated to cause 4.2 million premature deaths worldwide per year in 2016; this mortality is due to exposure to small particulate matter of 2.5 microns or less in diameter (PM2.5), which cause cardiovascular and respiratory disease, and cancers [2, 3]. Likewise, it should be noted that 91% of countries in the world exceed the guidelines established by WHO. The cases number of deaths increases to the most vulnerable are children under 5 years, pregnant women and older adults. Thus, in 2016, there were 600 million cases of deaths children due to exposure air conditions [4]. Fig. 1 shows the cases of deaths in 2018 of Latin American countries [5]. Therefore, there are world meetings to evaluate this problem and reach agreements between countries to limit air pollution.

There are entities that promote the care of the environment through the practice of cycling (transfer to work centers or study centers), recycling and other activities that reduce emissions in some way by human activity [6], in addition the technology has allowed the development of applications providing innovative ideas such as rapid notification of the

existence of natural disasters by emissions that damage the atmosphere, poor conditions of national centers and reserves.

Also, there are monitoring statistics for each region or country made by usually government organizations that show the conditions of air quality. Thus, the Air Quality Index (AQI) value for each pollutant of 100 generally corresponds to an ambient air concentration that equals the level of the short-term national ambient air quality standard for protection of public health. AQI values at or below 100 are generally thought of as satisfactory. When AQI values are above 100, air quality is unhealthy: at first for certain sensitive groups of people, then for everyone as AQI values get higher. The so-called AQI is divided into six categories that corresponds to a different level of health concern. Furthermore, each category is assigned a certain color, which makes it easy for people to quickly determine whether air quality is reaching levels of health risk in their geographic areas [7, 8].

Despite, the people are not informed of these conditions, so their health is at risk when exceeding the permissible ranges and put their lives in danger. In Fig. 2, a cause-effect map related to urban pollution is shown [9]. Therefore, it is necessary to inform the population about these conditions and governorates to take preventive measures. The health of the population should not be affected especially children, elderly, and pregnant women. Apply innovation tools that provide social support to affected populations.

The objective of the research is to monitor air conditions in urban areas using a service in the cloud computing named Adafruit IO by using the device for an analysis that allows spreading this information of the conditions of urban areas at risk. To do this, experimentally, data is collected that belongs to the city of Lima - Peru and is shown.

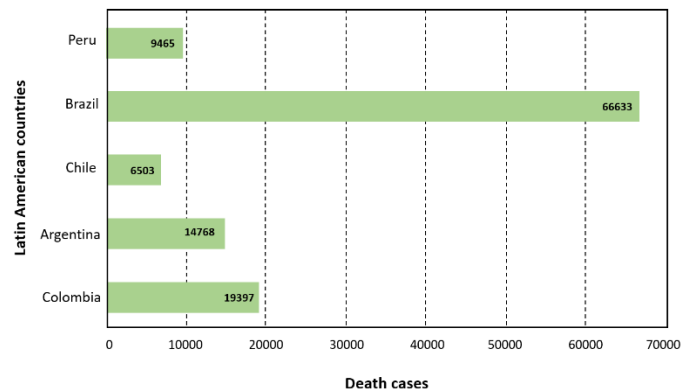


Fig. 1. Cases of Deaths in Latin American Countries.

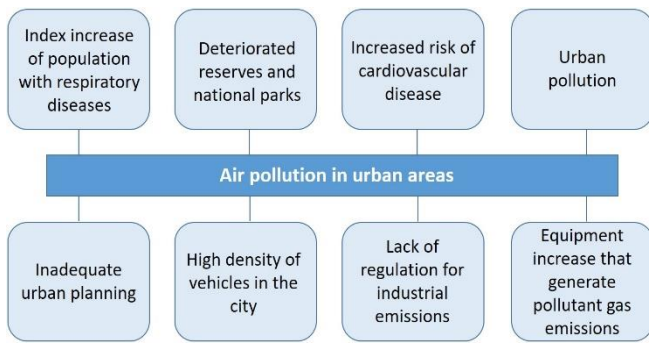


Fig. 2. Causes Effect of Air Pollution.

In what follows, the paper is organized as follows: Section II contains a review of the various proposals for air quality measurement systems, both indoors and outdoors. In Section III, the influence of polluting gases and particles in air suspension on human health will be described. Section IV covers aspects about the procedure for the development of the research project. Finally, the results obtained from the samples collected will be described.

## II. CURRENT CONTEXT OF THE AIR QUALITY MONITORING PLATFORMS

There are need for real-time air quality monitoring systems for micro, small and medium industries and air pollution in the streets so that timely decisions can be taken to avoid environmental degradation. IoT has been proven one of the effective ways for such systems and when merged with cloud computing provides a revolutionary method of management and analysis of data coming from sensors [10, 11].

There are some proposals for indoor monitoring as in [12], an indoor air quality monitoring platform based on IoT technology was developed. It consists of an air quality-sensing device and a web server on cloud computing to monitor indoor air quality. Collected air quality data is transmitted in real time to a web server via LTE mobile network. The IoT device includes pollutant detection sensors, microcontroller and LTE modem, and cloud computing allows visualize indoor air quality data. Likewise, [13] presents an end-to-end indoor air quality monitoring (IAQM) system where the gateway is in charge of processing collected air quality data and its transmission to end-users through a web-server. An adaptation of open-source web application Emoncms was made for live monitoring and long-term storage of the collected IAQM data through sensing technologies, wireless sensor networks (WSNs) and smart mobile.

The author in [14] includes a study where two gaps have been identified: short-term monitoring bias and IAQ data-monitoring solution challenges. The study addresses those gaps by proposing an Internet of Things (IoT) and Distributed Ledger Technologies (DLT)-based IAQ data-monitoring system. The solution helps the penetration of Industrial Internet of Things (IIoT)-based monitoring strategies in the specific case of Occupational Safety Health (OSH).

An air quality monitoring platform under the recommendations of the World Wide Web Consortium (W3C) about the Web of Things (WoT) was proposed in [15]. This

system is built based on a WoT capable of exposing its own Thing Description with resources that can serve requests providing measurements of the indoor attached sensors, depending on the underlying hardware and the application protocol selected. On the other hand, a study was carried out on the river Napo located on the border of Ecuador and Peru using a system applied to environmental monitoring to measure the pollution index in different parts of the area. Analyzing the emissions as the existence of toxic waste that damaged the animals that lived in the round of area [16].

An IoT Based Air Quality Monitoring System was proposed in [17]. It considers Arduino based Air Quality Monitoring setup using MQ135 and MQ7 sensors. The data collected by these sensors will be transmitted to the cloud on back end, here it uses Thing speak. Finally, data analysis was done taking into consideration the dataset from experimental tasks. This analysis helps in deeper understanding of the air quality status such that people will be aware of their effects. Furthermore, [18] presents an IoT platform for monitoring indoor air quality. For the implementation of the system the WoT concept is used to create IoT applications and also uses the CoAP protocol to collect data from multiple sensors. This proposal describes the developed platform: sensor hardware, firmware and software to collect the sensors information and to send to the dashboard for visualization.

An IoT based Air and Sound Pollution Monitoring System is presented in [19] to measure the air and noise pollution levels in industrial environment or by using wireless embedded computing system a particular area of interest. It includes IoT with sensing devices connected to the embedded computing system. Also, [20] considers the design and development of an IoT based industrial pollution real time monitoring, from water quality and air pollution and dangerous gas content like ethane and methane. Thus, it includes sensors and micro-controller and the cloud (AWS).

## III. INFLUENCE OF AIR POLLUTION ON HEALTH

Air pollution affects human health, due to the frequency of exposure to breathing pollutant gases in high concentration were the main effects of allergies and respiratory infections, asthma, bronchitis, premature deaths among others. It is also a factor associated with lung cancer is the cause of 5% of cases. Fig. 3, on the right side, shows a lung of an older adult where the change in color of the lung is manifested by inhaling particles and gases throughout his life.

Among the main concentrations that harm the health of people are carbon monoxide (CO) gas considered a silent killer belonging to the emissions by vehicles in urban areas especially where there is a high vehicle density or a large number of motorized vehicles with high emissions of these gases due to engine malfunction or wear. Table I shows the emissions of monoxide related by the ratings of the AQI.

In the case of particles in PM10 suspension, which are coarse smaller than 10 microns and that affect the respiratory system entering the pulmonary ducts evaluated in the following equation and evidenced in Table II evaluated in 24 hours:

$$I (PM_{10}) = [PM_{10}] * 100/150 \quad (1)$$

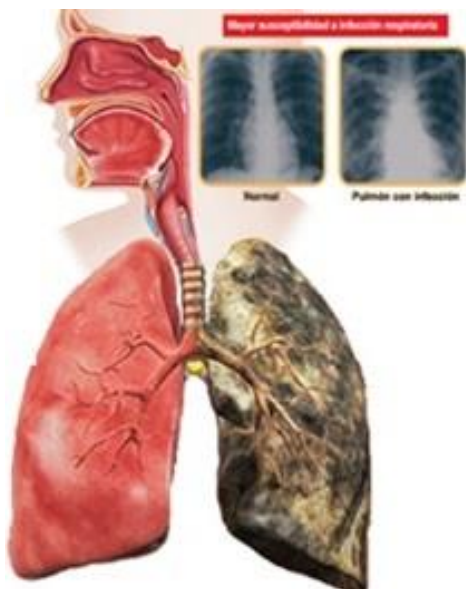


Fig. 3. Effects on the Respiratory System.

TABLE I. RELATED MONOXIDE CONCENTRATIONS THROUGH THE AQI

| AQI Value | Concentration (ppm) |
|-----------|---------------------|
| 0 -50     | 0.0 – 4.4           |
| 51 - 100  | 4.5 – 9.4           |
| 101 - 150 | 9.5 – 12.4          |
| 151 - 200 | 12.5 – 15.4         |
| > 201     | 15.5 – 30.4         |

TABLE II. PM<sub>10</sub> SUSPENDED PARTICLES EVALUATED IN 24 HOURS AQI

| AQI Value | Concentration (µg/m <sup>3</sup> ) |
|-----------|------------------------------------|
| 0 -50     | 0 – 75                             |
| 51 - 100  | 76 – 150                           |
| 101 - 150 | 161 – 200                          |
| 151 - 200 | 250 – 320                          |
| > 201     | > 321                              |

In the case of particles suspended PM<sub>10</sub> which are fine, less than 2.5 microns, affect the respiratory system entering the lungs to the lungs, evaluated in the following equation, and shown in Table III:

$$I (PM_{2.5}) = [PM_{2.5}] * 100/25 \quad (2)$$

Table IV shows the concentration of particles referring to the impact on the health of the person related to the inhalation of PM<sub>10</sub>.

TABLE III. PM<sub>2.5</sub> SUSPENDED PARTICLES EVALUATED IN 24 HOURS AQI

| AQI Value | Concentration (µg/m <sup>3</sup> ) |
|-----------|------------------------------------|
| 0 -50     | 0 – 12.5                           |
| 51 - 100  | 12.6 – 25                          |
| 101 - 150 | 25.1 – 125                         |
| > 151     | > 125                              |

TABLE IV. IMPACT ON HUMAN HEALTH

| Con. | Effects   | Impact      |
|------|---|-------------|
| 200  | Decreased respiratory function and the weakness of the skin       | Moderate    |
| 250  | Increase in respiratory and dermatological diseases               | Moderate    |
| 400  | It affects the entire population                                  | Severe      |
| 500  | Increased mortality in children, pregnant women, and older adults | Very severe |

The permissible limit according to WHO standards, PM<sub>10</sub> is 50 µg/m<sup>3</sup> and for PM<sub>2.5</sub> is 25 µg/m<sup>3</sup>.

#### IV. METHODOLOGY

The Internet of Things device is designed to monitor urban areas using a cloud service, which makes it possible for the population to visualize these parameters through access to the Adafruit IO platform. For this reason, the process of development of the IoT device is described:

##### A. Identification of Variables for Monitoring

In this part, the presence of gases in urban areas was related to the level of affectation in the vulnerable population, evaluated through the Ministry of Health of Peru (MINSA), relating the inhabitants and their departments. Also, the Ministry of the Environment of Peru (MINAM) and WHO have established levels of concentrations of contaminating particles.

##### B. Selection of Gas-Related Electrochemical Sensors

The use of low cost sensors will be used to determine the concentrations of gases in time, so the use of electrochemical sensors was applied, calibrated and analyzed using the equation of the sensitivity curve established by the datasheet. Fig. 4 shows the relationship and the calculation of the equation starting from the value of Rs / Ro for the detection of carbon monoxide (CO) using the MQ-9 sensor.

Fig. 5 shows the relation and calculation of the equation starting from the value of Rs / Ro for smoke detection using the MQ-2 sensor.

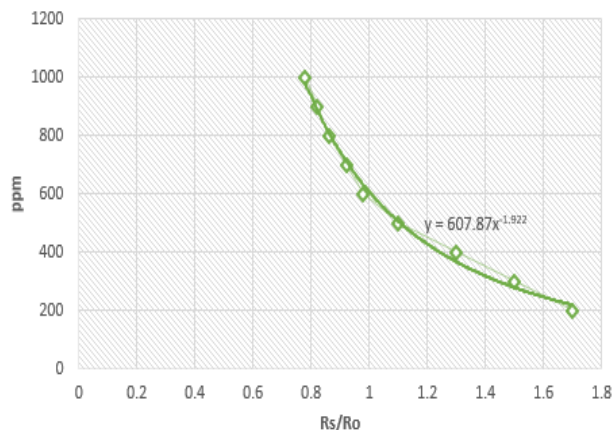


Fig. 4. Sensor Sensitivity Equation MQ 9 for the Detection of Carbon Monoxide (CO).

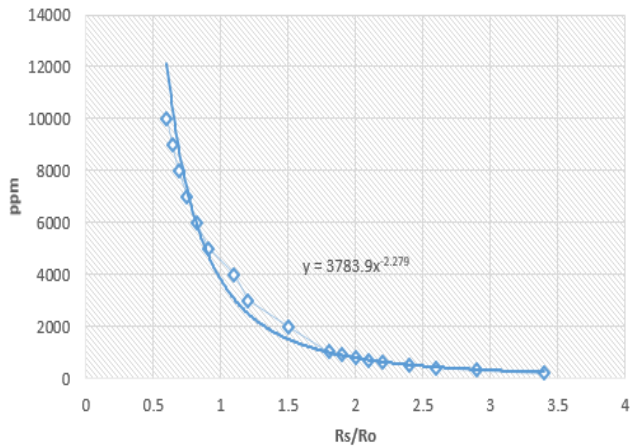


Fig. 5. Sensor Sensitivity Equation MQ 2 for Smoke Detection.

The relationship and calculation of the equation starting from the value of  $R_s / R_o$  for the detection of carbon dioxide using the MQ-135 sensor, shown in Fig. 6, is demonstrated.

### C. Design and Implementation of the Internet of Things Device

For the IoT platform development, the fixed characteristics and requirements are the following:

- Low consumption and low cost;
- Send data to the cloud computing in real-time;
- Allow the visualization of the air pollution conditions of the urban area;
- Small, portable and didactic;
- Monitoring display through an interface;

There are some IoT protocols, such as CoAP and MQTT. Comparatively, CoAP is highly competitive with MQTT, and it includes a mechanism of exploration and observation, but MQTT is much simpler to develop and implement and much more known [21, 22, 23]. Fig. 7 shows that MQTT transaction.

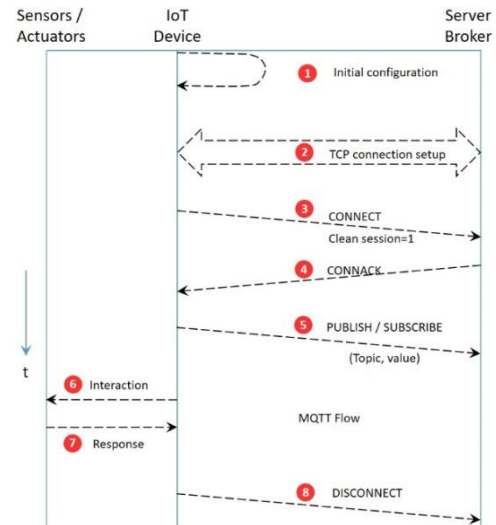


Fig. 7. MQTT Protocol.

The Broker element manages messages and transactions between clients which can be either a subscriber or a publisher. Data is carried in the payload of messages transmitted between clients, mainly a topic and its value. The publisher sends messages to the broker when it has an update message or a periodic message. The broker sends the messages to the subscribers of a specific topic.

Fig. 8 shows an IoT implementation stage. It includes the electrochemical sensors, the humidity, and temperature sensor DHT11, and the dust sensor are connected (obtain data on the density of dust in urban areas). Also, in this process it was decided to work with the Arduino MKR 1000, adapting a program where libraries were used to work with Internet of Things.

Fig. 9 shows the system architecture where IoT device is connected to Adafruit cloud service. It is possible to use Adafruit's IO client API libraries as they include support for MQTT. It allows display data in real-time, online, IoT internet-connected and connect IoT to web services like Twitter, RSS feeds, weather services, etc.

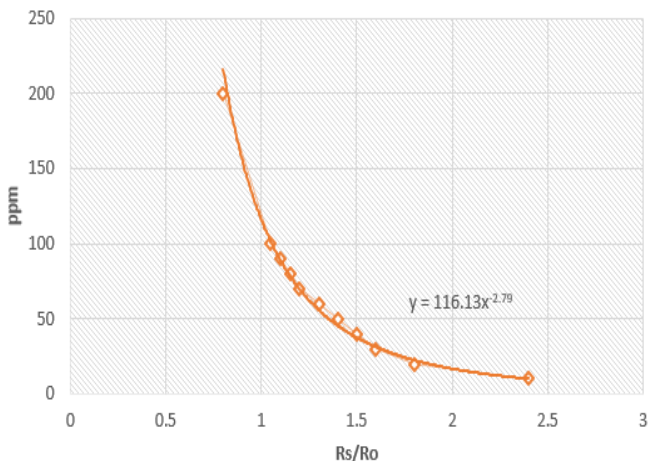


Fig. 6. Sensor Sensitivity Equation MQ 135 for the Detection of Carbon Dioxide (CO<sub>2</sub>).

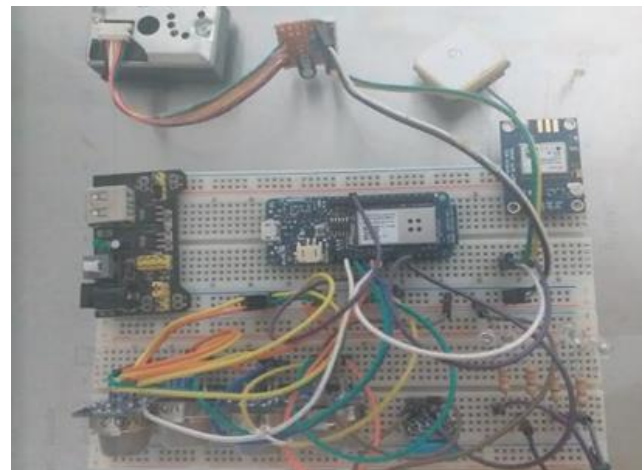


Fig. 8. The Electronic Circuit of the Internet of things Device.

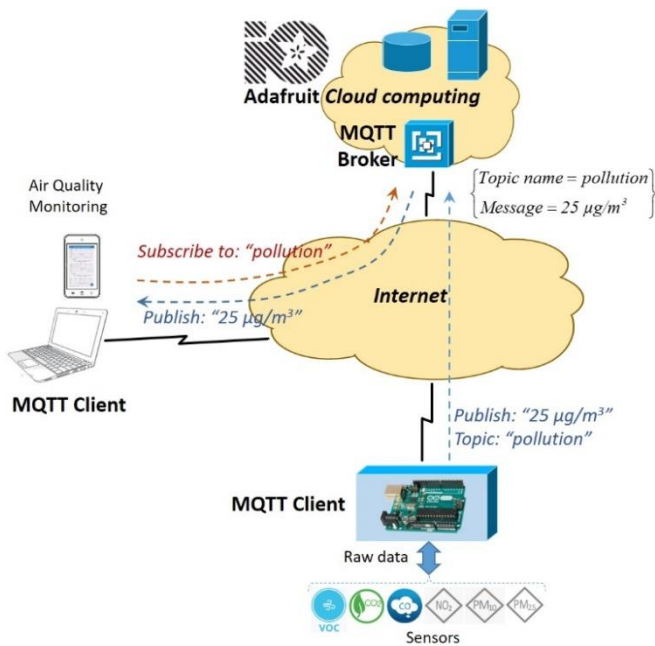


Fig. 9. System Architecture with the use of the Cloud Service.

In this part, the equations in the Arduino IDE were added for the measurement of the electrochemical sensors, to obtain said data in real-time evaluating each message sent by a sensor of 10 seconds. Also, to use the service in the cloud only work with the publication function, which allows visualizing each of the measurements of the sensors in real-time, through the application of an API KEY. Then, Fig. 10 shows some lines of code for reading data.

```

Adafruit_MQTT_Publish humo = Adafruit_MQTT_Publish(smqtt, AIO_USERNAME "/feeds/humo");
Adafruit_MQTT_Publish metano = Adafruit_MQTT_Publish(smqtt, AIO_USERNAME "/feeds/metano");
Adafruit_MQTT_Publish monoxido = Adafruit_MQTT_Publish(smqtt, AIO_USERNAME "/feeds/monoxido");
Adafruit_MQTT_Publish dioxido = Adafruit_MQTT_Publish(smqtt, AIO_USERNAME "/feeds/dioxido");
Adafruit_MQTT_Publish hum = Adafruit_MQTT_Publish(smqtt, AIO_USERNAME "/feeds/hum");
Adafruit_MQTT_Publish temp = Adafruit_MQTT_Publish(smqtt, AIO_USERNAME "/feeds/temp");
Adafruit_MQTT_Publish polvo = Adafruit_MQTT_Publish(smqtt, AIO_USERNAME "/feeds/polvo");
Adafruit_MQTT_Publish densidad = Adafruit_MQTT_Publish(smqtt, AIO_USERNAME "/feeds/densidad");
    
```

Fig. 10. IDE Programming for Reading Data.

### V. RESULTS

Fig. 11 shows the monitoring was carried out in different urban areas mainly in areas where industrial activities and avenues are carried out during hours of vehicular congestion. Visualization was carried out using tools from the Adafruit IO board, which varied depending on the concentration of the gas in the area, the measurements were made in varied periods, which are saved by the cloud service.

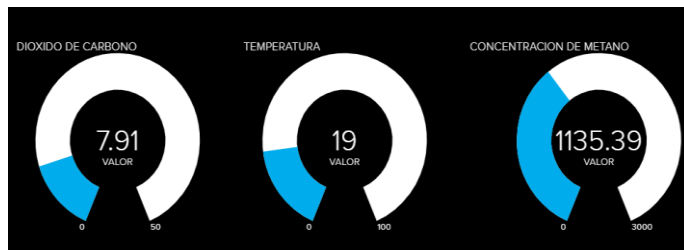


Fig. 11. Visualization of the Existing Concentrations in Ventanilla Area.

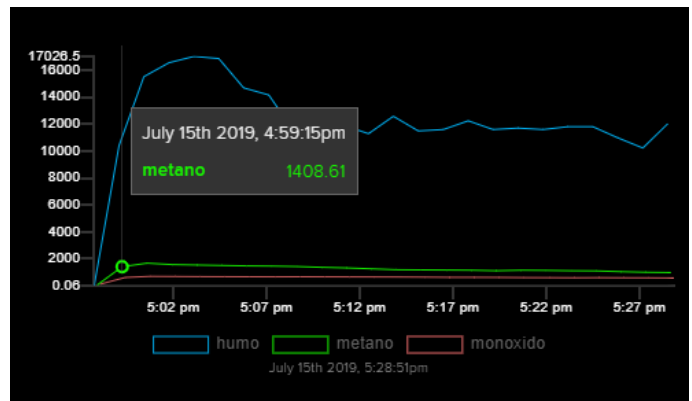


Fig. 12. The Concentration of Measured Gases in Real-Time.

Furthermore, to display several concentrations about time, several gases are configured in the same visualization tool, as shown in Fig. 12.

Through the use of the service, a comparison was made of these concentrations, evaluated in real-time and configured as required by the user, allowing to know the air conditions. Likewise, this information could be shared through a link with other users of the said population to contribute to the reduction of emissions by carrying out activities that could be modified and take care of the area where millions of people live.

### VI. CONCLUSION AND FUTURE WORKS

The Internet of Things has helped to solve tasks in different sectors by managing data and making decisions that could contribute to improving the health of society. For this reason, the use of this device will allow obtaining this information in real-time and publish it through a link so that the population becomes aware of these conditions considering the idea of regulation by monitoring the emissions by the activities human.

Air quality monitoring is a very important task since it would allow adequate policies and control about it in the country, and in this way prevent the affectation of pollution on people's health.

In the future, it is expected to create a network of sensors installed in different areas of Lima, promoting the care of the environment and health care. Also, with the data, this information would be provided to the entities in charge of the monitoring to be disseminated and this information allows knowing the air conditions at the national level.

#### REFERENCES

- [1] Charilaos Chourpiliadis and Abhishek Bhardwaj, Physiology, Respiratory Rate, StatPearls Publishing, January 28, 2019. [Online]. Available: <https://www.ncbi.nlm.nih.gov/books/NBK537306/>.
- [2] World Health Organization (WHO), Fact sheets: Ambient (outdoor) air pollution, 2 May 2018. [Online]. Available: [https://www.who.int/en/news-room/fact-sheets/detail/ambient-\(outdoor\)-air-quality-and-health](https://www.who.int/en/news-room/fact-sheets/detail/ambient-(outdoor)-air-quality-and-health).
- [3] World Health Organization (WHO), Ambient air pollution: a global assessment of exposure and burden of disease, 2016.
- [4] World Health Organization (WHO), Newsroom: More than 90% of the world's children breathe toxic air every day, 29 October 2018. [Online]. Available: <https://www.who.int/news-room/detail/29-10-2018-more-than-90-of-the-worlds-children-breathe-toxic-air-every-day>.

- [5] Expok Comunicación de Sustentabilidad, "Países más contaminados en América Latina – ExpokNews," *Comunicacion de Sustentabilidad y RSE*, 2019. [Online]. Available: <https://www.expoknews.com/paises-mas-contaminados-en-america-latina/>.
- [6] A. Budiarto and T. Febriana, "IoT device used for air pollution campaign to encourage cycling habit in inverleith neighborhood," *International Conference on Information Management and Technology (ICIMTech)*, Yogyakarta, pp. 356-360, 2017.
- [7] U.S. Environmental Protection Agency (EPA), *A Guide to Air Quality and Your Health*, EPA-456/F-14-002, February 2014.
- [8] U.S. Environmental Protection Agency (EPA), *Technical Assistance Document for the Reporting of Daily Air Quality – the Air Quality Index (AQI)*, EPA-454/B-18-007, September 2018.
- [9] Juan M. Rivera Poma, *Modelo de identificación de factores contaminantes atmosféricos críticos en Lima-Callao*, MSc Thesis, Universidad Nacional Mayor de San Marcos, 2012.
- [10] Sumi Neogi et al., "IoT Based Air Quality Monitoring Systems - A Survey," *Proceedings of International Conference on Computer Networks, Big Data and IoT (ICCBI)*, Madurai, Tamil Nadu – India, pp. 752-758, December 19-20, 2018.
- [11] Ashish Gupta and Rajesh Kumar, "An IOT Enabled Air Quality Measurement," *Indian Journal of Science and Technology*, Vol. 11, N° 46, December 2018.
- [12] JunHo Jo et al., "Development of an IoT-Based Indoor Air Quality Monitoring Platform," *Hindawi – Journal of Sensors*, Vol. 2020, January 2020.
- [13] Mohieddine Benammar et al., "A Modular IoT Platform for Real-Time Indoor Air Quality Monitoring," *MDPI - Sensors*, Vol. 18, N° 581, 2018.
- [14] Shengjing Sun et al., "Indoor Air-Quality Data-Monitoring System: Long-Term Monitoring Benefits," *MDPI - Sensors*, Vol. 19, N° 19, 2019.
- [15] Daniel Ibaseta et al., "An IoT Platform for Indoor Air Quality Monitoring Using the Web of Things," *WIT Transactions on Ecology and the Environment*, Vol 236, pp. 45-56, 2019.
- [16] Joel O. Aragon Valladares, *Diseño e implementación de una plataforma de gestión de una red de sensores aplicada a la monitorización de la calidad ambiental en la cuenca del río Napo*, Eng. Thesis, Pontificia Universidad Católica del Perú, 2014.
- [17] Kinnera Bharath Kumar Sai, Subhaditya Mukherjee and H. Parveen Sultana, "Low Cost IoT Based Air Quality Monitoring Setup Using Arduino and MQ Series Sensors With Dataset Analysis," *Procedia Computer Science*, Vol 165, pp. 322-327, 2019.
- [18] J. Esquiagola et al., "Monitoring Indoor Air Quality by using IoT Technology," *IEEE XXV International Conference on Electronics, Electrical Engineering and Computing (INTERCON)*, Aug. 2018.
- [19] Pradyumna Bapat et al., "IoT based Air and Sound Pollution Monitoring System," *International Journal of Research and Analytical Reviews (IJRAR)*, Vol. 6, N° 2, pp. 383-387, April-June 2019.
- [20] Raghava M S et al., "IoT based Industrial Pollution Monitoring," *Research Journal of Engineering and Technology (IRJET)*, Vol. 6, N° 5, pp. 5893-5899, May 2019.
- [21] Khalid Aloufi, "6LoWPAN Stack Model Configuration for IoT Streaming Data Transmission over CoAP," *International Journal of Communication Networks and Information Security (IJCNIS)*, Vol. 11, N°2, pp. 304–311, August 2019.
- [22] T. Yokotani and Y. Sasaki, "Transfer protocols of tiny data blocks in iot and their performance evaluation," *IEEE 3rd World Forum on Internet of Things (WF-IoT)*, pp. 54–57, 2016.
- [23] T. Yokotani and Y. Sasaki, "Comparison with HTTP and MQTT on Required Network Resources for IoT," *International Conference on Control, Electronics, Renewable Energy and Communications (ICCEREC)*, pp. 1–6, 2016.

# Modeling and Interpretation of Covid-19 Infections Data at Perú through the Mitchell's Criteria

Huber Nieto-Chaupis

Universidad Privada del Norte

Av Alfredo Mendiola 6062 Los Olivos, Lima, PERÚ

**Abstract**—In this paper, the criteria of Tom Mitchell based at the philosophy of Machine Learning have been used to interpret data of new cases per week of infections by Covid-19 at Perú. For this, it was constructed a mathematical scheme that encloses the Mitchell's criteria as well as the idea of propagation as commonly used in modern physics to attack complex problems of interactions. With this, both the 2009 season of AH1N1 flu outbreak and the ongoing Covid-19 data were analyzed in terms of task, performance and experience. In contrast with the AH1N1 case, the Covid-19 data do not exhibit any performance in terms of minimize infections at the first weeks of the beginning of the outbreak, suggesting that precise actions to reduce infections have not been taken appropriately.

**Keywords**—Covid-19; epidemiology; machine learning; Tom Mitchell; Monte Carlo

## I. INTRODUCTION

Recently, the unexpected apparition of Corona Virus Disease (Covid-19 in short) [1] has reconfigured the current policies of public health of global operators, forcing them to apply the more robust schemes of recovering and surveillance in the shortest times without an optimal usage of resources: Times, materials and human resources. Although to date, the first wave of pandemic is in most countries reaching its end, it is rather natural to ask about what we have learned from world-wide datasets.

In fact, as seen at all surveillance systems in all countries that are carrying out schemes of care, the understanding of data would exhibit imminent differences among them because the multicultural manifestations of societies as to face from the first moment the arrival of strain. It is also relevant the level of resilience of them for recovering as soon the conditions have shown a certain improvement.

**This paper tries to answer the question: To what extent the schemes of machine learning seen as an universal computational tool can be useful to understand recent data of data from infections by Covid-19?**

In order to answer that question, this paper has selected the Peruvian case that exhibits a remarked difference between current Covid-19 data and that of the 2009 AH1N1 season. Current data between March and July exhibit peaks and fluctuations, facts that would reinforce the hypothesis that in more cases (countries) the dynamics of spread and subsequent infections by Covid-19 appears to be strongly related to randomness. In this manner, this paper has assumed *a priori* that the time evolution of rate of infections is to some extent dictated by the rules that govern the propagation as commonly seen in physics and that was developed by Feynman [2]. In

consequence one can postulate that the action of spread and infection by virus follows the mathematical structure of a propagator integral than can be written down as:

$$\mathcal{P}_{1 \rightarrow 2}(t_2) = \int G(t_2 - t_1) \mathbf{H}(t_1) dt_1 \quad (1)$$

with  $G(t_2 - t_1)$  the causal Green's functions in the sense that  $t_2 > t_1$  and that plays the role as mathematical mechanism supporting the transition from the state **1** to **2**. In addition  $\mathbf{H}(t_1)$  the input function. Although in its original formulation, the physical propagation contains dependence on the space-time, at a first instance one can test this integration as a mathematical rule that engages the time evolution of current infections in large cities. Under the assumption that it is actually the tool that dictates the strain spread then any variation of kernel might be advantageous as to manage the rate of infections. Thus, under the scenario of Eq.(1) is applicable to the ongoing problem intercontinental infection, then the human intervention for alleviating the outbreak by Corona virus can be modeled through the kernel's free parameters. Once the problem of spread and infection is modeled through the propagator theory, this work has opted by the philosophy of Machine Learning in order to translate the *language* of dataset in terms of the view of Tom Mitchell [3] that states that all system can be universally described by actions, (i) task, (ii) performance, and (iii) experience. In this manner one can use this methodology to extract information from any statistical dataset, such as the ones recently have been taken due to the Covid-19 pandemic. The robustness of Machine Learning can also be used to carry out comparisons with previous pandemics such as the 2009 AH1N1[4] in order to find similarities or discrepancies as to the employed schemes that have been applied to optimize the actions taken by the public health systems. Although in principle one can claim that both AH1N1 and Covid-19 might no be associated each other from any angle of analysis, from the applied methodology in this paper, a noteworthy association between AH1N1 and Covid-19 suggests a possible link between the rate of infections and the public health policies that would determine the success about the management of a city or country in periods of crisis created by pandemic. In second section, the theoretical proposal based in the implementation of Green's functions and the possibility of a kind of entropy is presented. Here, the Mitchell's criteria are introduced in a mathematical manner. In third section, once the theory is build, then the applications of it projected onto the AH1N1 2009 Peruvian season and subsequently in the current 2020 Covid-19 Peruvian data is done. Therefore, the Machine Learning interpretation is done. Finally the conclusion of paper is presented.



## II. THE THEORETICAL PROPOSAL

### A. The Concept of Propagator and Green Function

In physics, the propagation between two space-time points is dictated by evolution operators that entirely depends on the dynamics and physical observables of system [5][6]. **Therefore the action of propagation must have a tangible cause, any action that produce changes to the system.** For example Eq. (1) can be extended from the time  $t_1$  to  $t_3$  as:

$$\mathcal{P}_{1 \rightarrow 3}(t_3) = \int G(t_3 - t_2)G(t_2 - t_1)\mathbf{H}(t_2)\mathbf{H}(t_1)dt_2dt_1 \quad (2)$$

by which the Green's functions  $G(t_3 - t_2)$  and  $G(t_2 - t_1)$  make possible the time propagation along the times from  $t_1$  to  $t_3$  passing through  $t_2$  by which in that time it was caused the last propagation. From this one can generalize for a large number  $L$  of propagations as written below:

$$\mathcal{P}_{1 \rightarrow L}(t_L) = \prod_{\ell=1}^L \int G(t_{\ell+1} - t_\ell)\mathbf{H}(t_\ell)dt_\ell. \quad (3)$$

While  $\mathcal{P}_{1 \rightarrow L}(t_L)$  encloses a chain of time propagation, it is perceived as a probability of a system undergoing a transition between the times  $t_1$  to  $t_L$ . Indeed one can assign to  $\mathbf{H}(t_\ell)$  the role of input function that is convoluted with the propagators. Actually, one has  $L - 2$  input function. It should be noted that the case of  $L = 3$  gives Eq.(2).

Consider for example a Gaussian profile that models the time propagation and its respective input function depending on the constant  $\tau_\ell$ , so that one can write down that:

$$\mathcal{P}_{1 \rightarrow L}(t_L) = \prod_{\ell=1}^L \mathbf{H}(\tau_\ell) \int_0^\infty \text{Exp} \left[ - \left( \frac{t_{\ell+1} - t_\ell}{\tau_\ell} \right)^2 \right] dt_\ell = \prod_{\ell=1}^L \mathbf{H}(\tau_\ell) \sqrt{\tau_\ell \pi} \quad (4)$$

where the change  $u = t_{\ell+1} - t_\ell$  and  $dt = dt_\ell$  was used. With the definition for example the input function can be written as:

$$\mathbf{H}(\tau_\ell) = \frac{h_\ell}{1 + \tau_\ell^2} \quad (5)$$

for the sake of simplicity one opts the assumption that all  $h_\ell = h$  and  $\tau = \tau_\ell$  have same value, that also means that the interactions of system have not effect along a complete cycle of interactions, so that the system has same chance to keep its initial state along the subsequent interactions. Thus one can write down:

$$\mathcal{P}_{1 \rightarrow L}(\tau_L) = \frac{h^L (\tau \pi)^{L/2}}{(1 + \tau^2)^L}. \quad (6)$$

This naive result is illustrated in Fig. 1 up to for 4 values of  $L$ . Due to the Lorentzian nature, all peaks are centered in a same value. The amplitudes have been varied with the incorporation of the constant  $(1.5)^L$  that multiplies Eq. (6).

The combination of the Gaussian and Lorentzian profiles can be combined in the sense that both can yield an approximated quantitative description of the evolution of a limited period of pandemic [7][8][9][10]. In this manner with

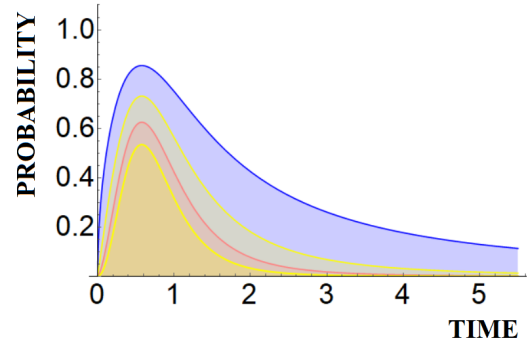


Fig. 1. Illustration of Eq. (6), the probability of propagation versus time. In this simple case all peaks are around of a same time

$\tau^2 \rightarrow 0.1(\tau - 2.5 * L)^2$  and  $h^L \rightarrow 0.85 * (0.2)^L$ , Eq. (6) is rewritten as:

$$\mathcal{P}_{1 \rightarrow L}(\tau_L) = \frac{0.85 * (0.2)^L (\tau \pi)^{L/2}}{(1 + 0.1(\tau - 2.5 * L)^2)^L}, \quad (7)$$

yielding the distributions as shown in Fig. 2. The color black arrows are indicating the decreasing of peaks in time. The fact that all peaks lost their initial value as indicated by the blue arrow in the first peak, is due to any action that in this case is due to the inclusion of term  $0.85 * (0.2)^L$  that is associated to the term  $h^L$  as given in Eq. (5) describes the deterioration of its amplitude along the different times where system experiences interaction.

For instance, one can assume that the curves denote the probability of having a certain number of infections or known as the rate of infection by time units. Thus, in this toy theory: the incorporation of  $0.85 * (0.2)^L$  can be interpreted as the decreasing of rate of infections imposed by the initial conditions of system. In fact, the positions where the arrows have been located would denote that of the times by which a decision has been imposed such as quarantine, curfew or social distance.

Therefore, the model yielding peaked distributions has emerged as one that can be seen a methodology to describe for example rate of infections once an outbreak has been confirmed. Thus, it is possible to define the number of infections as the product  $\mathcal{N} = n_0 \mathcal{P}$  with  $n_0$  the initial number of identified infections. So that the task is to reduce this number through concrete actions in according to the available technology that each public health operator manages in the affected countries. In praxis,  $\mathcal{N}$  would depend on a set of free parameters that features the intensity of pandemic such as population, human behavior and capacity to carry out the social rules after lethality of strain is identified.

### B. Pandemic as Entropy

From Eq. (3) the chain of propagators that introduce the concept of risk of pandemic can also be seen as a kind of Shannon's entropy. In fact, consider that

$$\mathbf{I} = \prod_{\ell=1}^L G(t_{\ell+1} - t_\ell)\mathbf{H}(t_\ell) = \prod_{\ell=1}^L G(\Delta t)\mathbf{H}(t_\ell) \quad (8)$$

$$= G^L(\Delta t)\mathbf{H}^L(t_\ell) = [G(\Delta t)\mathbf{H}(t_\ell)]^L \quad (9)$$

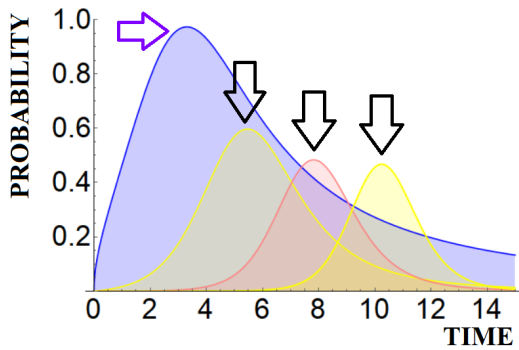


Fig. 2. Illustration of Eq. (7), Probability of Propagation. Arrows are indicating the times where actions have been established to minimize the Probabilities.

under the assumption that the time differences are same and equals to  $\Delta t$  for all the propagation in the sense that  $G(t_{L+1} - t_L) = G(t_L - t_{L-1}) = G(t_3 - t_2) \dots = G(t_2 - t_1)$ , as well as  $\mathbf{H}(t_1) = \mathbf{H}(t_2) = \mathbf{H}(t_3) = \dots = \mathbf{H}(t_{L-1}) = \mathbf{H}(t_L)$  (that can be perceived as a fast propagation of strain once the outbreak has initialized). In this manner Eq.(9) can be rewritten as:

$$\mathcal{P}_{1 \rightarrow L}(t_L) = \int [G(\Delta t)\mathbf{H}(t_\ell)]^L dt_\ell = \quad (10)$$

$$G^L(\Delta t) \int [\mathbf{H}(t_\ell)]^L dt_\ell \quad (11)$$

and by applying the logarithm, then the Shannon's entropy is given by:

$$\mathbf{S} = \text{Log} \mathcal{P} = L \text{Log} G(\Delta t) + \text{Log} \left[ \int [\mathbf{H}(t_\ell)]^\ell dt_\ell \right] \quad (12)$$

indicating that the entropy is only contributed by the propagators whereas the logarithm of integration over the input functions turns out to be  $\delta \mathbf{S} = \text{Log} \left[ \int [\mathbf{H}(t_\ell)]^\ell dt_\ell \right]$  the error of entropy. In the side of Epidemiology, Eq. (12) can also be interpreted as the disorder of the transportation mechanism of strain [11] that is dictated by nonlinearities [12] that leads to a kind of anarchy of system that actually would exhibit any city or place that faces the arrival of a virus causing the social [13] and economic disorder to some extent [14]. In order to go through Eq. (12) the input function is assumed to be a polynomial distribution so that a power series is applied, thus one gets with  $t_\ell \rightarrow u$ , the entropy error is given by:

$$\delta \mathbf{S} = \int [\mathbf{H}(u)]^L du = \int \left[ \sum_m^M C_m \frac{u^m}{m!} \right]^L du \approx \sum_m^M \left( \frac{C_m}{m!} \right)^L \int_0^u u^{mL} du \approx \sum_m^M \left( \frac{C_m}{m!} \right)^L \frac{u^{mL+1}}{mL+1}. \quad (13)$$

From this one can verify if this error is also an entropy. The Shannon's entropy associated to this is then written down as:

$$\text{Log}(\delta \mathbf{S}) \approx \sum_m^M L \left[ \text{Log} \left( \frac{C_m}{m!} \right) \right] + \sum_m^M [(mL+1)\text{Log}(u) - \text{Log}(mL+1)], \quad (14)$$

exhibiting that the two first terms of right side of Eq. (14) follows the structure of a Shannon's entropy. It should be noted that the term  $\left( \frac{C_m}{m!} \right)$  can be perceived as a kind of probability. For large values of integer  $m$  this probability turns out to be null. In this way there is a set of values for  $mL$  and  $u$  that cancels the term  $\sum_m^M L \left[ \text{Log} \left( \frac{C_m}{m!} \right) \right]$  so that one obtains an entropy to be null. Clearly it demands to find the best values of  $C_m$ ,  $M$ ,  $L$ , and  $u$ , by which it would be the task of Machine Learning algorithm.

### C. Theoretical Formulation of Mitchell's Criteria

As mentioned above, the philosophy of Machine Learning action [14], [15] can be resumed in the criteria postulated by Tom Mitchell [16] by which is assumed that the system has a (i) **task** to be done, such task demands to apply a (ii) **performance** that targets to optimize the system's parameters. After of carrying out the performance, the system acquires (iii) **experience** as to the obtained results.

Here one states the main argument of this paper by which is claimed that the probabilistic character of rate of infections as written in Eq.(3) can be translated in terms of the Mitchell's criteria. In fact from Eq. (3), it is feasible to formulate the following algorithm based on the Mitchell's criteria [17]. Consider the number of infections at the  $\ell$ th time as  $\mathcal{N}(t_{L+1}) = n_0 \mathbf{N}(t_{L+1})$ , so that Eq. (3) can be modified to:

$$\mathcal{N}(t_{L+1}) = n_0 \prod_{\ell=1}^L \int G(t_{\ell+1} - t_\ell) \mathbf{N}(t_\ell) dt_\ell. \quad (15)$$

Clearly it was assumed that  $\mathbf{N}(t_{L+1}) = \int G(t_{\ell+1} - t_\ell) \mathbf{N}(t_\ell) dt_\ell$  in where  $\mathbf{N}$  experiences a variation from the time  $t_\ell$  to  $t_{L+1}$ , through  $L$  interactions producing the subsequent propagation.

Thus, the task consists in to reduce the number of infections through  $L$  different periods of pandemic evolution. It is noteworthy that artificial intervention to the outbreak evolution can be applied in order to counteract the progress of pandemic, and therefore to minimize the infections.

In this manner, the minimization of  $\mathcal{N}(t_{L+1})$  is a genuine obligation of system [18][19][20]. To accomplish this, one requires the best strategy that in the picture of Tom Mitchell is known as the performance.

It requires to postulate the best representation of Green's function as seen in Eq. (3). It implies to find the best value for the integer number  $L$ . Once that this number have been determined, for example  $L' < L$  the one proceeds with the integrations. It is suitable to implement the Monte Carlo step that makes the decision of obtain an optimal and reduced number of infections. Thus, in the case that it is not in accordance to the desired number of infections *i.e.*:  $\mathcal{N}(t_{L'}) < \mathcal{N}(t_{L'-1})$  then it is opted that  $L' \rightarrow L' + 1$  to verify that there is a reduced number of infections. Thus, the action is considered as long as  $\mathcal{N}(t_{L'+1}) > \mathcal{N}(t_{L'})$ . In this manner, the process is stopped when it is verified the condition given an integer number  $n$  that verifies  $L' - n < n < L + 1$  then:

$$\frac{\mathcal{N}(t_{L'})}{\mathcal{N}(t_{L+1})} \ll 1, \quad (16)$$

for example if for  $t_{L+1}$  a pandemics yields 100 infections for a period of 10 days, then one expects 1000 for 100 days. Then,

under the usage of Mitchell's criteria one might to obtain a rate of 20 infections for 90 days. Then it implies approximately that  $20/1000 = 0.02$  that is fulfilling the condition of Eq. (16). Of course, for times  $> t_{L+1}$  one might to expect a certain nonlinearity since the rate of infections is partially governed by randomness more than deterministic laws [21].

Thus, once the path have been identified, the one can reconstruct the Green's functions of system, the one that the system has opted to yield a certain number of infections. Clearly, this reconstruction demands to know the involved free parameters and other unknown quantities that could not have been visible at the beginning of a pandemic. A crude estimate to reconstruct the Green's functions is done through the confrontation of data that displays the number of infections per unit of time. Thus one can write below the relations between data and the Green's function of system as:

$$\frac{1}{n_0 \mathbf{N}(t_\ell)} \frac{d^2 \mathcal{N}(t_3)}{dt_1 dt_2} = G(t_3 - t_2) G(t_2 - t_1), \quad (17)$$

where one can apply a fitting to the acquired data that is in essence the left side of this equation. Thus, the reconstruction of the product of propagators would depend entirely on the quality of fitting expressed in terms of  $\chi^2/\text{d.o.f}$ . One should note that after the fitting is done the Green's function can be written as:  $G \approx \sqrt{\frac{1}{n_0 \mathbf{N}(t_\ell)} \frac{d^2 \mathcal{N}(t_3)}{dt_1 dt_2}}$ .

### III. APPLICATIONS

#### A. The Peruvian 2009 AH1N1 Season

The so-called pig-flu strain [22][23][24] had its apparition in Perú [25] along the first week of May being through a people whom have been abroad. The infections started over Lima city being this the main place of strong spreading as seen in Fig. 3. In effect, infections reached its peak on June 20th at Lima city. Due to the outbreak, social measurements were imposed on people in order to block the strain mobility to avoid spreading in vulnerable population. Such social restrictions had an interesting effect as noted at the apparition of a secondary peak on July 2th. Clearly from this date the Lima's infections shown a descent behavior that can be associated to the social regulations. On the other side, data also exhibits for the rest of Provinces a first peak ion July 10th. Clearly one can ask about the why both distributions are not superimposed each other. Clearly data reveals us that the fast up of infections in Provinces is not in phase with Lima city due to human mobility that might be nonlinear. The why Lima city exhibit more infections might be entirely related to the total population. Thus, there is certain probability of a conjunction of external variables that would give the rise to the gap of the peaks of cases between Lima city and provinces.

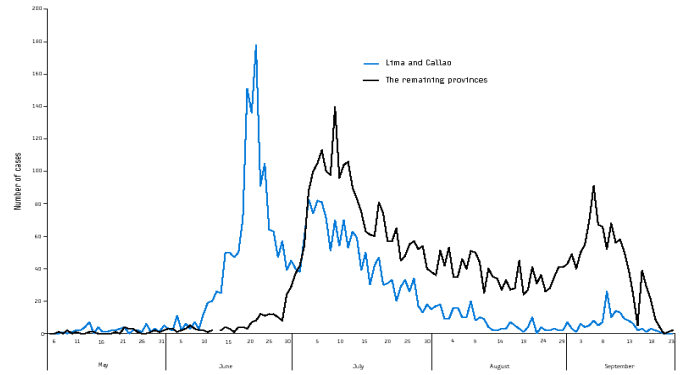


Fig. 3. The Time Evolution of Infections by 2009 Season AH1N1 [25] Outbreak at Lima City (Blue Line) and remaining Provinces (Black Line).

#### B. The Machine Learning Parameters

Eq. (7) is used to validate the Mitchell's parameters on the data of Fig. 3. Thus, it is assumed that Lima data of AH1N1 is a sum of up two different distributions. Therefore the law that models th infections is given  $\mathcal{N} = n\mathcal{P}_{1 \rightarrow 2}$ . To accomplish this, it was applied the change given by  $0.85(0.2)^L \rightarrow 850 \times (0.07)^L$ . The denominator has passed of  $(1 + 0.1(\tau - 2.5 \times L)^2)^L$  to  $(1 + 0.7 \times (x - (4 + 3 \times L))^2)^L$ . Thus, in the scenario of Machine Learning the quantity 850 denotes the expected number for a period of 10 weeks. By using the Mitchell's criteria the task consists in to reduce the first peak [26] that in turn it is equivalent to impose social restrictions that minimize the human contact . The performance is then focused on the different methods that would reduce the infections. It should be noted that during the AH1N1 pandemic in Perú, not any quarantine neither curfew was applied. Instead of that the closing of social activities was done. Mathematically speaking while the task is modeled by a Lorentzian distribution under a dependence on the term  $(x - (4 + 3 \times L))^2$  that exhibits the first two peaks indicating that after any action the second peak becomes reduced in its height, in conjunction to this, one expects the effect of the numerator in Eq.(7) to minimize the infections. In addition it should be noted that the term  $(0.2)^L$  plays a critic role. A slight variation of value  $0.2 + \delta$  with  $\delta$  a small number  $\ll 1$  yields abrupt changes at the morphology of spectrum. Thus for the present study  $(0.07)^L$  governs the behavior of Fig. 3. One can see that it is strongly correlated to the resulting integration of Gauss profiles. In fact the term  $(0.07)^L$  affects directly the value of  $x^{L/2}$  encompassing the role of propagators whose role here is that of minimize the first peak. In this manner, being the task modeled a Lorentzian distribution then one can anticipate available dates by which one expects the decreasing of infections.

In this way the experience is given by the second peak, after a period of decisions by the local public health operators.

Fig. 4 indicates the experience on the 9th week after implementation of Mitchell's criteria between 5th and 8th week. The arrow between the task and experience would denote the performance that is translated in the social restrictions to avoid the strain propagation in people. In this way, the management of AH1N1 pandemic in Lima city might be seen as efficient as well as sustainable to minimize the effects of the strain arrival.

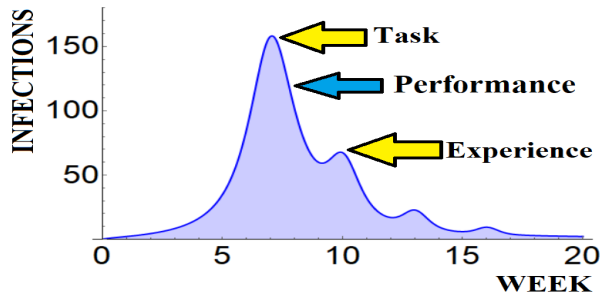


Fig. 4. Reconstructed curve of cases by using the 3 actions in according to the Mitchell's criteria.

### C. The Peruvian Covid-19 Pandemic

The current Covid-19 have assaulted in an unexpected manner the world-wide public health schemes, being to date Perú (date of submitting this paper) [28][29][30] as one of the more affected as to the number of new cases per day. In fact, in Fig. 5, the morphological composition displays up to phases being the first between the 1th and 13th week, and a second phase for the remaining data. While country government has dictated social restriction such as quarantine and curfew, even under this, the number of cases has been increasing in an unstoppable manner as seen at histogram. In fact, despite of the fact that social distancing and face protection were imposed, one can see that the new cases per week have shown a rapid growth of up to a 61% approximately (with respect to the 27686 cases of week 12th) as seen the jump from the 12th to the 13th week. It is noteworthy that it is perceived as the peak of first wave. Beyond of this, data exhibits a morphology that can be understood as the beginning of a second distribution as the consequence of all actions that were imposed before or after 13th week. Under the assumption that first peak has a substantial contribution from Lima city as reported by official data, then is feasible to state that a second distribution is due to cases from provinces. Under this view and by comparing to AH1N1 2009 season data then one can see that human mobility might be the main cause of the formation of second distribution. Thus, one can argue that infections were transported from Lima to provinces through actions of mobility as seen in the opening dates of the travels: July 1th (terrestrial) and July 15th (aerial). One can anticipate in a scenario of Mitchell's criteria that the performance to contain the infections could has been to some extent inaccurate.

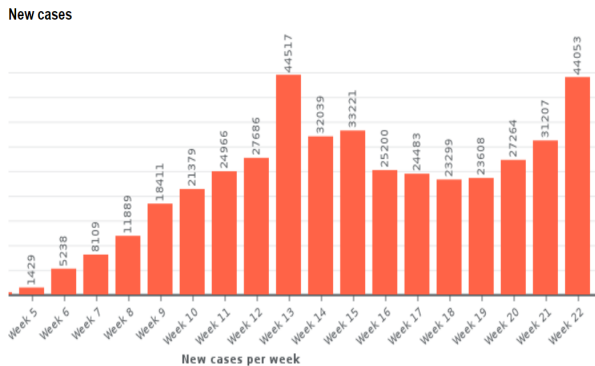


Fig. 5. Histogram showing new cases per week between 5th and 22th week at the ongoing Covid-19 outbreak in Perú [28][29][30].

### D. The Machine Learning Interpretation

Again, the Mitchell's criteria are applied to interpret the histogram of Fig. 5. Because the imminent presence of two well-define distributions having a similarity among them, one can see that the entire system do not exhibits not any flatness as require to claim the end of a first wave. Thus, it is fair to claim that while a possible first wave was ending on Lima city, the formation of a second one essentially due to new cases coming from provinces have been manifesting before the peak at the 13th week. In this way, the first peak is perceived as the task of system to reduce it. However it is clear the the superposition from provinces might add a kind of bias to data. Once the **task** has been identified one can pass to apply a strategy or performance whose target is to decrease the peak on the subsequent weeks. The **performance** as modeled by a Gaussian profile appears to be rather limited as to provoke a fast decreasing of the number of new infections. In fact, in Fig. 6 the Machine Learning reconstruction of Fig. 5 is plotted. The task and performance are indicating the possible dates by which can be matched with the official data. For this end, the equation  $\mathcal{N}_{\text{COVID}} = \sum_{L=1}^2 \frac{n(5/2+0.05*L)}{50+(x-(3+12*L))^2}$  was used, with  $\mathcal{N}$  the new cases per week and  $n$  the expected number of infections. In this manner, the new infections  $n$  becomes a free parameter of system. Although performance is not identified in data, from Fig. 5 one can see that the performance is not visible as a tangible action that has caused variations on the curve. An argument of the why performance cannot be identified on the data is because the superposition of Lima and provinces data that would generate a kind of unrecognizable noise that would affect the national data. Because of this, performance turns out to be a constant and are modeled by a Gaussian profile containing a width that is randomly fixed. It implies that  $\tau_\ell \rightarrow \beta$  a constant. Thus, one has that  $\int_0^\infty \text{Exp} \left[ -\left( \frac{t_{\ell+1}-t_\ell}{\tau_\ell} \right)^2 \right] dt_\ell = \sqrt{\beta\pi}$ . Thus **experience** acquires same morphological shapes from task. This is seen in Fig. 6 where task and experience are modeled by two Lorentzian shapes separated by a gap of 18 weeks approximately.

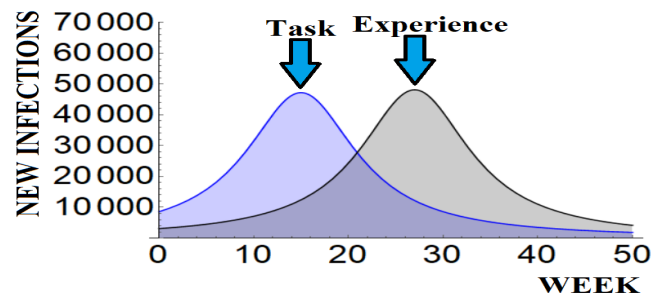


Fig. 6. Reconstructed Curve of Cases by using the Mitchell's Criteria without the Identification of Performance.

## IV. DISCUSSION AND CONCLUSION

The fact that  $\mathcal{N}_{\text{COVID}}$  does not exhibits the Mitchell's performance but encompasses to some extent to Fig. 5, then it is interpreted as follows: Performance was applied before the first peak of 13th week, so that it could has been broken as effect of the end of national quarantine as well as the stopping of curfew at the noon. On the other side, while the apparition

of peaks can also be seen as the resulting outputs after the implementation of imposed actions on the subsequent weeks once the strain was recognized, at the language of Mitchell's criteria, task as a focused fact is not reflected from data. In this manner, inputs Lorentzian distributions were not affected by a constant propagation in contrast to the AH1N1 by which the Mitchell's criteria fits well to data. Indeed a constant performance in terms of Feynman propagator is interpreted as the system has not any scheme (or strategy) to experience variation in time. Thus, one can argue that for the ongoing Covid-19 pandemic in Perú, its translation in terms of Machine Learning could not involve the action of performance, a crucial step to manage the system evolution. While the mathematical probability as given by Eq. (7) was calculated through the usage of a Gaussian profile that models the propagation, then with this one can conclude that the apparition of a second peak is due to a very limited and almost invisible performance. Here, one can ask about the usage of a different propagator distribution. However, it would demand to introduce a set of free parameters that might not be fitted to data, so that it put apart the Mitchell's criteria far from realistic interpretation that must be adjusted to the ongoing acquired data.

#### REFERENCES

- [1] Hannah Stower, Spread of SARS-CoV-2, Nature Medicine 26, 465–465, 9 April 2020.
- [2] Feynman, Richard P, Space-Time Approach to Quantum Electrodynamics, Physical Review 76 (6): 769-789, 1949.
- [3] Tom M. Mitchell, The Discipline of Machine Learning, Machine Learning Department technical report CMU-ML-06-108, Carnegie Mellon University, July 2006.
- [4] Gilberto Gonzalez-Parra, Modeling the epidemic waves of AH1N1/09 influenza around the world, Epidemiology Volume 2, Issue 4, December 2011, Pages 219-226.
- [5] R. P. Feynman, Space-Time Approach to Non-Relativistic Quantum Mechanics, Rev. Mod. Phys. 20, 367 – 1th April 1948.
- [6] John Archibald Wheeler and Richard Phillips Feynman, Classical Electrodynamics in Terms of Direct Interparticle Action, Rev. Mod. Phys. 21, 425 – Published 1 July 1949.
- [7] Huber Nieto-Chaupis, Feynman-Theory-Based Algorithm for an Efficient Detaining of Worldwide Outbreak of AH1N1 Virus, IEEE CHILEAN Conference on Electrical, Electronics Engineering, Information and Communication Technologies CHILECON, 2019.
- [8] Huber Nieto-Chaupis, Face To Face with Next Flu Pandemic with a Wiener-Series-Based Machine Learning: Fast Decisions to Tackle Rapid Spread, IEEE 9th Annual Computing and Communication Workshop and Conference, CCWC, 2019.
- [9] Hannah Stower, Functionally assessing coronavirus entry, Nature Medicine 26 , 465–465, 9 April 2020.
- [10] Peter Horby, Improving preparedness for the next flu pandemic, Nature Microbiology volume 3, pages 848–850 (2018).
- [11] Peng Zhou and *et.al*, A pneumonia outbreak associated with a new coronavirus of probable bat origin, Nature 579, 270–273, 3 February 2020.
- [12] Fan Wu and *et.al*, A new coronavirus associated with human respiratory disease in China, Nature 579, 265–269, 3 February 2020.
- [13] John Tregoning, Coronavirus diaries: all the things we do not do, Nature Comments and Opinion, 8 May 2020.
- [14] David Westgarth, Coronavirus pandemic: A time for reflection, BDJ In Practice 33 , 4–4, 4 May 2020.
- [15] T. Tanaka; T. M. Mitchell, Embedding learning in a general frame-based architecture, IEEE International Workshop on Tools for Artificial Intelligence, Year: 1989 Pages: 77 - 84.
- [16] Tom M. Mitchell, Machine Learning, McGraw Hill, 1997.
- [17] Tom M. Mitchell and *et.al*, Never-Ending Learning, Communications of the ACM, 61(5), pp. 103-115, May 2018.
- [18] Tom M. Mitchell and Michael I. Jordan, Machine Learning: Trends, Perspectives, and Prospects, Science 349, 255, 2015.
- [19] Tom M. Mitchell, Mining Our Reality, Perspective, Science, 326, December 2009.
- [20] Tom M. Mitchell, Machine Learning and Data Mining, Communications of the ACM, Vol. 42, No. 11, November 1999.
- [21] Tom M. Mitchell, Does Machine Learning Really Work?, Invited paper, AI Magazine, Vol. 18, Number 3, AAAI Press, Fall p.11-20, 1997.
- [22] L. A. Angelova, Long Term Immunity to Influenza A(H1N1) in Humans, Ann. Vir. (Inst. Pasteur) 1982, 133, E, 267-272.
- [23] ShakerAl Faress, Identification and characterization of a late AH1N2 human reassortant in France during the 2002–2003 influenza season, Virus Research Volume 132, Issues 1–2, March 2008, Pages 33-41.
- [24] Seth J. Sullivan, 2009 H1N1 Influenza, Mayo Clinic Proceedings Volume 85, Issue 1, January 2010, Pages 64-76.
- [25] Laguna-Torres, *et.al*, Influenza-like illness sentinel surveillance in Peru. PLoS One. 2009;4(7):e6118.
- [26] Natalia Goni and *et.al*, Bayesian coalescent analysis of pandemic H1N1 influenza A virus circulating in the South American region, Virus Research Volume 170, Issues 1–2, December 2012, Pages 91-101.
- [27] Wolfram Mathematica: [www.wolfram.com](http://www.wolfram.com).
- [28] Cases confirmed by Covid-19 coronavirus amount to 375,961 in Peru (Communique N 182) "(in Spanish). 24 July 2020. Retrieved 24 July 2020.
- [29] Minsa: Cases confirmed by Covid-19 coronavirus amount to 407,492 in Peru (Communique N 191) "(in Spanish). Ministry of Health. 30 July 2020. Retrieved 30 July 2020.
- [30] Minsa: Cases confirmed by Covid-19 coronavirus amount to 439 890 in Peru (Communique N 196) ". Ministry of Health (in Spanish). 4 August 2020. Retrieved 4 August 2020.