# Editorial Preface

*From the Desk of Managing Editor...*

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

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

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

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

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

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

**Thank you for Sharing Wisdom!**

# Editorial Board

# CONTENTS

<cn># Transcription

# Analyzing Personality Traits and External Factors for Stem Education Awareness using Machine Learning

Sang C. Suh[1]

Department of Computer Science
Texas A&M University-Commerce
Commerce, TX–USA

Anusha Upadhyaya B.N[2]

Department of Computer Science
Texas A&M University-Commerce
Commerce, TX-USA

Ashwin Nadig N.V[3]

Department of Computer Science
Texas A&M University-Commerce
Commerce, TX-USA

*Abstract*—**The purpose of the paper is to present the personality traits and the factors that influence a student to pursue STEM education using machine learning techniques. STEM courses have high regard because they play a vital role in global technology, inventions and the economy. Educational Data Mining helps us to identify patterns and relationships in a large educational database. On the other hand, Machine Learning facilitates decision making process by enabling learning from the dataset. A survey comprising of an extensive variety of questions regarding STEM education was conducted and the opinions of students from various backgrounds and disciplines were collected. A dataset was generated based on the responses from students. Machine Learning algorithms (one class-SVM and KNN) applied on this dataset emphasizes variety of courses offered, research-oriented learning, problem-solving approach, a good career with high paying job are some of the factors which may influence a student to choose STEM course.**

*Keywords—Educational Data Mining (EDM); Science Technology Engineering Management (STEM); Machine Learning (ML); K-Nearest Neighbor (KNN); One class–Support Vector Machine (one class - SVM)*

## I. INTRODUCTION

### A. Significance of STEM Education

Education is now one of the fundamentals of living. It is no more just a tool to spread knowledge. Every day in this modern world, a new scientific invention or technology is being introduced. Science, Technology, Engineering and Mathematics have integrated into all aspects of life and it has been more accessible than ever before. Students of all areas of study are directly or indirectly interacting with these aspects, no matter what is the field of study there is some technology to store it, share it and enhance its possibility in various dimensions.

As this era is turning out to be an era of automation and machines are becoming more capable of exhibiting their intelligence and skill, current and next generation of students have to compete with machines also. So the focus of the research is on the attitude, knowledge and abilities of the students regarding STEM discipline [1]. As an individual chooses his/her path for a career based on the interests which evolved from his/her childhood. So it is essential to analyze the distinct quality of each person who is in multiple year of study in college and his/her perspective towards STEM education.

The teaching Science, Technology, Engineering and Mathematics also emphasize integrating other areas of study to be a part of them. A person who teaches these subjects should be able to identify research and explain the significance and uniqueness of this integrity among his/her students. Integration of the skills acquired from STEM courses will play a vital role in understanding and re-structuring of various aspects of life [2]. Thus, a student should be able to understand the local politics and know what is happening around the globe and ready to articulate his/her views on it. Inquiry-based learning model can help students to become scientists in their way to explore new information. STEM courses are potent tools for a student to contrive solutions. It also helps in understanding, being creative and innovative in approach to any given challenge. Furthermore, STEM helps to acquire the necessary skills [1].

### B. Machine Learning in Educational Research

The influence of Machine Learning and Artificial Intelligence is already noticeable on the global economy. Thus, it is driving much attention from analysts. ML deals with several algorithms to enhance their performances. It tackles multiple aspects of regression or classification issues in research related to data analytics [2]. *Educational Data Mining* (EDM) deals with extracting necessary information from massive data sets, which are related to students. EDM involves various algorithms like Decision trees, K-Nearest Neighbor, Neural Networks, Support Vector Machines and other Machine Learning algorithms. Analysis, prediction and data-driven problem-solving techniques have shown effective results in forecasting consumer behaviors, fraud detection, intrusion detection and various assessments. The education system also can be one of the significant areas where one can implement data-driven techniques, where it will help to analyze various patterns associated with students [3]. At the same time, statistical methods make it difficult to identify and comprehend.

In this modern world that is driven by data, there is an exponential growth of textual chronicles in the field of education. Every device connected to the internet is an ocean of knowledge. As the online tutoring and skill sharing are becoming more popular each day, this process is producing a massive amount of data daily. Classification of this enormous amount of data is the biggest challenge in data analytics and there have been introduced a variety of methods to handle data volumes. One such method of classification is K-Nearest

Neighbor algorithm which is a supervised learning method [10]. On the other hand, there is an unsupervised learning method called one-class SVM, where a model is trained on the data with only one class. It gleans the features of typical cases and from these features it can predict the cases which are not normal. The main aim of these classifications on educational data is to help the authorities and teachers to improve the performance of a student and assist them in exploring learning and career options with early predictions which are based on previous data.

In this paper, Machine Learning algorithms, one-class SVM and KNN are applied on the data collected from a survey conducted on the students of Texas A&M University-Commerce, where the students from various disciplines gave their responses. The survey consisted of questions regarding their knowledge and opinion about STEM education. Thus, one can able to predict the factors that may influence a student to take up a course associated with STEM or not.

In Section II, we have presented various works which are related to Educational data mining and Machine Learning. Section III describes Machine Learning techniques. Section IV has experimental analysis. Section V has the conclusion and followed by the future work.

## II. RELATED WORK

Educational data mining is the new popular topic among the researchers who are conducting researches on higher education in various institutions. Data mining in the education sector helps us to understand multiple hidden traits among the students, and it helps to analyze the factors influencing a student to take up a particular course and dedicate all their time in learning it. STEM courses are one of the major streams in education, where the STEM courses will help a student to have a successful career. So the awareness regarding this plays a vital role in a students' life.

F. Sciarrone [2] conducted Educational data mining (EDM) with Machine Learning techniques to extract the data. Moreover, he described various models of Learning Analytics with EDM.

Thakar, Mehta, and Manisha [3] have described the techniques used in data mining of Educational data and have given a brief description of the work that has been conducted regarding Educational Data Mining.

Bhardwaj and Pal [4] described the techniques to analyze student performance using decision trees, based on the data extracted from student assessments and their final exam results.

Romero and Ventura [5] surveyed the studies carried out on EDM, which used the computational approaches to analyze the educational data to study educational questions, also it elaborately analyzes the various educational environments and data produced by it.

Meyrick [6] compared traditional college methods and STEM-based program and gave an insight into the positive aspects which influence students to attend STEM Courses.

Popenici and Kerr [7] explored the emergence of artificial intelligence techniques in learning and teaching. It explores the connection between education and emerging technologies by the teaching methods of institutions and how students evolve.

B. Yildirim [8] gave an overview of research studies that are conducted on STEM education and the attributes of students associated with it. Several studies are considered for meta-synthesis method. The results of the study emphasize that the students of STEM discipline are more creative and interested in solving problems.

Amra and Maghari [11] gave a student performance prediction model by applying KNN and Naïve Bayesian on an educational dataset of secondary schools, extracted from the ministry of Gaza strip in 2015, where they compared two of the techniques.

Imdad et al. [12] proposed a method for student result classification based on two traditional algorithms KNN and artificial neural network using the data from the Pakistan education board.

Manevitz and Yousef [13] gave a detailed description of different version of SVMs for one class classification with the context information retrieval.

Kruengkrai and Jaruskulchai [14] presented a comprehensive approach of relevant sentence extraction using only positive samples for training. Here they have applied methodology of support vector machines for one class classification. The fundamental goal of one class SVM is to modify data into feature space compatible with Kernel and then discrete them from original with the highest margin.

Ciolacu et al. [16] presented their case study with the primary goal of predicting the final score of students before attending the exams. Authors proposed an early recognition system with actual data extracted from an embedded learning curriculum with the personalized test for each before the start of the semester.

To analyze the personality traits of the students regarding the knowledge about STEM education, a survey was conducted at the Texas A&M University-Commerce, where students were given a set of statements and questionnaire to respond. The poll was designed to explore the knowledge of a student regarding STEM and to know according to them what the positive and negative aspects of STEM education are. Moreover, the survey helped in understanding the perspective of a student regarding STEM education along with various external factors such as the latest technological trends, world economy, etc.

## III. MACHINE LEARNING TECHNIQUES

Machine Learning techniques are categorized into two types, supervised learning and unsupervised learning. In supervised learning, the outcome can be predicted using the previous input and output data, whereas in unsupervised learning, the algorithm recognizes the hidden patterns or the internal structures of the input data. In our research one

algorithm in supervised learning (KNN) and one algorithm in unsupervised learning (One-class SVM) were considered.

### A. One-Class Support Vector Machine (One-Class SVM)

Schölkopf et al. [9] presented a support vector method for novelty detection. The data points are separated from the origin and distance is maximized from this hyperplane to the origin. Results were shown as in binary function where regions are captured in the input space with the probability density if data is present. So the function yields +1 for the small region and -1 for the rest of it.

$$\min_{w,\xi_i,\rho} \frac{1}{2}\|w\|^2 + \frac{1}{vn}\sum_{i=1}^{n}\xi_i - \rho \tag{1}$$

Subject to:

$$\left(w \cdot \phi(x_i)\right) \geq \rho - \xi_i \quad \text{for all } i=1..., n \ \xi \geq 0 \quad \text{for all } i=1., n$$

In equation (1), parameter $v(nu)$ characterizes the solution. It decides the upper bound on the fraction of outliners and the number of training examples is the lower bound used as support vectors.

By a kernel function for dot product calculations, the final decision function is given in equation (2):

$$f(x) = \text{sgn}\left(\left(w \cdot \phi(x_i)\right) - \rho\right) = \text{sgn}(\sum_{i=1}^{n}\alpha_i K(x, x_i) - \rho) \tag{2}$$

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

K- Nearest Neighbor (KNN) is one of the supervised machine learning techniques. It is a non-parametric lazy learning algorithm. The main aim of the KNN algorithm is to utilize a database wherein which the data points are divided into various classes to predict the classification of a new sample point.

One of the challenging aspects regarding the KNN algorithm is to finding proper value of k [10], For instance, if k equal to 1, then it is defined as the nearest neighbor algorithm. It is simple and straight forward to implement, where it requires only two parameters, stores all the given cases and classifies new cases depending on the similarity measure. In this research, Minkowski distance metric is given by equation (3).

$$\left(\sum_{i=1}^{k}\left(|x_i - y_i|\right)^q\right)^{\frac{1}{q}} \tag{3}$$

### IV. Experimental Analysis

As described in the previous section, the survey responses are given to identify various perspectives of students regarding STEM. The survey was successful in identifying some of the key factors which influenced them to pursue STEM. Those key factors include the attitude towards STEM, big five personality traits of an individual [15] and the external factors like career opportunities and professional growth by studying STEM.

The responses given by the students were recorded and used to create a dataset. Further, applying the Machine Learning techniques on the dataset, one will be able to predict the factors influencing a student to take up STEM or not.

To apply the two machine learning techniques on this dataset, the following parameters were considered. For one-class SVM, MinMaxScaler is used for feature scaling and radial basis function kernel (RBF) and the $v(nu)$ parameter value is 0.025. Parameters used in KNN algorithm is, StandardScaler is used for feature scaling, Minkowski distant metric for measuring distance and the number of neighbors (k) is 1. Results of the experiment are given in Table I.

TABLE I. Accuracy of Classification

| Algorithm | Precision | Recall | Accuracy | F1 Score |
|---|---|---|---|---|
| One-Class SVM | 98.54% | 98.5% | 98.5% | 98.49% |
| KNN | 87.34% | 87.25% | 87.25% | 87.08% |

### V. Conclusion

The paper has examined the accuracy of features, which are considered to be the influencing factors on a student to choose STEM courses. The analysis of the data using machine learning highlighted the influence of certain personality traits on students to opt STEM courses, for instance, a student who is confident in science is likely to take up math for study. Likewise, the person who is interested in a particular area of studies such as robotics, engineering, psychology, statistics and the person who has ability and passion on technical problem solving is more likely to choose STEM rather than any other course. There are other external factors such as the variety of courses offered, research-oriented learning, extensive area of study and most importantly the opportunities to explore and a good career with high paying jobs have elevated impact on choosing STEM courses for their education.

### VI. Future Work

The survey is limited to a particular university, and it reflected the attitude of a certain set of students. In the future, this process can be conducted on a large scale to determine factors influencing current and future students with the involvement of students from various state and educational institutions. This will help to understand the educational and career perspectives of students on national and international levels.

### References

[1] Q. C. Pham, R. Madhavan, R. Chatila, L. Righetti, and W. Smart, "The Impact of Robotics and Automation on Working Conditions and Employment [Ethical, Legal, and Societal Issues]," IEEE Robotics & Automation Magazine, vol. 25, no. 2, pp. 126–128, June 2018.

[2] F. Sciarrone, "Machine Learning and Learning Analytics: Integrating Data with Learning," in 2018 17th International Conference on Information Technology Based Higher Education and Training (ITHET), April 2018.

[3] P. Thakar, A. Mehta, and Manisha, "Performance Analysis and Prediction in Educational Data Mining: A Research Travelogue," International Journal of Computer Applications, vol. 110, no. No. 15, January 2015.

[4] B. K. Baradwaj and S. Pal, "Mining Educational Data to Analyze Students' Performance" (IJACSA) International Journal of Advanced Computer Science and Applications, vol. 2, no. 6, 2011.

[5] Cobal Romero and S. an Ventura, "Educational Data Mining: A Review of the State of the Art," IEEE Transactions On Systems Man And Cybernetics, vol. 40, no. 6, pp. 601–618, NOV 2010.

[6] K. M. Meyrick, "How STEM Education Improves Student Learning," Meridian K-12 School Computer Technologies Journal, vol. 14, no. 1, 2011.

[7]    S. A. D. Popenici and S. Kerr, "Exploring the impact of artificial intelligence on teaching and learning in higher education," in Research and Practice in Technology Enhanced Learning, 2017.

[8]    B. YILDIRIM, "An Analyses and Meta-Synthesis of Research on STEM Education." Journal of Education and Practice, vol. Vol.7, no. No.34, 2016.

[9]    B. Schölkopf, R. C. Williamson, A. J. Smola, J. Shawe-Taylo, and J. C. Platt, "Support Vector Method for Novelty Detection," in Advances in Neural Information Processing Systems '12.

[10]   Moldagulova and R. B. Sulaiman, "Using KNN algorithm for classification of textual documents," in 2017 8th International Conference on Information Technology (ICIT). IEEE May 2017.

[11]   I. A. A. Amra, A. Y. A. Maghari "Students performance prediction using KNN and Naïve Bayesian" ICIT 2017 — 8th Int. Conf. Inf. Technol. Proc., pp. 909-913 2017.

[12]   Imdad, Ulfat, et al. "Classification of students results using KNN and ANN." 2017 13th International Conference on Emerging Technologies (ICET). IEEE, 2017.

[13]   Manevitz, Larry M., and Malik Yousef. "One-class SVMs for document classification." Journal of Machine Learning research 2. Dec (2001): 139-154.

[14]   Kruengkrai, Canasai, and Chuleerat Jaruskulchai. "Using one-class SVMs for relevant sentence extraction." International Symposium on Communications and Information Technologies. 2003.

[15]   Ciolacu, Monica, et al. "Education 4.0-Artificial Intelligence Assisted Higher Education: Early recognition System with Machine Learning to support Students' Success." 2018 IEEE 24th International Symposium for Design and Technology in Electronic Packaging(SIITME). IEEE, 2018.

[16]   L. M. P. Zillig, S. H. Hemenover, and R. A. Dienstbier, "What Do We Assess When We Assess a Big 5Trait?: A Content Analysis of the Affective, Behavioral, and Cognitive Processes Represented in Big 5 Personality Inventories," Personality and Social Psychology Bulletin, vol. 28, no. 6, pp. 847–858, June 2002.

# High-Speed FPGA-based of the Particle Swarm Optimization using HLS Tool

Ali Al Bataineh[1], Devinder Kaur[3]

Department of Electrical Engineering and Computer Science
University of Toledo Ohio
USA

Amin Jarrah[2]

Department of Computer Engineering Hijjawi
Faculty for Engineering Technology
Yarmouk University, Irbid, Jordan

*Abstract*—The Particle Swarm Optimization (PSO) is a heuristic search method inspired by different biological populations on their swarming or collaborative behavior. This novel work has implemented PSO for the Travelling Salesman Problem (TSP) in high-level synthesis to reduce the computational time latency. The high-level synthesis design generates an estimation of the hardware resources needed to implement the PSO algorithm for TSP on FPGA. The targeted FPGA of this algorithm is the Xilinx Zynq family. The algorithm has been implemented for getting the best route between 5 given cities with given distances. The research has used 7 number of particles for a different number of iterations for generating the best route between those 5 cities. The overall latency has been reduced due to the applied optimization techniques. This paper also implemented and parallelized the same algorithm on CPU Intel I7 Processor; the result shows the FPGA implementation gives better results than CPU on the comparison of performance.

*Keywords—FPGA; High Level Synthesis; Particle Swarm Optimization; Travelling Salesman Problem (TSP)*

## I. INTRODUCTION

Particle Swarm Optimization (PSO) is an algorithm which is been adapted from the unpredictable flight of the bird's flock. This algorithm was proposed by Eberhart and Kennedy in 1995. This algorithm is also called as population-based stochastic optimization technique [1]. It uses a model inspired from flying birds in the flock or flock of fish. In the flock, bird or fish (called particle) will search and identify the whole space guided by both its previous best position (pbest) and the best position of the swarm (gbest) or global best position [2]. PSO is applied for multiple fields including scheduling applications, for finding best routes or planning routes and the optimization problems [3-4].

Qang et al. [3] have implemented PSO for job scheduling application. They encoded each particle with a natural number vector and have developed an own approach to move particles in the solution space. They also compared the genetic algorithm (GA) with the PSO for job scheduling application and they found that PSO is very competitive with the GA.

The PSO algorithm with simulated annealing is implemented for optimization of the TSP problem is done by fang et al. [4]. This implementation uses simulated annealing (SA) method for slow down the degeneration of the swarm and increase the swarm diversity. They compared the PSO with SA, basic Genetic Algorithm (GA) and two other algorithms for solving TSP problem in which the PSO with SA gives the superior results than the other methods.

The embedded method of PSO and the SA gives the faster solution than the PSO method in the small and medium-size problem. The SA algorithm is capable to search on a subspace of the whole search space by means of individual particle that result in faster solution and accurate [5].

Hybrid PSO with SA gives the better performance than the adaptive particle swarm optimization and the genetic chaos optimization algorithm [6]. The combination of PSO and SA can narrow down and speed up the field of the search. This strength of the PSO with SA can also optimize the TSP problem with the better search result and speed up.

Hassan et al. [7] have done the comparison of Genetic Algorithm (GA) and the PSO in terms of its effectiveness for finding the global optimal solution and computational efficiency. This research has done the comparison by implementing statistical analysis and formal hypothesis testing.

The goal of this research is to develop and optimize the Travelling Salesman Problem (TSP) on FPGA using High Level Synthesis. Different optimization techniques will be applied such as loop unrolling, loop pipelining, dataflow, loop merging and others. This will help in finding the best route in a high speed.

## II. TRAVELLING SALESMAN PROBLEM

Travelling Salesman Problem (TSP) is a problem of finding the best route for traveling between multiple cities. In the TSP, one salesman wants to visit n cities, the objective of TSP is to identify the shortest Hamilton cycle through which the salesman can visit each city only once and finally return to the starting position or city. The TSP problem is solved using different algorithms as Ant Colony Optimization, Genetic Algorithms, Neural Network, and others [1].

For solving the TSP, the Ant System (AS) [8] and the Particle Swarm Optimization (PSO) [9] is the preferred method due to its optimized solution for the TSP problem. The first implementation of PSO for solving TSP is done by Maurice Clerc in 2000 [9]. At that implementation, results show that PSO was feasible but not very efficient for solving the TSP PSO and AS are implemented for TSP problem for the surveillance mission by Barry R. Secrest [10]. The work is on the planning the best route for the surveillance mission with

different types of aircraft. This implementation is targeted for the Mission Route Planning (MRP) for the Unmanned Air Vehicle (UAV) [10].

General description of TSP can be done as: particles have to identify the shortest path that covers all cities along. Let G= (V;E) be a graph where V is a set of vertices and E is a set of edges. Let C=(cij), which is the distance or cost matrix associated with E. The particles need to identify the minimum cost path or Hamilton cycle between the cities [11].

### III. PARTICLE SWARM OPTIMIZATION

PSO algorithm is inspired from flocks of birds, schools of fish and herds of animals to adapt their environment, find rich source of food and secure themselves from predators by the information sharing approach. Therefore, the PSO algorithm mimics the social behavior of natural organisms, which consists of action of individual member and the effect of other individuals within the group [7]. Each particle in PSO is considered [12],

- to have specific position and a velocity;

- to knows its own position and the value associates with it;

- to knows the best position it has ever achieved, and the value associated with it; and

- knows its neighbors, their best positions and their values.

PSO algorithm gives high performance for different search and pathfinding problems. Therefore, it has been implemented for solving and optimized a wide range of problems. For the optimization on a solution with PSO, the computation cost and the precision are considered as the main variable.

PSO algorithm use Eq. (1) and (2) to calculate the new velocity and position at each iteration:

$$v_i(t+1) = w\ v_i(t) + c_1(p_i(t) - x_i(t)) + c_2\ (g(t) - x_i(t)) \qquad (1)$$

where:

w: real value coefficient (inertia).

$c_1$, c2: real value coefficients (the personal influence and the global influence factors).

$v_i(t)$: current velocity of particle i.

$v_i(t+1)$: next velocity of particle i.

$x_i(t)$: current position of particle i.

$p_i(t)$: personal best known position of particle i.

$g(t)$: global best known position of the whole swarm.

$$x_i\ (t+1) = x_i\ (t) + v_i(t+1) \qquad (2)$$

where:

$x_i(t)$ : current position of particle i.

$x_i(t+1)$ : next position of particle i.

$v_i(t+1)$ : next velocity of particle i.

The pseudocode of PSO is showed as follows [2],

```
Initializing the whole swarm randomly
        for(i = 0; i < swarm size; i + +)
Evaluate f(xᵢ)
    while(termination criteria is not atisfied){
for(i = 0; i < swarm size; i + +){
        if(f(xᵢ) > f(pbestᵢ)) pbestᵢ = xᵢ;
        if(f(xᵢ) > f(gbestᵢ)) pgestᵢ = xᵢ;
Update (xᵢ, vᵢ)
Evaluate f(xᵢ)}}
```

*where* $f(x_i)$ is the fitness function for estimating the quality of the solution. *pbest* is the particle local best position and the *gbest* is the global best position of the entire flock. The update step is performed by Eq. 3.

$$v_{id}^{k+1} = w * v_{id}^k + c_1 rand(\ ) * \left(pbest_{id}^k - x_{id}^k\right) + $$
$$c_2 rand(\ ) * \left(gbest_d^k - x_{id}^k\right) \qquad (3)$$

where $i$ is the number of particle and $d$ is the number of dimensions, $w$ is inertia weight and it decide how much the pre-velocity will affect the new one. *c1* and *c2* are constant values, which are also called as learning factors. These constant values decide the degree of affection of *pbest* and *gbest*. $rand(\ )$ denotes to a random number between 0 and 1 [2].

The particles fly towards a new position based on Eq. (4). PSO algorithm is implanted for a certain number of iteration until the stopping criteria will give the solution that lies in the global best.

$$x_{id}^{k+1} = x_{id}^k + v_{id}^{k+1} \qquad (4)$$

### IV. HIGH LEVEL SYNTHESIS

High Level Synthesis (HLS) is the methodology of implementing algorithm on high-level language and targeting the algorithm on the hardware or called Field Programmable Gate Array (FPGA). High-level synthesis program allows writing the algorithms on the high-level language as C/C++ and OpenCL. HLS tool converts the algorithm from this high-level language (C/C++/OpenCL) to the Hardware Description Language (HDL) level [13].

Xilinx has a high-level synthesis tool called VIVADO HLS, which synthesize and converts algorithm written in C/C++/OpenCL into VHDL/Verilog and System C [13]. VIVADO HLS has a number of built-in functions and libraries for video processing, math, linear algebra, digital signal processing and IP (Intellectual Property) Design [13].

Fig. 1 shows the HLS design flow, with the VIVADO HLS we can write out the algorithm on C, C++ or System C or OpenCL which will converted by HLS into VHDL/Verilog/System C or IP format. The HLS tool can also export the design into other formats also which include System Generator, P-Core or XPS format [13].

Fig. 1. High Level Synthesis Design Flow.

In HLS, there are many optimization strategies for optimizing the latency and resources. One of the main strategies for optimization of latency and resource consumption is pragma directives. These pragma directives instruct the compiler for performing the specific operation while compilation [14]. In this PSO implementation, different optimization techniques are incorporated such as loop pipelining, loop unrolling, dataflow, loop merging and others.

## V. HARDWARE PLATFORM

Our PSO implementation is done for targeting the Xilinx FPGA hardware called ZedBoard which is Zynq 7000 family of FPGA. This ZedBoard has xc7z020clg484-1 FPGA device. The Zynq 7000 architecture consists of Processing System (PS), which is programmable dual-core ARM Cortex A9 Processor and Programmable Logic (PL), which is the 7 series Xilinx FPGA Core with resources as LUT, FIFO, BRAM, DSP and IO's [15]. This Zynq has following resources on the PL section, 53200 logic implementable block called, look up table (LUT), 106400 Flip-Flop (FF), 220 DSP blocks and 280 Block RAM [16].

## VI. PSO ANALYSIS AND OPTIMIZATION

The goal of this research is on solving the simple traveling salesman problem (TSP) with the PSO in High Level Synthesis, HLS allows writing an algorithm on C/C++ or OpenCL language. For the TSP problem, there is predefined number of cities with the predefined distance between those cities. The PSO algorithm is used to solve the TSP problem to find the shortest path between cities so that each city must visited only once. For the calculation of the shortest path between the cities, we have taken the 5 number of cities, which is represented in Fig. 2.



Fig. 2. Traveling Salesman Problem with 5 Cities and Distances.

For the random number generation on HLS Catalin Baetoniu [17] has stated the one of the best methods that can generate the true random number in high speed in Xilinx FPGA. Linear Feedback Shift Register (LFSR) is another method of random number generation that uses the shift register to take input as the function from the previous shift register. LFSR method has least feedback than the counters, so it can also be implemented as the fast counter. Eq. (5) represents the mod-2 polynomial function as the feedback to the input to the LFSR [18]:

$$input\ bit = x^3 + x^2 + 1 \qquad (5)$$

LFSR method is used in our research for generating random number; this random number is used by the particles for selection of cities from the given 5 cities.

The PSO algorithm is used for getting the best and shortest path between cities for the Traveling Salesman Problem (TSP). The PSO algorithm can be represented on flowchart, which is depicted in Fig. 3.

We use pipelining technique in the initialization stage, which will help in concurrent operations inside the loop and so increasing the throughput of executing the algorithm and reducing the latency.

At the beginning, the personal best position and fitness of the particle are the current position and fitness of it. Then, the personal best position of this particle is determined by comparing the previous position with the current position. The global best position depends on the function of fitness we choose. If this function aims to find the minimum value and choosing it to be the best choice, then the minimum value of personal best array will be the global best one. The pseudo code of this calculation is shown as follows:

```
for i = 1 to number_of_particles do
    if fitness [i] ≤ personal_best [i] then
    update personal_best_position [i]
    if personal_best [i] ≤ global_best [i]
    then
    update global_best_position [i]
    repeat
```

Fig. 3. Algorithm of Particle Swarm Optimization.

We use pipelining and loop unrolling which will result in creating multiple independent operations. This technique will reduce the latency and improve the throughput of the algorithm. In addition, we use pipelining technique in this stage.

Position and velocity of the particle will be updated every iteration. The new values will become the present position and present velocity in the next iteration. The pseudo code for this updating is shown as follows:

```
for i = 1 to number_of_particles do
update velocity and position of particle i
until stopping criterion
```

In the global best update, we use pipelining technique in inner loop, which led to increase throughput and decrease execution time. We use pipelining and loop unrolling techniques in this loop, which led to enhance the final results.

The mathematical model of particles motion described in Eq. (1) and (2) forms the velocity equation by sum of three parts. The first part is parallel to the previous velocity and equals to $WVi(t)$; called Inertia component. The second part is parallel to the vector connecting Xi to Pi and equal to $C1(Pi(t)-Xi(t))$; called Cognitive Component. The third part is parallel to the vector connecting Xi to G(t) and equal to $C2(G(t)-Xi(t))$; called Social Component. The three parts of Eq. (1) combined to update velocity vector to new one and this update will cause the particle position to be updated to the new position in the problem search space as in Eq. (2). The main stages in PSO algorithm are problem definition, initialization, and core work of PSO.

Problem Definition: in this stage the optimization problem will be define to solve it by PSO and we define a function of X which return the solved values of the problem, which called the cost value. We also define the numbers of variables that exit in the problem, and the ranges of these variables must be determined. The size of matrix that that will be used of these variables should also be defined, and the parameters of PSO like number of iteration, swarm size and the values of *W*, *C1* and *C2* determine at this stage.

Initialization: In this stage, a group of steps must be taken place to start PSO like create the initialization positions of particles. Then, evaluate them and initialize personal based and global based values of particles, and other initialization tasks needed to run PSO in right ways as shown in Fig. 4.

In this stage the information and data for all particles stored in array of structures. Every swarm particle represented by one structure, all needed field for these particles must be stored in these structures, like positions, velocities, cost values, and personal best position. The particles positions with random values and with zero's velocity values should be initialized. Then, we evaluated the cost values of the particles with its positions and saved it in cost fields. After that, we update the personal best position and save it, the global best initialize firstly to the some value that far from the request solution then it replaced by the best particle personal value.

PSO core work: this stage represents the main loop or work of the PSO algorithm as shown in Fig. 5. For every iteration, the particles new velocities using Eq. (1) is calculated and its new positions using Eq. (2) and update its cost values. A comparison of the current particles position and cost values to its personal best ones that stored in memory is also performed. if the new values better than the personal best values we update the personal best values with the new current values, after that we compare the particle best values with the global best values and replace the global best value with the personal best value if it better than it.



Fig. 4. Initialization of the PSO Algorithm.

Fig. 5. Core Work of the PSO Algorithm.

## VII. RESULTS

We have implemented the PSO algorithm for TSP problem in C language in the VIVADO HLS tool. We also have implemented the same algorithm on the Code:Block Compiler. The design is simulated and tested on the HLS as well as on Code Block environment. On the testing of the PSO for TSP with 5 cities with defined distances is explained above. We have run the test for iteration =10 with the number of particles=7. The best path identified for the TSP is 1-2-5-3-4-1 with the total distance of 52. This best path in terms of distance is shown in following Fig. 6.



Fig. 6. Minimum Cost Path Obtained from the PSO Algorithm, Path 1-2-5-3-4-1.

We have compared the latency and resource utilization of the PSO for TSP with different scenarios. This PSO for TSP on VIVADO HLS algorithm is targeted for the Xilinx Zynq FPGA having the FPGA device of xc7z020clg484-1.

For the optimization of latency and resource, we have implemented the pragma directives on the HLS. We have tested the PSO for TSP algorithm by changing the number of iterations and placing the number of particles fixed as 7.

Table I is the resource utilization table while implementing the PSO for TSP on VIVADO HLS. The resource utilization of the implementation is shown in number and percentage. The implementation consumes 6% of BRAM_18K, 12% of DSP48E, 9% of FF and 35% of LUT of the targeted Zynq device.

Table II shows the comparison of latency with respect to the number of iteration and particles while implementing the PSO for TSP on HLS. This comparison is without the optimization with a pragma. In this table, the number of particles is placed constant and the number of iterations is varied from 10 to 100.

The latency of implementation with the number of iteration -10 and particles=7 is smallest then the other implementation. While increasing the iteration as the usual the latency has increased.

Table III shows the latency of the implementation with the utilization of pragma directives for the optimization. This optimization shows the latency of implementation with a different number of iteration and the constant number of particles. While comparing Table II and III, the optimization methodology reduces the overall latency than the without optimization. While the number of iterations =100 and the number of particles=7, the latency with optimization is 1.8 times less than the without optimization on HLS implementation.

TABLE I. RESOURCE UTILIZATION OF PSO FOR TSP WITH OPTIMIZATION AND FLOAT DATA TYPE AT ITERATION=10 AND NUMBER OF PARTICLES=7

| Resources | Available Resources | Utilization by float data type | Utilization by float datatype (%) |
|---|---|---|---|
| BRAM_18K | 280 | 19 | 6 |
| DSP48E | 220 | 28 | 12 |
| FF | 106400 | 10026 | 9 |
| LUT | 53200 | 18937 | 35 |

TABLE II. COMPARISON OF MAXIMUM LATENCY WITH RESPECT TO THE DIFFERENT NUMBER OF PARTICLES AND ITERATION WITH TARGETED CLOCK=10NS AND WITHOUT OPTIMIZATION

| Iteration No.(i) | Number of Particles (a) | Latency (max.) | Time (in Second) |
|---|---|---|---|
| 10 | 7 | 844124 | 0.00844124 |
| 20 | 7 | 960584 | 0.00960584 |
| 30 | 7 | 1077044 | 0.01077044 |
| 40 | 7 | 1193504 | 0.01193504 |
| 50 | 7 | 1309964 | 0.01309964 |
| 100 | 7 | 1892264 | 0.01892264 |

TABLE III.    COMPARISON OF MAXIMUM LATENCY WITH RESPECT TO THE DIFFERENT NUMBER OF PARTICLES AND ITERATION WITH TARGETED CLOCK=10NS AND WITH OPTIMIZATION

| Iteration No.(i) | Number of Particles (a) | Latency (max.) | Time (in Second) |
|---|---|---|---|
| 10 | 7 | 792924 | 0.00782924 |
| 20 | 7 | 808584 | 0.00808584 |
| 30 | 7 | 834244 | 0.00834244 |
| 40 | 7 | 859904 | 0.00859904 |
| 50 | 7 | 885564 | 0.00885564 |
| 100 | 7 | 1013864 | 0.01013864 |

Table IV is the table for resource utilization of resources with the pragma directives for the optimization of resources and latency. The BRAM and DSP is constant while varying the number of iteration and making fixed the number of particles. While the number of FF and LUT has increased respectively when increasing the number of iterations from 10 to 100. The resource utilization is increased because of while the number of iteration is increased, the number of FF and LUT needed for processing arraySubtraction_float(), arrayAddition_float() and multiplyArrayWithScalar().

Table V shows the latency and the best path identified by the PSO for TSP on the Intel 6700HQ Processor and Code::Block Compiler. This Intel x86 processor has 3.50 GHz of frequency, 4 cores and 16GB RAM with Windows 10 Operating System.

TABLE IV.    RESOURCE UTILIZATION REPORT  WITH RESPECT TO THE NUMBER OF ITERATION AND PARTICLES AND WITH OPTIMIZATION

| Iteration No.(i) | Number of Particles (a) | BRAM_18K | DSP48E | FF | LUT |
|---|---|---|---|---|---|
| 10 | 7 | 19 | 28 | 10026 | 18937 |
| 20 | 7 | 19 | 28 | 10496 | 21190 |
| 30 | 7 | 19 | 28 | 10966 | 23300 |
| 40 | 7 | 19 | 28 | 11436 | 24826 |
| 50 | 7 | 19 | 28 | 11906 | 26296 |
| 100 | 7 | 19 | 28 | 14256 | 33790 |

TABLE V.    RESULT WHILE RUNNING PSO WITH DIFFERENT NUMBER OF ITERATION AND PARTICLES ON X86 PC [ INTEL I7 6700HQ PROCESSOR

| Iteration No.(i) | Number of Particles (a) | Identified best route | Minimum cost-based distance | Total time spent (sec) |
|---|---|---|---|---|
| 10 | 7 | 0-1-4-2-3-0 | 52 | 0.001 |
| 20 | 7 | 0-1-4-2-3-0 | 52 | 0.001 |
| 30 | 7 | 0-1-4-2-3-0 | 52 | 0.002 |
| 40 | 7 | 0-1-4-2-3-0 | 52 | 0.002 |
| 50 | 7 | 0-1-4-2-3-0 | 52 | 0.002 |
| 100 | 7 | 0-1-4-2-3-0 | 52 | 0.002 |

We have exported the HLS implementation of PSO as the IP format to the VIVADO program. Then we have integrated our PSO IP with other necessary blocks for implementing on the Zynq FPGA board. The interconnection of PSO IP with other blocks is done for instructing the PSO IP with the number of iteration and number of particles from the Zynq Processing System. From this design, Zynq PS can accept the instruction of number of iteration and particles from the UART terminal and then the PS configures the information to the PSO IP. The PSO IP run the information of iteration and number of particles and then reply back the result to the PS.

## VIII.    CONCLUSION

In this paper, the Particle Swarm Optimization algorithm has been implemented on the HLS methodology. This work has implemented PSO for traveling salesman problem (TSP) in the C programming language. The number of cities is 5 and the number of iteration and particles are varied. The HLS algorithm also been optimized with the pragma directives. While optimizing the algorithm the number of resources needed has been increased because of the latency optimization, we used the LOOP_FLATTEN, LOOP_UNROLL and PIPELINE pragma directives. Moreover, the VIVADO block design has been implemented and the processor configuration is implemented.

REFERENCES

[1]    L. Diosan and M. Oltean, "Evolving the Structure of the Particle Swarm Optimization Algorithms," in Evolutionary Computation in Combinatorial Optimization-EvoCOP, Lecture Notes in Computer Science-LNCS, Heidelberg, Springer-Verlag, 2006, pp. 25-36.

[2]    W.-h. Zhong, J. Zhang and W.-n. Chen, "A novel discrete particle swarm optimization to solve traveling salesman problem," presented at IEEE Congress on Evolutionary Computation, Singapore, September 25-28, 2007.

[3]    Q. Kang, H. He, H. Wang and C. Jiang, "A Novel Discrete Particle Swarm Optimization Algorithm for Job Scheduling in Grids," at Fourth International Conference on Natural Computation, Jinan, China, October 18-20, 2008.

[4]    L. Fang, P. Chen and S. Liu, "Particle swarm optimization with simulated annealing for TSP," proceeding of International Conference on Artificial Intelligence, Knowledge Engineering and Data Bases, Corfu Island, pp 206-210, vol.6, 2007.

[5]    G.H. Shakuri , K. Shojaee and H. Zahedi, "An Effective Particle Swarm Optimization Algorithm Embedded in SA to solve the Traveling Salesman Problem," proceeding of  Chinese Control and Decision Conference (CCDC 2009), pp 5581-5586, 2009.

[6]    X.-H. Wanf and J.-J. Li, "Hybrid Particle Swarm Optimization with Simulated Annealing," Proceedings of the Third International Conference on Machine Learning and Cybernetics, Shanghai, pp. 2402-2405, vol. 3, 2004

[7]    R. Hassan, B. Cohanim, O. d. Weck and G. Venter, "A Comparision of Particle Swarm Optimization and the Genetic Algorithm," American Institute of Aeronautics and Astronautics, Austin, TX, April 18-21, 2004.

[8]    M. Dorigo and T. Stützle, Ant Colony Optimization, 1st ed. London, U.K: MIT Press, 2004, ch.3, pp.65-81.

[9]    M. Clerc, Particle Swarm Optimization, 1st ed. London, U.K: ISTE, 2006, pp.172.

[10]    B. R. Secrest, "Travelling Salesman Problem for Surveillance Mission using Particle Swarm Optimization", M.S Thesis, Dept. of the air force, Air Force Institute of Technology, Air University, Ohio, 2001.

[11]    X. Yan, C. Zhang, W. Luo, W. Li, W. Chen and H. Liu, "Solve Traveling Salesman Problem Using Particle Swarm Optimization Algorithm," International Journal of Computer Science, vol. 9, no. 6, pp. 264-271, 2012.

[12] E. F. G. Goldbarg, G. R. de Souza and M. C. Goldbarg, "Particle Swarm for the Traveling Salesman Problem," EvoCOP 2006: Evolutionary Computation in Combinatorial Optimization, vol. 3906, pp. 99-110, April, 2006.

[13] Vivado Design Suite (High Level Synthesis)-UG902, Xilinx, Inc., 2017, pp. 5-14.

[14] Vivado HLS Optimization Methodology Guide-UG1270, Xilinx, Inc., 2017, pp. 9-32.

[15] Zynq-7000 All Programmable SoC, Technical Reference Manual (UG585), Xilinx, Inc., 2017, pp.26-40.

[16] ZedBoard hardware user guide, Avnet, Inc.,2012, pp. 3-29.

[17] C. Baetoniu, "High Speed True Random Number Generators in Xilinx FPGAs," Xilinx, Inc., San Jose, CA, 2004.

[18] Mentor Graphics Corporation, High-Level Synthesis-Blue Book, Mentor Graphics Corporation, 2010, pp. 142-143.

# Exploring Factors Associated with Voucher Program for Speech Language Therapy for the Preschoolers of Parents with Communication Disorder using Weighted Random Forests

Haewon Byeon[1]
Department of Speech Language Pathology, Honam University
Gwangju, Republic of Korea

Sulki Cha[2]
Dept. of Rehabilitation
Graduate School, Honam University
Gwangju, Republic of Korea

KyoungYuel Lim*[3]
Department of Speech Language Pathology, Honam University
Gwangju, Republic of Korea

*Abstract*—It is necessary to identify the demand level of consumers and recognize the support target priority based on it in order to provide efficient services with a limited budget. This study provided baseline data for spreading the use of consumer-oriented voucher service by exploring factors associated with the demand of the Voucher Program for Speech Language Therapy for preschool children. This study were analyzed 212 guardians living with children (≤5 years old) who resided in Seoul from Aug 11 to Oct 9, 2015. The outcome variable was defined as the demand (i.e., required and not required) of the Voucher Program for Speech Language Therapy. The results of the developed prediction model were compared with the results of a decision tree based on classification and regression tree (CART). The prediction performance of the developed model was evaluated using a confusion matrix. Among the 212 subjects, 112 (52.8%) responded that the Voucher Program for Speech Language Therapy was necessary. The weighted random forest-based model predicted five variables (i.e., whether preschooler caregiving services were used or not, economic activity after childbirth, the awareness of Seoul's welfare counselor operation, mean monthly living expenses, and whether welfare related information was obtained) as the variables associated with the demand of the Voucher Program for Speech Language Therapy and the accuracy was 72.1%. It is needed to develop systematic policies to expand consumer-oriented language therapy services based on the developed prediction model for the Voucher Program for Speech Language Therapy.

*Keywords*—*Weighted random forests; CART; speech language therapy; prediction model; voucher program*

## I. INTRODUCTION

The number of people with disabilities is increasing in South Korea. Korean Act on Welfare of Persons with Disabilities divides disabilities into 15 types including physical disability and hearing impairment. As of 2017, the population of people with disabilities is estimated as 2,660,000 people, which is an increase of more than 20% compared to 2005 (2,140,000 people) [1]. The most common disability is a physical disability, followed by visual impairment, hearing impairment, low intelligence, and autism [2]. Among them, hearing impairment, low intelligence, and autism affect language development to lead communication problems [3].

When people with communication disorders get married and have a child, the language development of the child is more likely to be delayed even if the child does not have a disability [4]. It is because the communication restriction of the parents with language impairment limits the language development support needed for the child. Therefore, Korean Ministry of Health and Welfare has implemented "the Voucher Program for Speech Language Therapy" for children without disabilities raised by parents with disabilities since 2009 in order to support the successful language development of children and strengthen the competence of families with disabilities [5]. The key components of the program are to provide language rehabilitation services such as language development and aural development to children (<12 years old) without disabilities and having at least one parent with a disability such as hearing impairment or language impairment. It is supported in the form of a voucher. The target amount is 220,000 KRW per month as of 2015, and there may be copayment according to the income of the target family based on the mean nationwide household income.

A voucher means a subsidy that is provided for somebody and it limits purchasing power by allowing a user to select goods or services within a limited range [6]. It is defined as "tied demand side subsidy" [5]. Although the South Korean government actively supports the voucher system, the Voucher Program for Speech Language Therapy, providing vouchers for users who have a desire for the service instead of supporting social welfare service institutes, experiences restrictions on access to services (e.g., reduced project support personnel) due to budgetary deficit [7,8].

It is necessary to identify the demand level of consumers and recognize the support target priority based on it in order to provide efficient services with a limited budget. Nevertheless, previous studies on the Voucher Program for Speech Language Therapy have mainly focused on how to understand voucher, a new policy, and user satisfaction of it [9,10,11,12]. Moreover, their research methods mostly aimed to conduct a factual survey and identify the characteristics of consumers [8]. The

*Corresponding Authors

fact that previous studies have rarely discussed voucher with taking into account the unique characteristics of social welfare services indicates that there is insufficient baseline data, which is needed to run a voucher system effectively, and necessary discussions have not been made yet. As far as we are aware of, there is no study evaluating the factors affecting the demand of the Voucher Program for Speech Language Therapy with considering socio-demographic factors, knowledge of voucher services, methods of obtaining welfare-related information, and education policy satisfaction using data mining techniques.

It is necessary to analyze the demand level of consumers in-depth in order to operate "the Voucher Program for Speech Language Therapy" efficiently and the analysis results will be useful to suggest ways to improve the special education service support in the future. This study provided baseline data for spreading the use of consumer-oriented voucher service by exploring factors associated with the demand of the Voucher Program for Speech Language Therapy for preschool children. The composition of this study is as follows. Chapter 2 will explain the algorithm and model development procedure of weighted random forests as well as study subjects and variable measurement. Chapter 3 will compare the results of the developed prediction model with those of CART model. Lastly, chapter 4 will present conclusions and future research directions.

## II. METHODS AND MATERIALS

### A. Target Subjects

This study selected children (≤5 years old) and guardians residing with them who conducted Seoul Welfare Study, which targeted the local population who resided in Seoul from Aug 11 to Oct 9, 2015. The population of this study was households living in Seoul at the time of census among target households of Statistics Korea's 2010 Population and Housing Census (complete enumeration). Systematic sampling was used. Three thousand household were planned to be sampled so that the maximum limit of error could be approximately 1.8% at the 95% confidence level. Sample households per plot were ten households and 300 sampling plots were selected. The computer-assisted personal interviewing method was used and, for this method, interviewers visited the target households in person and input the responses to the structured questionnaires to the portable computer directly. Fifty-three interviewers were trained from Aug 10 to Aug 13, 2015, prior to the survey. When it was hard to conduct an in-person interview due to speech impairment, hearing impairment, and other difficulties, the spouse of the target was surveyed. Among 3,019 target subjects (households), 2,807 subjects were excluded from the analysis because they did not live with a preschooler (≤5 years old). This study analyzed 212 subjects.

### B. Measurements

The outcome variable was defined as the demand (i.e., required and not required) of the Voucher Program for Speech Language Therapy. When a target subject responded that he or she did not know "the Voucher Program for Speech Language Therapy, an interviewer explained it ("Voucher Program for Speech Language Therapy is a program to provide language development services for children who are raised by parents with hearing impairment or language impairment") and identified the demand.

The explanatory variables included age, number of children, mean monthly living expenses, whether a household was eligible for National Basic Living Security (yes or no), whether a subject use child care services (yes or no), family life satisfaction (dissatisfied, okay, or satisfied), economic activity after giving birth (yes or no), whether a subject knew the welfare counselor service of Seoul (don't know or know), welfare-related information acquisition (none, inquiring to a welfare facility in person, inquiring to a community service center in person, inquiring to a local resident, call center, internet search, and friend/relative/friend), satisfaction of education policy (satisfaction, okay, or dissatisfaction).

## III. ANALYSIS METHODS

### A. Model Development and Evaluation

Data were divided into training data (70%) and test data (30%) in order to develop a model to predict the demand of the Voucher Program for Speech Language Therapy. The development of the prediction model was based on weighted random forests algorithms. The results of the developed prediction model were compared with the results of a decision tree based on classification and regression tree (CART). The prediction performance of the developed model was evaluated using a confusion matrix. Moreover, the importance of variables and major risk factors were compared.

### B. Bagging Tree

Bootstrap aggregating (Bagging) is an ensemble technique that combines multiple bootstrap samples and predicts outcome variable. It is mainly used for models with a small bias and a large variance [13]. Bootstrap means sampling with replacement with having the same sample size for various data. The bth bootstrap sample can be calculated as shown in Equation (1).

$$Z^{(b)} = \left( z_1^{(b)}, ..., z_N^{(b)} \right), where\ z_i^{(b)} = \left( x_i^{(y)}, y_i^{(b)} \right), i = 1, ..., N. \quad (1)$$

Bagging tree uses a decision tree model for each bootstrap sample [14]. The decision tree model has a large variance because a tree has a completely different structure according to the first divided variable (j) and division point (s) [15]. Therefore, it is possible to reduce the variance of an unstable tree model by obtaining the mean after constructing multiple tree models through bagging. Each tree model using bootstrap uses a tree model without pruning to minimize bias.

### C. Random Forests

Random forests are one of the ensemble techniques that makes a tree model using bootstrap samples and predicts by integrating all models [16]. Random forests conduct division by randomly selecting m-dimension, which is smaller than p-dimension, explanatory variables rather than p-dimension explanatory variables. Random forests have the advantage of using out of bag (OOB) samples because they use bootstrap samples [17,18]. The importance variables score can be calculated easily through permutation [19], and the mean square error (MSE) of the OOB sample is calculated using the regression tree model generated from the bootstrap samples.

Fig. 1. Concepts of Weighted Random Forests Algorithm.

### D. Weighted Random Forests

Random forests are one of the ensemble techniques and it conducts model averaging by applying the same weight to each tree model. It is possible that random forests generated by bootstrap have good models and bad models [20]. If model averaging is carried out by applying higher weights to good tree models, it can have better prediction power than the classical random forests giving the same weight to all tree models [21]. The weighted random forests algorithm is developed based on this concept (Fig. 1).

The weighted random forests also use OOB samples. When $b = 1, . . . , B$ and MSE $e(b)$ of OOB sample $(O(b))$ was calculated using a tree model generated by the bth boostsrap sample $(Tr(fb))$, it is assumed that a large $e(b)$ means a bad tree model and a small $e(b)$ means a good tree model. The weighted random forests are defined as a model averaging technique using the weight, which is given to each tree model $(Tr(fb))$ using the calculated $e(b)$. In this model, Akaike weights were used for AIC model selection [22].

## IV. RESULTS

### A. General Characteristics of Subjects According to the Demand of Voucher Program for Speech Language Therapy

Table 1 shows the general characteristics of the subjects according to the demand of the Voucher Program for Speech Language Therapy. Among the 212 subjects, 112 (52.8%) responded that the Voucher Program for Speech Language Therapy was necessary. The results of chi-square test revealed that there a significant ($p<0.05$) difference in welfare-related information acquisition between the subjects who responded that the Voucher Program for Speech Language Therapy was needed and those who responded that it was not needed. The demand for the Voucher Program for Speech Language Therapy was higher in the group (53.7%) that obtained welfare related information from local residents.

TABLE I. GENERAL CHARACTERISTICS OF THE SUBJECTS BASED ON DEMAND OF THE VOUCHER PROGRAM FOR SPEECH LANGUAGE THERAPY, N (%)

| Variables | Demand of the Voucher Program for Speech Language Therapy (n=212) | | p |
|---|---|---|---|
| | Not required (n=100) | Required (n=112) | |
| Age, mean±SD | 38.9±8.1 | 38.1±7.3 | 0.446 |
| Number of children, mean±SD | 1.6±0.8 | 1.5±0.8 | 0.962 |
| Mean monthly living expenses (KRW), mean±SD | 372.3±980.1 | 244.3±103.4 | 0.171 |
| Whether a household was eligible for National Basic Living Security | | | 0.370 |
| Yes | 1 (25.0) | 3 (75.0) | |
| No | 99 (47.6) | 109 (52.4) | |
| Whether a subject use child care services | | | 0.461 |
| Yes | 63 (49.2) | 65 (50.8) | |
| No | 37 (44.0) | 47 (56.0) | |
| Family life satisfaction | | | 0.533 |
| Dissatisfied, | 4 (36.4) | 7 (63.6) | |
| Okay | 45 (51.1) | 43 (48.9) | |
| Satisfied | 51 (45.1) | 62 (54.9) | |
| Economic activity after giving birth | | | 0.317 |
| Yes | 33 (42.9) | 44 (57.1) | |
| No | 67 (50.0) | 67 (50.0) | |
| Whether a subject knew the welfare counselor service of Seoul | | | 0.128 |
| Know | 7 (31.8) | 15 (68.2) | |
| Don't know | 93 (48.9) | 97 (51.1) | |
| Welfare-related information acquisition | | | 0.011 |
| None | 39 (60.9) | 25 (39.1) | |
| Inquiring to a welfare facility in person | 4 (66.7) | 2 (33.3) | |
| Inquiring to a community service center in person | 16 (43.2) | 21 (56.8) | |
| Inquiring to a local resident | 0 (0.0) | 2 (100.0) | |
| Call center | 5 (71.4) | 2 (28.6) | |
| Internet search | 29 (34.1) | 56 (65.9) | |
| Friend/relative/friend | 7 (63.6) | 4 (36.4) | |
| Satisfaction of education policy | | | 0.718 |
| Satisfaction | 39 (44.8) | 48 (55.2) | |
| Okay | 37 (48.1) | 40 (51.9) | |
| Dissatisfaction | 18 (52.9) | 16 (47.1) | |

## B. Results of Weighted Random Forests Model Development

The model to predict the demand of the Voucher Program for Speech Language Therapy was developed through the weighted random forests and the predictive power was compared with the results of CART (Table 2). Weighted random forests had higher classification accuracy than CART in both training and test data. The analysis results of training data showed that the classification accuracy was 72.5% for weighted random forests and 71.2% for CART. For test data, it was 72.1% for weighted random forests and 70.8% for CART.

## C. Comparison of Language-Related Factors by Model

Table 3 shows the results of constructing prediction models based on CART and weighted random forests using 10 explanatory variables for predicting the demand of the Voucher Program for Speech Language Therapy. In this study, the weighted random forests model estimated the key variables using the decrease of the GINI coefficients [23]. In the CART-based model, four variables (i.e., economic activity after childbirth, the awareness of Seoul's welfare counselor operation, mean monthly living expenses, and whether welfare related information was obtained) were predicted as the factors associated with the demand of the Voucher Program for Speech Language Therapy and the accuracy was 70.8%. The weighted random forest-based model predicted five variables (i.e., whether preschooler caregiving services were used or not, economic activity after childbirth, the awareness of Seoul's welfare counselor operation, mean monthly living expenses, and whether welfare related information was obtained) as the variables associated with the demand of the Voucher Program for Speech Language Therapy and the accuracy was 72.1%.

TABLE II.    THE PREDICTION PERFORMANCE OF THE DEVELOPED MODEL

| Data | Model | Accuracy (%) |
|---|---|---|
| Training data | Classification and regression tree | 71.2 |
| | Weighted random forests | 72.5 |
| Test data | Classification and regression tree | 70.8 |
| | Weighted random forests | 72.1 |

TABLE III.    COMPARISON OF LANGUAGE-RELATED FACTORS BY MODEL

| Model | Factors | Characteristics |
|---|---|---|
| Classification and regression tree | 4 | Economic activity after childbirth, the awareness of Seoul's welfare counselor operation, mean monthly living expenses, and whether welfare related information was obtained |
| Weighted random forests | 5 | Whether preschooler caregiving services were used or not, economic activity after childbirth, the awareness of Seoul's welfare counselor operation, mean monthly living expenses, and whether welfare related information was obtained |

## V. DISCUSSION

The establishment and expansion of Voucher Programs for Speech Language Therapy are very important in the aspect that it can enhance the language development of children in the high-risk communication disorder group and the quality of family's life. This study developed a model to predict the demand for language therapy service targeting preschooler without a disability and under parents with language or hearing impairment using the weighted random forest algorithm.

The weighted random forest-based prediction model showed that mean monthly living expenses (reflecting the mean household income) and whether welfare-related information was obtained or not were important factors to predict the demand for language therapy service. On the other hand, it was confirmed that whether receiving the National Basic Living Security or not, reflecting the low-income status of a household, was not a key factor. These results posed two meanings. First, it is necessary to choose the support target based on the actual demand survey rather than prioritizing the low-income class to expand the language therapy service. Additionally, it is needed to increase the budget and alleviate the income criteria for application. Korean Ministry of Health and Welfare allocated 83 billion KRW for development and rehabilitation service projects in 2019, which is 6.7 billion KRW increase from 2018 budget and predicted that it would support 57,094 children [24]. However, selection does not guarantee that all expenses for developmental rehabilitation services will be covered, and the grade will be determined according to the income level and the subsidy varies accordingly. It is because the demand for language therapy services far exceeds the supply. It is mainly because the government did not accurately estimate the demand for the project. In other words, the government tends to make a rough estimate from allocating budget and makes it similar to the previous year's, and the number of target subjects always exceeds the actual demand to exhaust the budget and malfunction the service continuously. Additionally, although each municipal receives application every year, it is already full and the existing applicants are given priority. Therefore, the entry barriers are too big for the new applicants and it requires a counterplan. In the case of Gwangju metropolitan city, the budget for language therapy of it is 4,751,074 USD in 2019, which means that 140 applicants will not be supported when they receive graded payment based on the number of registered people with disabilities [25]. Therefore, actual demand should be surveyed based on mean household income, not on low-income households.

Second, it is necessary to expand the language therapy service institutions and actively advertise the system in order to successfully expand the service. Under the existing unequal service supply system, relatively unfavorable groups should be set as service priority subjects and service institutions should be secured not to exclude them from services. Since it is possible that the access to services may not be guaranteed due to a serious information gap [26], supplying institutions and local governments, service and administrative agencies, should develop active promotion strategies not to make the information gap lead to inequality in social welfare services.

Another finding of this study was that the accuracy and prediction power of weighted random forests was higher than those of CART. It is believed that the weighted random forests had higher accuracy than CART because the former is based on the bagging algorithm that generates diverse decision trees from 500 bootstrap samples [27,28]. CART can be used for both regression and classification problems, and it is widely used because it is simple yet has strong prediction power [29]. The decision tree has a small bias, but the variance of the model is large because the structure of a tree model varies greatly according to the first divided variable. When N-1 divisions are performed for data (n=N), each area contains only one datum and a large tree model generated by it generally has an overfitting problem. Pruning is performed to prevent it, and a suitable size tree model is selected as a final model to conduct a prediction. However, the decision tree is still highly likely to have a overfitting problem.

Hothorn & Lausen (2003) [14] proposed the bagging tree to overcome the overfitting problem in this prediction model. The bagging tree has the advantage of minimizing the bias using each tree model created by bootstrap samples and effectively reducing the variance of the model at the same time [14]. However, bootstrap samples are positively correlated because they are sampled with replacement from the same data. In other words, there is a positive correlation between tree models, so the prediction value of the bagging tree has a higher variance. To complement this, Breiman (2001) [30] proposed random forests that can reduce the correlation for each tree model created by bootstrap samples. Although random forests is also an ensemble technique that predicts result variables by generating many tree models using bootstrap samples, just like the bagging tree. However, the algorithm for building the tree model of it is different from that of the bagging tree. Unlike the general tree model, which starts from dividing all the explanatory variables in the p-dimension in the growth process, the tree models constituting random forests use only randomly selected m ($\leqslant$ p) variables in each division process. The key idea of random forests is to give randomness in the growth process of the tree model. It can reduce the correlation for each tree model through it to make it have better prediction power than the bagging tree [31]. Therefore, the results of this study suggest using weighted random forests for developing a highly accurate prediction model.

## VI. Conclusion

The weighted random forests had higher accuracy than the decision-making trees because they maintained the bias of trees and reduce the variance. Random forests that extract many training datasets, generate trees, and predict a target variable, are suitable to construct a prediction model using data containing many variables such as the Voucher Program for Speech Language Therapy or big data. Particularly, weighted random forests that give higher weights to better-performing tree models would have better prediction power than existing random forests granting the same weight to all tree models. Future studies need to seek ways to enhance the performance of weighted random forest models. It is needed to develop systematic policies to expand consumer-oriented language therapy services based on the developed prediction model for the Voucher Program for Speech Language Therapy.

## References

[1] Ministry of Health & Welfare, Disability registration status, Korea Institute for Health and Social Affairs, Sejong, 2017.

[2] S. J. Kwon, Health care and health status of people with disabilities: policy issues. Health and Welfare Policy Forum, vol. 263, pp. 21-33, 2018.

[3] R. Paul and C. F. Norbury, Language disorders from infancy through adolescence, Elsevier Health Sciences, Oxford, 2012.

[4] D. V. Bishop, and L. Leonard, (Eds.), Speech and language impairments in children: causes, characteristics, intervention and outcome. Psychology press, Hove, 2014.

[5] H. Kang, Achievement of social service voucher. Health·welfare Issue & Focus , vol. 171, pp. 1-8, 2013.

[6] R. J. Daniels, and M. J. Trebilcock, Rethinking the welfare state: Government by voucher, Routledge, Abingdon-on-Thames, 2013.

[7] H. S. Bae, The meaning and assignments of introducing voucher system in social welfare services of Korea. Social Welfare Policy, vol. 31, pp. 319-342, 2007.

[8] Y. I. Jang, A study on the childcare voucher as a new system of state's subsidy for childcare: in the perspective of children's right. The Korea Association of Child Care and Education, vol. 58, pp. 189-217, 2009.

[9] J. A. Son, Rethinking about Publicness of Social Welfare. Journal of Critical Social Welfare, vol. 62, pp. 131-155, 2019.

[10] H. S. Jo, Prospects for 2019 Health and Welfare Policy. Health and Welfare Policy Forum, vol. 267, pp. 2-5, 2018.

[11] S. J. Yu, Qualitative research on paenting stress & homeostatic factors of parents with disabilities in vocational rehabilitation service(voucher) of disabled children in community welfare centers – grounded theory approach. Journal of Regional Studies, vol. 27, no. 1, pp. 123-143, 2019.

[12] S. Kang and J, Moon, Differences between sociodemographic characteristics, instrumental activities of daily living, and healthcare needs in disabled persons with and without language. Journal of The Korean Society of Integrative Medicine, vol. 7, no. 1, pp. 37-45, 2019.

[13] L. Breiman, Bagging predictors. Machine learning, vol. 24, no. 2, pp. 123-140, 1996.

[14] T. Hothorn, and B. Lausen, Bagging tree classifiers for laser scanning images: a data-and simulation-based strategy. Artificial Intelligence in Medicine, vol. 27, no. 1, pp. 65-79, 2003.

[15] T. G. Dietterich, An experimental comparison of three methods for constructing ensembles of decision trees: Bagging, boosting, and randomization. Machine learning, vol. 40, no. 2, pp. 139-157, 2000.

[16] H. Byeon, Developing a model for predicting the speech intelligibility of South Korean children with cochlear implantation using a random forest algorithm. International Journal of Advanced Computer Science & Applications, vol. 9, no. 11, pp. 88-93, 2018.

[17] C. Strobl, A. L. Boulesteix, A, Zeileis, and T. Hothorn, Bias in random forest variable importance measures: illustrations, sources and a solution. BMC Bioinformatics, vol. 8, no. 25, pp. 1-21, 2007.

[18] R. Genuer, J. M. Poggi, and C. Tuleau-Malot, Variable selection using random forests. Pattern Recognition Letters, vol. 31, no. 14, pp. 2225-2236, 2010.

[19] H. Byeon, Developing a model to predict the occurrence of the cardio-cerebrovascular disease for the Korean elderly using the random forests algorithm. International Journal of Advanced Computer Science & Applications, vol. 9, no. 9, pp. 494-499, 2018.

[20] S. J. Winham, R. R. Freimuth, and J. M. Biernacka, A weighted random forests approach to improve predictive performance. Statistical Analysis and Data Mining: The ASA Data Science Journal, vol. 6, no. 6, pp. 496-505, 2013.

[21] A. Booth, E. Gerding, and F. Mcgroarty, Automated trading with performance weighted random forests and seasonality. Expert Systems with Applications, vol. 41, no. 8, pp. 3651-3661, 2014.

[22] E. J. Wagenmakers, and S. Farrell, AIC model selection using akaike weights. Psychonomic Bulletin & Review, vol. 11, no. 1, pp. 192-196, 2004.

[23] H. Pham, and S. Olafsson, On Cesaro averages for weighted trees in the random forest. Journal of Classification, pp. 1-14, 2019.

[24] Ministry of Health & Welfare, Developmental Rehabilitation Service, Korea Institute for Health and Social Affairs, Sejong, 2019.

[25] Gwangju City Hall, Budget scale of Gwangju metropolitan city, Gwangju City Hall, Gwangju, 2019

[26] H. Kang, Issues in Social Services Policy. Health and Welfare Policy Forum, vol. 125, pp. 6-22, 2007.

[27] H. Byeon, A prediction model for mild cognitive impairment using random forests. International Journal of Advanced Computer Science & Applications, vol. 6, no. 12, pp. 8-12, 2015.

[28] H. Byeon, H, Jin, and S. Cho, Development of Parkinson's disease dementia prediction model based on verbal memory, visuospatial memory, and executive function. Journal of Medical Imaging and Health Informatics, vol. 7, no. 7, pp. 1517-1521, 2017.

[29] G. De'ath, and K. E. Fabricius, Classification and regression trees: a powerful yet simple technique for ecological data analysis. Ecology, vol. 81, no. 11, pp. 3178-3192, 2000.

[30] L. Breiman, Random forests. Machine learning, vol. 45, no. 1, pp. 5-32, 2001.

[31] R. O'Brien, and H. Ishwaran, A random forests quantile classifier for class imbalanced data. Pattern Recognition, vol. 90, pp. 232-249, 2019.

# ZigBee Radio Frequency (RF) Performance on Raspberry Pi 3 for Internet of Things (IoT) based Blood Pressure Sensors Monitoring

Puput Dani Prasetyo Adi[1], Akio Kitagawa[2]

Micro Electronics Research Laboratory (MeRL), Kanazawa University, Kanazawa, Ishikawa, Japan

*Abstract*—**Wireless Sensor Network has grown rapidly, e.g. using the Zigbee RF module and combined with the Raspberry Pi 3, a reason at this research is building a Wireless Sensor Network (WSN). this research discusses how sensor nodes work well and how Quality of Service (QoS) from the Sensor node being analyzed and the role of Raspberry Pi 3 as an internet gateway will sending a blood pressure data to the database and displayed in real-time on the internet, from this research it is expected that patients can check the blood pressure from home and don't need to the Hospital even data can be quickly and accurately received by Hospital Officers, doctors, and medical personnel. the purpose of this research is make a prototype to providing a blood pressure (mmHg) real-time data from systolic and diastolic data patient's that determine patients suffering from symptoms of certain diseases, i.e, anemia, symptoms of hypertension and even more chronic diseases. this research discusses how sensor nodes work well and how Quality of Service (QoS) from the Sensor node being analyzed and the role of Raspberry Pi 3 as an internet gateway will sending a blood pressure data to the database and displayed in real-time on the internet. Furthermore, Zigbee has the task of sending Blood pressure (mmHg) data in real-time to the database and then sent to the internet from Zigbee end-device communication to ZigBee coordinator. Zigbee communication at a distance of 5 meters, RSSI simulations show a value of -29 dBm and the experiment shows a value of -40 dBm, at a distance of 100 m, RSSI shows a value of -55 dBm (simulation) and -86 dBm (experiment).**

*Keywords*—*Zigbee; Raspberry Pi 3; IoT; blood pressure*

## I. INTRODUCTION

At present, the Internet of Things (IoT) is growing rapidly, with the development of various platforms and security systems i.e, the ZigBee RF module, WiFi Module, Bluetooth (BLE) module, and Micro-electromechanical system (MEMS). the Internet of Things (IoT) processing data on the internet, i.e, sensor data. IoT allows automation and convenience so that data automatically will be sent to the internet and can be viewed in real-time. i.e, on a Personal Computer Platform, a tablet and a Smartphone has an internet connection. WiFi modules are compatible with Internet of Things Sensors such as Fingerprint, references [1] investigated WiFi Positioning based on fingerprinting and Quality of Service (QoS) and Receive Signal Strength Indicator (RSSI) from sending sensor data using the Positioning of WiFi module. With IoT Technology, human life becomes easier and more flexible, especially in this research, IoT in the health sector i.e, monitoring Blood Pressure using a ZigBee device as

data senders from the Sensor Node to the Coordinator node located on the Raspberry Pi 3 as an internet gateway.

However, IoT devices are faced with the problem of attacks from malicious software [2] so a Security tools method is needed that can protect IoT devices to continue working stably. accordingly, The latest version of Bluetooth devices, BLE is a device with Low Energy that guarantees IoT device Life Time. Bluetooth devices like RN-42 [3] are needed as devices that can communicate Master-Slave, this Bluetooth type can be combined with Internet Gateway with the same Raspberry Pi 3 in the research conducted in this paper. In the study [4] examined about Collaborative, Seamless and Adaptive Sentinel for the Internet of Things (COSMOS) Apps to protect the smart environment from cyber threats. In addition to using the Zigbee RF Module, Bluetooth standards such as the RN-42, or Bluetooth Low Energy (BLE), sending sensor data can use the Global System for Mobile Communication (GSM) device, which includes a Wireless communication system with a wide range of research sensors, namely ultrasonic sensors using GSM SIM900A [5] to smartphone or mobile devices, the similarity of this research is the sensor data transmission system through the Wireless-based System platform. Apart from GSM other data sender technology devices are Radio Frequency Identifier (RFID), 3G, UMTS, WiFi, BLE, infrared and Zigbee which are distinguished by the Object Abstraction Layer.

Accordingly from the data share of dominant IoT application project by 2025, the biggest potential economic impact of sized IoT Applications are Health Care 41%, manufacturing 33%, Electricity 7%, Urban infrastructure 4%, Security 4%, Resource Extraction 4%, Agriculture 4 %, Vehicle 2% and Retail 1%. [6]. therefore, this is the biggest reason this paper was created is to develop the role of IoT in the field of Health Care using Zigbee Performance Pipe which is rarely used as an IoT support device. e.g. on reference [7] examined the application of IoT to optimize the performance of vehicle tracking in cloud servers. Device Healthcare support platforms such as e-health can be combined with IoT and can be analyzed for the level of security, reliability, and architecture they have [8], moreover, e-health is very compatible with many sensors such as SPo2, blood pressure, and pulse sensors. In another study, an approach to Pacemaker Pulse detection [9] was carried out in certain patients who needed Pacemaker, this sensor was able to detect Pacemaker surges up to 12 cm from Pacemaker leads. Dimiter V. Dimitrov, MD, Ph.D. [10] in his research tried to see and

analyze complex data / Big Data on IoT in the field of Healthcare, which consisted of Clinical Data, EHR Data, Molecular omics data, and Wearables data that entered the Repository data, furthermore, a Big Data The warehouse that was then analyzed, happened Dimensionality Reduction, segmentation and ended with Big Data with the number of Billion in the form of Medical Apps developers.

## II. RELATED STUDIES

Giorgio Biagetti, Paolo Crippa, Laura Falaschetti and Claudio Turchetti [11], in this research installing a device called surface electromyographic (sEMG) and an accelerometer on the arm of a person doing fitness then sEMG data and accelerometer sent with Wireless Sensor Network devices, then classified according to level Fusion, consequently, the accuracy of the results of the SEMG data transmission and accelerometer is 82.6% of all types of styles during fitness activities, pada research ini data yang dihasilkan adalah sEMG data and accelerometer, in this research the data produced is sEMG data and accelerometer, in this research can be developed using Kalman Filtering and Learning Algorithms. The position of the sensor is like this research, which is on the wrist that takes the movement of the arteries.

Muhammad Niswar, Amil Ahmad Ilham, Elyas Palantei, Rhiza S. Sadjad, Andani Ahmad, Ansar Suyuti, Indrabayu, Zaenab Muslimin, Tadjuddin Waris, Puput Dani Prasetyo Adi [12], in this research using Zigbee as an RF Module to transmit pulse data patients in different distances, in this research, pulse data is sent from 5 sensor nodes simultaneously to one ZigBee node coordinator. Data received by ZigBee experiences Packet Loss when sending sensor node data to 4 sensor nodes, consequently, that 80% of the pulse data is received by the Coordinator node, this occurs because of the bottleneck factor in the ZigBee coordinator, bottleneck due to the multiplexing factor that occurs in Zigbee coordinator. accordingly, This is due to simultaneous data transmission by 5 sensor nodes, data can be received in stages but a decrease factor in Quality of Service (QoS), Receiver Signal Strength Indicator and Pathline when sending data bits per second (bps). Packet loss can be calculated from the reduction in bits sent by the receiver. In this research Radio Frequency devices used are the same as research [12], the difference is in the type or type of XBee e.g, XBee S2C and XBee S1 Pro.

Muhammad Anwar, Abdul Hanan Abdullah, Ayman Altameem, Kashif Naseer Qureshi, Farhan Masud, Muhammad Faheem, Yue Cao, and Rupak Kharel [13], in their research tried to implement the Wireless Body Area Network (WBAN) on the patient's body and analyze its routing system using the energy-aware link efficient routing approach (ELR-W), therefore, the goal is to save energy from the Wireless Body Area Network (WBAN). accordingly, in Wireless Sensor Network (WSN), WBAN, M2M Communication, it is very important for energy efficiency factors, i.e, batteries used, so the use of dynamic Power Supply such as PMFC is important, Plant-Microbial Fuel Cell (PMFC) as an energy source is especially appropriate in designing indoor systems or Outdoor environments [14], In this research, it refers to WBAN technology [14], by

developing a Blood Pressure sensor with a light type that is convenient to be used by Patients.

Moh. Khalid Hasan, Md. Shahjalal, Mostafa Zaman Chowdhury and Yeong Min Jang [15], in this research using Bluetooth Low Energy-based Wireless Sensor Network devices and the board used, is e-Health connected to the sensor electrocardiogram (ECG) using the hybrid OCC / BLE System. furthermore, From this research, a comparison of Quality of Service (QoS) in the OCC, BLE and Hybrid Schemes was produced. ECG has more complex output than Blood Pressure, ECG placement is also in patients 'chest and sensor placement which is at several points on patients' chest, while blood pressure only takes on the base of the wrist or wrist tip.

Gaël Loubet, Alexandru Takacs, Ethan Gardner, Andrea De Luca, Florin Udrea and Daniela Dragomirescu [16], using LoRa as a Wireless Sensor Network device to monitor patient health, LoRAWAN is Vehicular-based communication [17]. This LoR works on the 868 MHz ISM frequency radio. in this research the communication system built is Mesh Network with an approach to battery-free LoRaWAN sensing and communication nodes. accordingly, The power density of the electromagnetic waves is higher than 0.5 µW / cm2. So the Energy Consumption (EC) factor is important to be reduced to get a Long Life factor and stabilize the Packet Delivery Ratio (PDR) factor on reference [18] by evaluating root mean square error and deep learning method can produce 98% accuracy on PDR and EC Predictions. The use of LoRAWAN RF will be a good comparison with ZigBee RF in terms of QoS and Ability at long distances and hilly locations. In future research, it is necessary to compare and analyze the two types of RF Modules to send sensor data to Health monitoring.

Adolfo Di Serio, John Buckley, John Barton, Robert Newberry, Matthew Rodencal, Gary Dunlop 3 and Brendan O'Flynn in reference study [19], use the Zigbee to sending a heart rate data (BPM) and arterial Oxygen saturation (SpO2) data. ZigBee works on the 2.4 GHz frequency, but several other Wireless Sensor Network (WSN) devices have lower frequencies, for example the ISM Band 915 MHz frequency. The 915 MHz frequency shows a 10 dB return loss bandwidth of 55 MHz, with a gain value of -2.37 dB in free-space and -6.1 dBi on-body. In the references [20] using IEEE 802.15.4 Protocol simulation using Q-Learning to improve the performance of MAC Protocol and obtained a comparison of values from average latency, average backoff, channel access ratio, and transmission success ratio. The board used is the STM32F769 board and STM32L486 board which has a Cortex M7 processor (216 MHz and 120 MHz) and Cortex M4 (80 MHz). The SpO2 sensor is precise when combined with a Blood Pressure sensor in this Research. Complex data will make the patient's examination more detailed and accurate.

## III. MATERIAL AND METHOD

### A. Block Diagram that Represents Research as a Whole

Overall this research system is shown in Fig. 1. There are 3 parts shown in the dashed line showing three parts to be analyzed, the first part is the ZigBee RF Module connection

on the End Device and ZigBee RF Module on the Raspberry Pi 3. This section can be analyzed from Quality of Service (QoS) which includes the Power Receiver (PRx) and Receiver Signal Strength Indicator (RSSI). The construction of a star, tree, and mesh network is in the first part of the analysis, furthermore, then in the second analysis is Python programming and its connection using PuTTY and Web Server using XAMPP or WAMP followed by Website-based programming languages using HTML, PHP and Javascript or JSON can emerge Blood Pressure charts in real time in various platforms, at this stage, access to the domain is needed in order to be able to connect to the internet in realtime which is depicted at the 3rd analysis stage. furthermore at this research will be using the tool editor in HTML, Javascript or JSON using Dreamweaver Creative Cloud (CC) Software.

### B. LM358N Operational Amplifier

LM358N Op-Amp is a Low Power IC, easy to use on a dual channel op-amp. The function of an LM358N Operational Amplifier is a signal amplifier in AC and DC currents and as a high input impedance differentiation amplifier and a low impedance output amplifier. accordingly, IC LM358N Operational Amplifier can handle 3-32 Volt DC supply and Source up to 20 mA per channel. In Fig. 1 Schematic of LM358 N Op-amp shows Voltage Controlled Oscillator (VCO). In this case, the LM358N Operational Amplifier is applied to the sensor node, in Fig. 2.



Fig. 1. Block Diagram that Represents Research as a Whole.



Fig. 2. IC and Schematic of LM358N Op-Amp.

## C. Blood Pressure Sensor

Accordingly, on this research, The pressure sensor used is the MPS20N0040D-S type, this is a kind of solid pressure sensor, using MEMS technology, high reliability, and low cost. The pressure range is 0-5.8 Psi (40 kpa), the electricity supply is 5 volts DC and Constant Current is 1 mA. the input impedance of 4-6 Ω. bias voltage ± 25mV, full-scale output voltage 50-100 mV. In Fig. 3 shows the module of the Pressure Sensor used in this research. furthermore, Fig. 4 is a Pressure Sensor Block Type MPS20N0040D-S used in this research.



Fig. 3. Pressure Sensor.



Fig. 4. Block Diagram and Dimensions of Pressure Sensors.

In Fig. 6. It shows that the classification of blood pressure is divided into three, i.e, Ideal blood pressure, pre-high blood pressure, and high blood pressure. And there are two types of terms used to measure the 3 classifications e.g, Systolic (Top number) and Diastolic (Bottom Number). furthermore, Data in Fig. 6. In accordance with data Systolic and diastolic is as follows, e.g, 150/80 mmHg.

In Fig. 5. shows a Blood Pressure sensor diagram, in this diagram it will be developed using the MCU with Micro Arduino or Nanoduino so that the light version is convenient to be used by Patient.

The heart is the most important part of the body that is responsible for pumping blood throughout the body. accordingly, from the results of the examination can be concluded that the patient suffers from a particular disease, such as symptoms of stroke or heart disease. It is necessary to know the unit to state the amount of blood pressure is Millimeters of Mercury (Hydrargyrum) and then abbreviated mmHg. While KiloPascals (kPa) is a unit based on Standard International (SI) to express the amount of Pressure value, so to change mmHg to Kpa is 1 mmHg equal 133,322 Pascals (Pa) then 1 kPa equal 1000 Pascals (Pa) then mmHg Value x 133,322 Pa equal kPa value x 1000 Pa or mmHg value equal kPa value x 7.50062. while Psi stands for Pound Per square inch (Psi), Psi is used to express the value of pressure other than using kPa, 1 Psi equal 6.89475729 kilopascals.

In Fig. 5 a Blood Pressure Sensor block diagram contains several important components i.e, MCU is part of the Data Controller and Processor, can also be called a Microcontroller, in this research used Arduino Micro Microcontroller. Then the Output section is LCD, the LCD used for research development is an 8x2 bit LCD. With an actuator, a DC motor which is assigned to provide pressure in the form of air and a Pressure sensor component (Fig. 5) are connected to the amplifier circuit and Band Pass Filter. The amplifier circuit with the Band Pass Filter functions as a filter for analog frequencies. In detail, the Union body of the blood pressure sensor can be seen in Fig. 7.



Fig. 5. Block Diagram of a Blood Pressure Sensor.

## D. Pass Band Filter Circuit

Circuits that are built using Op-amps (Operational Amplifiers) and Capacitor and Resistor circuits serve to pass only high signals or only low ones or The Lower Frequency cutoff and the Higher Frequency cutoff [Fig. 8]. Result of 1st

filter for Lower Frequency cutoff is 0.278 Hz and the result of the 1st filter is higher. The frequency cutoff is 5.837 Hz. While the Result of 2nd filter for Lower Frequency cutoff is 0.278 Hz and the result of the 2nd filter is higher. Frequency cutoff is 21.80 Hz. whereas for the value of the mid-band gain of the first filter (A) is -12,156 and the mid-band gain of the Second filter (A) is 32,549.

For the Hz result change the value to International Standar (SI), e.g, 56 uF equal $56 \times 10^{-3}\ F$, and *10.2K* Resistor equal 10200 ohm.Pi value is 22/7 or 3.14, then 2 pi is 6.28.

- The Low Bandpass 1st Filter = flow = 1 / (2 pi (C4)(R3))

- The Low Bandpass 2ndFilter = flow = 1 / (2 pi (C3)(R1))

- The High Bandpass 1st Filter = fhigh= 1 / (2 pi (R4)(C2))

- The High Bandpass 2ndFilter = fhigh= 1 / (2 pi (R2)(C1))

- The mid-band gain of the first filter (A) = - R4/R3

- The mid-band gain of the Second filter (A) = - R2/R1

### E. Zigbee RF Module

In Fig. 9 shows 802.15.4 architecture, this shows there are 2 Layers on 802.15.4 architecture i.e, Zigbee and 802.15.4, in fact, the setting of ZigBee module, consist of two types ie, Zigbee S1 and Zigbee S2, Zigbee S1 is a type of Zigbee module with the function of start communication ability and Zigbee S2 to a tree and mesh communication. Therefore Zigbee S1 is IEEE 802.15.4 protocol and Zigbee S2 is Zigbee with the dynamic communication.



Fig. 6. Blood Pressure Classification for Adult.



Fig. 7. Uni Body of Blood Pressure Sensor MPX10 Type.



Fig. 8. Bandpass Filter Stage.



Fig. 9. Zigbee Architecture.

Fig. 10. Zigbee RF Module.

Zigbee is a wireless device, often referred to as a Wireless Sensor Network device, its specifications are small and Low Power and Low data rate. nevertheless, it is compatible and suitable for handling sensor Nodes, ZigBee has transmitted Power of 1 mW (0 dBm) with a data rate of 250 kbps, receiver sensitivity up to -92 dBm (1% packet error rate) and 100 dBm (1% packet error rate).[21] furthermore, ZigBee has a type i.e, XBee S1, Xbee S2, XBee Pro S1, and XBee Pro S2, with this specification XBee can send data in a tree, star or mesh depending on this type, S2 type can communicate between ZigBee with tree or mesh types. In this report XBee Pro S1 is used because the communication is a star network, in Fig.10 shows the ZigBee RF Module XBee Pro S1 device used in this research.

### F. Power Receiver (Prx (dBi)) and RSSI (dBm)

ZigBee is Radio Frequency that has a Power Transmitter of 1mW RF Power, so that if it is converted to dBm to 1 mW equal 0 dBm equal to -30 dB. So when we look for the Received Signal Strength Indicator (RSSI) first determine the value of the Power Receiver (Prx) in dBi units, accordingly the theory, ZigBee is an isotropic antenna type or omnidirectional antenna whose transmitter gain value (Gtx) and receiver gain (Grx) are -3 dBi or -3dB (decibel). As equation (1) is the equation to determine the value of the Received Signal Strength Indicator (RSSI) in units of dBm.

$$RSSI\ (dBm) = 10 \log (Prx) \tag{1}$$

and to determine the value of the Power Receiver (Prx), it is necessary to know the values of the Transmitter Gain (Gtx) in dB, Gain Receiver (Grx) in dB, Power Transmitter (Ptx) in dBi.

While the value of the wavelength λ is the result of the division of the speed value of light 3x108 with the value of the magnitude of the ZigBee is 2.4 GHz frequency or equal to 24x108 Hz so the result is 0.125 m. accordingly, equation 2 will determine the value of the Power Receiver (dBi).

$$Prx\ (dBi) = \frac{Ptx.Gt.Gr.\lambda^2}{(4\ \pi\ R)^2} \tag{2}$$

R states the distance in units (m), with the value π is 3.14 in short if we enter the values to find RSSI (dBm) at R or 3 m distance in the calculation equation [1] produces the value Prx (dB) equal $29.7x10^{-3}$ dB and if included in the equation (2), the result is *RSSI = 10 log (Prx)* then 10 log ($29.7x10^{-3}$), and the result is -25 dBm. furthermore, on the results of evaluation and analysis in this research, the results will be compared between measurements in experiments and simulations with a distance of 1– 00 m in the free space.

The RSSI relationship with d (distance) can be represented in equation (3) and equation (4). The RSSI value on the

Wireless Sensor Network of Zigbee module can be obtained with several models, one of which is the calculation of RSSI in the Free space propagation model.

$$RSSI\ (d) = Pt(D0) - 10n_p \log (D/D0) \tag{3}$$

Where *RSSI (d)* is the value of RSSI in dBm at distance D (meters), np is the Path Loss exponent, Pt (D) is a strength of a transmitter in dBm, D0 is a D at the beginning of the transmitter Pt (D0) at a distance of 1 meter.

$$D = 10^{[(P0\ -\ Fm\ -\ Pr\ -\ 10*np*\log\_10(f)+30*np-32.44)/10*np]} \tag{4}$$

Where *D* is distance (*m*), $P_0$ is a Power transmitter (*dBm*) at 0 distance, *Pr* is a signal Receiver, *F* is Frequency (*Hz*), $n_p$ is Path Loss exponent (Table I).

TABLE I. PATH LOSS EXPONENT (NP)

| Environment | Path Loss Exponent ($n_p$) |
|---|---|
| Free Space | 2 |
| Urban area cellular radio | 2.7 to 3.5 |
| Shadowed urban cellular radio | 3 to 5 |
| In Building Line-of-site | 1.6 to 1.8 |
| Obstructed in building | 4 to 6 |
| Obstructed in factories | 2 to 3 |

### G. Blood Pressure Node Sensor Test

In Fig. 11 is the process of testing Blood Pressure sensors using the ATmega 328p Microcontroller, this examination was successfully carried out by storing Systolic and Diastolic sensor data (mmHg) in the MySQL Database. furthermore, the blood pressure data processed by the ATmega 328p Microcontroller. furthermore, a ZigBee sends the data to the ZigBee Coordinator on the Raspberry Pi 3. The connection between the Raspberry Pi 3 and Zigbee RF module is shown in Fig. 12.



Fig. 11. Blood Pressure Connectivity Testing.



Fig. 12. Raspberry Pi 3 and Zigbee RF Module Connectivity.

## H. Blood Pressure Sensor Pseudocode–1

Pseudocode-1 shows how the Blood Pressure Sensor can work and provide Systolic and Diastolic (mmHg) values in the Blood Pressure sensor node made in this research. In the serial output monitor, the MAP (Mean Arterial Pressure) value is different, furthermore, this is the output produced at the sensor node, i.e, the Blood Pressure value in mmHg.

```
Blood Pressure Sensor Pseudocode-1
    a.   Data Type Analyzes
    b.   Float type
PressureMin=-15; //psi
PressureMax= 15;  // psi
Vsupply=5; // voltage supply
volta=0; maxvolt=0; volt=0; pressure=0;
MAP=0; maxv=0;
    c.   Integer type
analogInPin = A0, i;
    d.   BoudRate and Output Pins Analyzes
Boudrate = 9600 bps
Digital Pin Output = 3
    e.   Logic and Looping Process
If Digital Pin 3 = High / ON
Then
for(i=0;i<40;i=i+1){
volta = analogRead(analogInPin);
    f.   running the equation
volt=(volta*Vsupply)/(pow(2,10)-1);
maxv=max(abs(volt-2.5), maxvolt);
maxvolt=abs(maxv-2.5);
    g.   Give the delay time
delay(250);
    h.   Pressure ON and Equation of Pressure
pressure=(((maxvolt)-
.1*Vsupply)/((.8*Vsupply)/(PressureMax-
PressureMin)))+PressureMin;//psi
MAP=-1*(14.7-pressure*-1)*51.7-3.16/maxvolt;
//mmHg
    i.   Digital Pin 3 = LOW
digitalWrite(3,LOW);
    j.   Output on Serial Monitor
  Serial.print(MAP);
  Serial.println(MAP*1.1);
  Serial.println(MAP*0.8);
```

## I. ZigBee Blood Pressure in Python Pseudocode–2

Whereas in *Pseudocode-2*, the input process is from Python Code to MySQL database, accordingly, input data comes from Serial Port port = '/ dev / ttyUSB0' and this is the value of the Pressure Sensor node captured by the ZigBee Coordinator Node.

```
ZigBee Blood Pressure Python Pseudocode-2

1.  Import the Python Extension
import pymysql, time, serial
2.  serial data initialization
ser
=serial.Serial(port='/dev/ttyUSB0',baudrate=9600,
bytesize=serial.EIGHTBITS,
parity=serial.PARITY_NONE, timeout=6)
3.  Database Connection initialization
connection=pymysql.connect(host='localhost',
user='root',password='',db='zigbee',
charset='utf8mb4',
cursorclass=pymysql.cursors.DictCursor)
if(ser.isOpen()):
4.  Data Type Analyzes
Integer type = a, b;
Data serial = a,b;
```

```
5.  Cursor initialization
with connection.cursor() as cursor:
6.  Enter the Blood Pressure data to the MySQL
    Table on the zigbee Database
sql="insert into
zigbeebloodpressure(systolicdiastolic)VALUES (%s)"
cursor.execute(sql, (a))
connection.commit()
with connection.cursor() as cursor:
7.  Read The Single Record
sql="select 'id', 'systolicdiastolic' from
zigbeebloodpressure WHERE 'systolicdiastolic' =
%s"cursor.execute(sql, (a))
result=cursor.fetchone()
8.  Print the result and Close
finally: connection.close()
```

## J. Javascript Object Notation (JSON) Pseudocode-3

In *Pseudocode-3*, Javascript Object Notation (JSON) will form a Chart that can be displayed on Web Page based on HTML and PHP, this data comes from the MySQL database, i.e, Blood Pressure Sensor data.

```
Javascript Object Notation (JSON) Pseudocode-3
    1.  Connection Initialize
$connect = mysqli_connect("localhost", "root", "",
"zigbee");
    2.  Query Initialize
$query = '
SELECT sensors_bloodpressure_data,
sensors_bloodpressure2_data,sensors_data_date,
sensors_data_time -> Desc
$result = mysqli_query($connect, $query);
    3.  Array Rows and Tabel Created
$rows = array();
$table = array();
'label' => 'Date Time', 'Systolic (mmHg)',
'Diastolic (mmHg)',
    4.  Date and Time Initialized
$datetime = explode(".", $row["datetime"]);
"v" => 'Date(' . $datetime[0] . '000)'
"v" => $row["sensors_bloodpressure_data"]
"v" => $row["sensors_bloodpressure2_data"]
    5.  Calling Javascript
src="https://www.gstatic.com/charts/loader.js"
src="ajax.googleapis.com/ajax/libs/jquery/1.10.2/jqu
ery.min.js"
    6.  Name the Chart Header
<h2 align="center">Display Chart of Blood Pressure
Real-Time Monitoring (MeRL) </h2>
```

## IV. Result and Analysis

### A. Blood Pressure Sensor Testing

Experiments a, b, c, and d [Fig. 13] are giving different treatments to the Pressure sensor. It is noted that Pulse sensors work or are active, in experiments a. a straight line on the value of Systolic and diastolic (mmHg) is the condition when the mini-Pump is off. In experiment b. the graph shows an increase at one time, not every time, only 1-time increase, this is because there is no supply voltage from the amplifier IC or amplifier. accordingly, In experiment c. sensors such as losing a stable position, this is because of the Pressure sensor no supply input from a Microcontroller or Arduino Analog Pin (A0). Experiment d is a graphical difference when the sensor pressure is turned off or the mini-pump off is then turned on after a few seconds, consequently, that there is a trigger when

the Pressure sensor is on or the mini-pump is on.The blue line shows the value of Systolic (mmHg) and the red line shows the value of Diastolic (mmHg), the experiment is done at the same time, the processor is the Arduino Microcontroller, and the plot or graph is formed from the Arduino Serial Port, furthermore, Arduino Microcontroller has 2 Output Functions i.e. Serial Plotting and Serial monitor, data in the form of plotting is made through the Arduino Plot series, while serial data shows the values of processing the Arduino Microcontroller in the form of Systolic and Diastolic (mmHg) values. The value in the graph shows the value of mmHg, can be converted into an international unit, i.e, Psi, kPa and atm for more complex data needs.

- 1 atm=760 mmHg = 760 torr = 101.3 KPa= 14.7 psi

- 400 mmHg = 400 mmHg x (1 atm/ 760 mmHg) = 0.52 atm

- 400 mmHg = 400mmHg x (101.3 kPa / 760 mmHg) = 53.3 kPa

### B. Receiver Signal Strength Indicator (RSSI dBm)) of ZigBee RF Module

Receiver Signal Strength Indicator (RSSI) is stated in Fig. 14. There are two comparisons, i.e, RSSI Simulation and RSSI Experiment at the field, in the results of these experiments and calculations, Zigbee communication at a distance of 5 meters, RSSI simulations show a value of -29 dBm and the experiment shows a value of -40 dBm, at a distance of 100 m, RSSI shows a value of -55 dBm (simulation) and -86 dBm (experiment).

RSSI (-dBm) experiment data retrieval is using DIGI X-CTU software by sending Blood Pressure data using Zigbee from different distances, and this data is recorded every meter. RSSI (dBm) in Fig. 14 shows a decrease in signal strength.The distance of (Tx) and Receiver (Rx) determine of RSSI(-dBm) value, the farther the distance Tx-Rx, the greater the value of RSSI (-dBm). The experiment was carried out from 1-50 meters in the Free Space area.

### C. Output Graph on the Web Page

In Fig. 15 and Fig. 16 shows the Real-time display of the MySQL data database. This data is Systolic and Diastolic (mmHg) data, then this data is used as an indicator parameter that shows the patient's condition. This data is seen on the WEB Page, the Localhost system will be upgraded to the Domain level, so that data can be seen on all platforms, not only that JSON is used so that the quality of HTML can be used on smartphone and tablet platforms with different WEB Page views on Personal Computer, in this case, is expected to be a flexible WEB Page, so that data can be viewed easily by the user. The blue line shows the value of Systolic (mmHg) the value of the upper part of the Blood Pressure and the red line is Diastolic (mmHg) is the value of the bottom, the graph shows the ups and downs of Systolic and Diastolic values because of several checks and different results. Furthermore, the more checks are carried out, the more blood pressure data on the MySQL database and the graphs generated will be more complex and details.



Fig. 13. Blood Pressure Test Consisting of a, b, c and d.

Fig. 14. RSSI (dBm) of Zigbee RF Module.



Fig. 15. Graph of Systolic (mmHg).



Fig. 16. Graph of Diastolic (mmHg).

## V. CONCLUSION AND SUGGESTION

Sensor data on a ZigBee module can be sent properly as indicated by the RSSI parameter, at a distance of 5 meters RSSI simulations show a value of -29 dBm and -40 dBm in the following experiments so that it gets smaller. furthermore at a distance of 50 m, the Receiver Signal Strength Indicator (RSSI) shows around –70 dBm in field measurements, while in simulations of equations 1 and 2 the RSSI results are quite

good at -50 dBm, this result is better than field measurements, at a distance 100 m, the RSSI value produced is -55 dBm (simulation) and -86 dBm (experiment). The sensor node works very well, blood pressure can be stored in the MySQL database using the Python programming language that works on the Raspberry Pi 3 model B. And the data can be displayed on the Web Page using JSON. Accordingly, this Research still uses Localhost, it needs to be improved by moving localhost to the domain so that data can be seen on all platforms such as smartphones or tablets connected to the internet.

## VI. DISCUSSION

Some points that need to be added for future research development are (1) The prototype that is already in the form of Light type and fix on PCB and design Product, to make it comfortable for Patients (2) Programming Languages that support Graphical User Interface (GUI) output from Blood Pressure and other sensors on the mobile platform devices (3) The Security Method of IoT Protocol needs to be added to this research. Research development in monitoring patient health based on Internet of Things (IoT) is to use Algorithms for data confidentiality and security, e.g, RC4 Algorithm. RC4 Algorithm is one of the algorithms used for security on IoT using the encryption method. Improvements are also needed in this research by utilizing a programming language that supports the appearance of a capable Graphical User Interface (GUI) for display flexibility in all internet connected platforms/devices.

## REFERENCES

[1] Muhammad Usman Ali, Soojung Hur and Yongwan Park, "Wi-Fi-Based Effortless Indoor Positioning System Using IoT Sensors", Sensors 2019, 19, 1496; doi:10.3390/s19071496

[2] Eyal Ronen, Adi Shamir, Achi-Or Weingarten, Colin O'Flynn, "IoT goes nuclear: Creating a ZigBee chain reaction", 2017 IEEE Symposium on Security and Privacy, DOI: 10.1109/SP.2017.14

[3] Puput Dani Prasetyo Adi and Akio Kitagawa, "Performance Evaluation WPAN of RN-42 Bluetooth based (802.15.1) for Sending the Multi-Sensor LM35 Data Temperature and RaspBerry Pi 3 Model B for the Database and Internet Gateway" International Journal of Advanced Computer Science and Applications (IJACSA), 9(12), 2018. DOI: dx.doi/10.14569/IJACSA.2018.091285

[4] Pantaleone Nespoli, David Useche Pelaez, Daniel Díaz López and Félix GómezMármol, "COSMOS: Collaborative, Seamless and Adaptive Sentinel for the Internet of Things", Sensors 2019, 19, 1492 ; doi:10.3390/s19071492

[5] Puput Dani Prasetyo Adi and Rahman Arifuddin, Design of Tsunami Detector Based Sort Message Service Using Arduino and SIM900A to GSM/GPRS Module, JEEMECS (Journal of Electrical Engineering, Mechatronic and Computer Science) Volume 1, No.1. 2018, DOI: doi/10.26905/jeemecs.v1i1.1982

[6] Ala Al-Fuqaha, Mohsen Guizani, Mehdi Mohammadi, Mohammed Aledhari, and Moussa Ayyash, "Internet of Things: A Survey on Enabling Technologies, Protocols, and Applications", IEEE Communication Surveys & Tutorials, VOL. 17, NO. 4, Fourth Quarter 2015, DOI: 10.1109/COMST.2015.2444095

[7] Yixin Mei, Fan Li, Lijun He and Liejun Wang, "Joint Source and Channel Rate Allocation over Noisy Channels in a Vehicle Tracking Multimedia Internet of Things System", Sensors 2018, 18, 2858; doi:10.3390/s18092858

[8] Aitor Agirre, Aintzane Armentia, Elisabet Estévez and Marga Marcos, "A Component-Based Approach for Securing Indoor Home Care Applications", Sensors 2018, 18, 46; doi:10.3390/s18010046

[9] Emilio Andreozzi, Gaetano D. Gargiulo, Antonio Fratini, Daniele Esposito and Paolo Bifulco, "A Contactless Sensor for Pacemaker Pulse Detection: Design Hints and Performance Assessment", Sensors 2018, 18, 2715; doi:10.3390/s18082715

[10] Dimiter V. Dimitrov, MD, PhD, "Medical Internet of Things and Big Data in Healthcare", A Digital Achive and Reference Linking Platform of Korean Medical Journals, Published online July 31, 2016. https://doi.org/10.4258/hir.2016.22.3.156

[11] Giorgio Biagetti, Paolo Crippa, Laura Falaschetti and Claudio Turchetti, "Classifier Level Fusion of Accelerometer and sEMG Signals for Automatic Fitness Activity Diarization", Sensors 2018, 18, 2850; doi:10.3390/s18092850

[12] Muhammad Niswar, Amil Ahmad Ilham, Elyas Palantei, Rhiza S. Sadjad, Andani Ahmad, Ansar Suyuti, Indrabayu, Zaenab Muslimin, Tadjuddin Waris, Puput Dani Prasetyo Adi, "Performance evaluation of ZigBee-based wireless sensor network for monitoring patients' pulse status", 2013 International Conference on Information Technology and Electrical Engineering (ICITEE) DOI: doi/10.1109/ICITEED.2013.6676255

[13] Muhammad Anwar, Abdul Hanan Abdullah, Ayman Altameem, Kashif Naseer Qureshi, Farhan Masud, Muhammad Faheem, Yue Cao, and Rupak Kharel, "Green Communication for Wireless Body Area Networks: Energy Aware Link Efficient Routing Approach", Sensors 2018, 18(10), 3237; https://doi.org/10.3390/s18103237

[14] Edith Osorio de la Rosa , Javier Vázquez Castillo , Mario Carmona Campos, Gliserio Romeli Barbosa Pool, Guillermo Becerra Nuñez, Alejandro Castillo Atoche and Jaime Ortegón Aguilar, "Plant Microbial Fuel Cells–Based Energy Harvester System for Self-powered IoT Applications", Sensors 2019, 19, 1378; doi:10.3390/s19061378

[15] Moh. Khalid Hasan, Md. Shahjalal, Mostafa Zaman Chowdhury and Yeong Min Jang, " Real-Time Healthcare Data Transmission for Remote Patient Monitoring in Patch-Based Hybrid OCC/BLE Networks ", Sensors 2019, 19, 1208; doi:10.3390/s19051208

[16] Gaël Loubet, Alexandru Takacs, Ethan Gardner, Andrea De Luca, Florin Udrea and Daniela Dragomirescu, "LoRaWAN Battery-FreeWireless Sensors Network Designed for Structural Health Monitoring in the Construction Domain", Sensors 2019, 19, 1510; doi:10.3390/s19071510

[17] José Santa, Ramon Sanchez-Iborra, Pablo Rodriguez-Rey, Luis Bernal-Escobedo and Antonio F. Skarmeta, "LPWAN-Based Vehicular Monitoring Platform with a Generic IP Network Interface", Sensors 2019, 19, 264; doi:10.3390/s19020264

[18] Muhammad Ateeq , Farruh Ishmanov, Muhammad Khalil Afzal, and Muhammad Naeem, "Multi-Parametric Analysis of Reliability and Energy Consumption in IoT: A Deep Learning Approach", Sensors 2019, 19, 309; doi:10.3390/s19020309

[19] Adolfo Di Serio, John Buckley, John Barton, Robert Newberry, Matthew Rodencal, Gary Dunlop 3 and Brendan O'Flynn, "Potential of Sub-GHzWireless for Future IoT Wearables and Design of Compact 915 MHz Antenna", Sensors 2018, 18, 22; doi:10.3390/s18010022

[20] Jiheon Kang and Doo-Seop Eom, "Offloading and Transmission Strategies for IoT Edge Devices and Networks", Sensors 2019, 19, 835; doi:10.3390/s19040835

[21] Puput Dani Prasetyo Adi, Analisis kinerja jaringan sensor nirkabel untuk monitoring denyut nadi pasien, April 2018, DOI: 10.13140/RG.2.2.29145.83040.

# Storage Consumption Reduction using Improved Inverted Indexing for Similarity Search on LINGO Profiles

Muhammad Jaziem bin Mohamed Javeed[1], Nurul Hashimah Ahamed Hassain Malim[2]

School of Computer Sciences, Universiti Sains Malaysia, Penang, Malaysia

*Abstract*—**Millions of compounds which exist in huge datasets are represented using Simplified Molecular-Input Line- Entry System (SMILES) representation. Fragmenting SMILES strings into overlapping substrings of a defined size called LINGO Profiles avoids the otherwise time-consuming conversion process. One drawback of this process is the generation of numerous identical LINGO Profiles. Introduced by Kristensen et al, the inverted indexing approach represents a modification intended to deal with the large number of molecules residing in the database. Implementing this technique effectively reduced the storage space requirement of the dataset by half, while also achieving significant speedup and a favourable accuracy value when performing similarity searching. This report presents an in-depth analysis of results, with conclusions about the effectiveness of the working prototype for this study.**

*Keywords—Simplified Molecular-Input Line-Entry System (SMILES); LINGO profiles; similarity searching; inverted indexing*

## I. INTRODUCTION

Rapid advances in technology over the past few years have allowed for many virtual screening experiments to be conducted extensively [1]. In ligand-based screening, large chemical databases consisting of small molecules are effectively screened by a query molecule as to identify molecules with similar biological activity, applying the well-known similarity principle that "structurally similar molecules are likely to have similar properties" [2,3,4,5,6]. The query structure itself normally exhibits a potentially useful level of biological activity and might be, for example, a competitor's compound or a structurally novel hit from an initial high-throughput screening (HTS) experiment [7]. Both the query and database molecules are characterized by descriptors.

Simplified Molecular-Input Line-Entry System (SMILES) is a type of 1D representation [8] which represents molecular structures in strings format [9,10]. The SMILES specialized algorithm known as LINGO [11] is introduced in the field as it delivers a required level of simplicity for retrieving the molecules from database. LINGO representation avoids the necessity for producing an explicit model of the chemical structure in the form of either a graph or a 3D structure because it generates the representation of a molecule directly from canonical SMILES [12].

The continuing rise in the number of compounds to be processed is one of the common challenges which have to be confronted in this field, in terms of the accompanying demand for higher processing power and storage costs [13]. Small libraries can take up to 10^5 compounds, while commercially available datasets have approximately (2 x 10^7 compounds) in their libraries [14]. Many research studies have been conducted to address this problem by developing a coherent technique to store the compounds, but this has been limited only to compounds represented in 2D fingerprints [15]. This situation has, consequently, led to the necessity of introducing a data structure efficient enough to store the compounds represented in LINGO Profiles. An inverted index is a type of index data structure which is commonly used to encode data in string format [16,17,18]. It allows for term-based searches to be more effective [19,20]. This study seeks to ascertain whether the introduction of inverted indices actually achieves any reduction in storage and processing costs when performing similarity searching. Therefore, the rest of the paper is organised as follows: Section II presents several related studies pertaining to similarity searching methods. Next, Section III elaborates the research methodology in terms of implementation and experimental design, while Section IV discusses the analyses outcomes. Lastly, this paper ends with a conclusion depicted in Section V.

## II. BACKGROUND REVIEW

The search for compounds similar to a given target ligand structure and compounds with defined biophysical profiles are two main important principles in modern drug discovery process [21]. Both tasks make use of molecular descriptors with different complexity (atomic, topographic, sub structural fingerprints, 3D, biophysical properties, etc.) leading to different representations of the same molecule [22]. In general, structural representation, also known as molecular descriptor is used in describing the characteristics of compounds [23].

Ozturk and co-workers [24] used a state-of-the-art algorithm; the Weighted Nearest Neighbor-Gaussian Interaction Profile (WNN-GIP) with which to evaluate the performance between 1D SMILES representation and 2D representation-based descriptors in the protein-drug interaction task. Their investigation successfully demonstrated that SMILES-based methods [25] of molecular similarity comparison perform as well as 2D-based methods. Moreover, SMILES-based kernels were found to be computationally faster and more flexible than their 2D competitors.

In a different experiment, comparisons were examined between 2D fingerprints such as Daylight, MOLPRINT 2D, MACCS, and Open Babel with 3D shape-based methods,

typically SHAEP, PARAFIT and ROCS, in order to measure the efficiency of the similarity searching method across a range of virtual screening methods [26]. Results showed in the past [26][27] that 3D shape-based methods could not perform as well as a simple fingerprint similarity search, in spite of giving conformational information (i.e. shape information) and atomic coordinates of a compound.

Most previous drug-target interaction prediction tasks involving LINGO have utilized the Tanimoto coefficient. Vidal and colleagues [28] used a bioisostere dataset to compute intermolecular similarity between bioisosteric molecules and some randomly sampled pairs of molecules using an integral Tanimoto coefficient. The average similarities (LINGOsim) obtained effectively demonstrated that important information about a molecule is stored in LINGO Profiles. On the other hand, LINGO-DOSM, introduced by Hentabli et al [29], outperformed other descriptors such as EPFP4, GRFP, MACCS etc. LINGO-DOSM is the integral set derived from a given DOSM string. DOSM allows rigorous structure specification by implementing a small and natural grammar. The positive performance of LINGO-DOSM is not only limited to the top 5% for MDDR but it also gives best results for the top-1% for MDDR. This is mainly due to limiting the selection of LINGO length to just four characters. Finally, using the Briem and Lessel benchmark, Andrew and colleagues concluded that LINGO generated from isomeric SMILES can offer better retrieval rates, compared to non-isomeric SMILES. In addition, when LINGO was compared with more complex approaches (Daylight fingerprint) [25], it managed to identify active compounds better for two activity classes (ACE and TXA2).

The effectiveness of LINGO in predicting the property/activity of one molecule compared with another molecule similar to it, however, has a limitation [30]. This technique is associated with the length of the substrings obtained from the fragmentation of a canonical SMILES string, requiring the manipulation of the string and meaning that the processing cost will increase linearly along with the SMILES length [28]. Since search efficiency is progressively more vital with the ongoing expansion of these databases, scalability problems naturally arise when virtual compounds or recently synthesized compounds are added accordingly [31]. A variety of data structures and algorithms were consequently introduced throughout the years to accelerate this process by reducing the search, i.e. by rapidly eliminating the molecules that are not homogenous to the query, without computing their similarity to the query [32].

Imran and co-workers [33] presented a new algorithm known as the SIML ("Single-Instruction, Multiple-LINGO") to measure the similarity between molecules. Each multiset in a molecule is represented in 32-bit integers and it is stored in a sorted vector of 4-Lingos (represented as integers). A new algorithm, (1), was derived based on the vector representation of the multisets. This sparse vector algorithm speeds up the

computation involved, as for every Tanimoto calculation only the intersection size $\langle A, B \rangle$ needs to be calculated.

$$T_{AB} = \frac{\langle A,B \rangle}{\langle A,A \rangle + \langle B,B \rangle - \langle A,B \rangle} \tag{1}$$

Outside the field of cheminformatics, numerous information retrieval communities in general have been conducting experiments for decades on searching text in large datasets [34]. State-of-the-art algorithms from general information retrieval, known as inverted indices, are considered applicable for use in cheminformatics, as both domains arrived at the same similarity measure and representation [35] independently from one another. Features are associated with each respective list of documents contained in a given database, as shown in Fig. 1.

The features-documents association guarantees the reduction of the similarity computations between database molecules and the query as it removes database molecules which are irrelevant to the desired list. This approach can also be applied directly to SMILES string representations for molecules.

Kristensen et al. [36] proposed performing a similarity search between a target and database compounds represented using LINGO multisets by representing the database as inverted indices. The idea was to keep the LINGO multisets as a vector, where every cell in the vector is assigned to hold one of the LINGO identifiers (ID) from the verbose representation. Unlike SIML which uses two arrays to represent a LINGO multiset, verbose representation utilizes only an array to store the whole multisets including duplicate LINGO represented using multiple different IDs as shown in panel (a) of Fig. 2.



Fig. 1. Molecules Represented in Fingerprint Format are Stored in Inverted Index Data Structure.

| LINGO | LINGO ID |
|-------|----------|
| c0cc | 19 |
| cc0L | 23 |
| c0cc | 41 |
| ccc0 | 10 |
| cccc | 15 |
| cccc | 34 |

a.) Verbose Representation



b.) Inverted indices Representation

*S = Compounds SMILES string

Fig. 2. (a) Each LINGO is Associated with their Respective IDs. (b) LINGO and their Reference to their Original SMILES String in Inverted Indices Representation.

Input from panel (a) is used to create inverted indices (panel (b)) listing all the multisets associated with a given ID. Similarities are computed based on the value stored in the counting vector after the inverted indices are traversed. This strategy, however, led to a drawback as multiple occurrences of the similar LINGO in a compound will consume more storage space. It is certainly not feasible for huge datasets (e.g. ChEMBL). In addition, the construction of the inverted indices necessitates a search of the largest ID in the dataset. These situations will cause the increase in the processing time and consume high amount of resources, when performing similarity searching process. Besides, Kristensen work is only practical for chemical dataset such as Maybridge and ZINC.

Instead of finding a new method for indexing a database, a small modification of the inverted indexing scheme introduced by Kristensen et al. [36] is proposed in this study. Verbose representation is eliminated by the introduction of a pattern matching approach to resolve a query. This modification is made to increase the available storage space and to minimize the time taken to search a LINGO. A brief explanation of how the indexing method for this study was implemented is discussed in the following section.

## III. METHODOLOGY

The work was conducted purely on the 102,540 MDDR dataset compounds, where searches were focused only on selected structures from eleven activity classes. The first experiment of this study aimed at measuring the recall values obtained by LINGO Profiles on MDDR dataset, comparing it with various other fingerprints. The second experiment of this study intended to perform similarity searching based on the proposed indexing method, which as discussed earlier in the literature. The time taken and the storage consumption for both experiments were to be computed along before presenting a full discussion of these results in the next section.

### A. Performing Similarity Searching in Sequential Manner

A q-LINGO is a q-character string which may include letters, numbers, and symbols such as "(",")", "[", "]", "#", etc. and which is obtained by stepwise fragmentation of a canonical SMILES molecular representation [28]. Before the LINGOs are created from a compound, the compound ring numbers must be substituted for "0". If atoms such as "Cl" and "Br" are present, they will be replaced by "L" and "R", respectively. Raw MDDR Dataset (file A) stores all possible LINGOs for similarity searching after it is being fragmented and modified from the original MDDR dataset. It is attached together with its respective ID in sequential manner. Fig. 3 shows the whole process in generating LINGO Profiles.

Using the raw MDDR dataset (file $A$), to obtain LINGO for our query string (compound $A$) and a MDDR database compound (compound $B$), the ID of the compounds was compared with file $A$. Next, the LINGO$sim$ function was used to calculate the similarities between the two compounds. Based on a comparison of the LINGOs of the two compounds, $A$ (query compound) and $B$ (MDDR database compound) any intermolecular similarities were computed using the integral Tanimoto coefficient. $N_{Ai}$ represents the number of LINGOs of type ($i$) in compound $A$, while $N_{Bi}$ represents the number of LINGOs of type ($i$) in compound $B$, and ($l$) is the number of LINGOs contained in either compound $A$ or $B$.

### B. Performing Similairty Searching using Proposed Indexing Scheme

Two columns existed in this inverted indexing scheme ("Word" and "Documents") allow the query to perform similarity searching via random access [37]. "Word" column in Table I can be referred to as the unique LINGO Profiles obtained from the MDDR dataset and the "Document" column signify the compound IDs which contains the respective LINGO [37].

From the list (file $A$) generated earlier, it is possible to map the LINGO and IDs into the indexing scheme. 409,752,8 LINGO Profiles contained in file $A$ are compared with each other and if two or more identical LINGO Profiles is found, then their respective ID are appended together with the LINGO Profile in the list. In the end, the indexed database would only have 4054 unique LINGO Profiles. Fig. 4 summarizes the whole process.

Fig. 3.   Modifying and Fragmenting LINGO Profiles from MDDR Dataset.



Fig. 4.   Comparing LINGO Profiles and Eliminating Redundant LINGO Profiles on file a; as to Generate Indexed Dataset.

Calculating similarity values using our proposed method differed from the conventional method because it was based on a pattern-matching technique. Whenever the LINGOs in the query compound were found in the indexed database, the IDs in the "Document" column were retrieved and the frequency of occurrence is accumulated and calculated accordingly. The ranked list obtained were then sorted in a descending order to calculate the recall values. The whole process is illustrated in Fig. 5.

Table II shows the activity classes which were used in both experiments. Activity classes that were used in the experiments are slightly different in nature. The diversity was determined

using the main pairwise Tanimoto similarity (MPS) and it is included in Table I. Structurally homogenous classes such as Renin and ATI has high MPS value as compared to COX and PKC which have low MPS value since they are structurally diverse.

TABLE I.      STRUCTURE OF AN INVERTED INDEXING SCHEME

| Word | Documents |
|------|-----------|
| Cow | Document 1, Document 4, Document 6, Document 9, Document 15 |
| The | Document 2, Document 5, Document 8 |
| Hello | Document 12 |
| Cat | Document 7 |

TABLE II.      ACCURACY COMPARISON BETWEEN LINGO AND OTHER DESCRIPTORS (TOP : ACTIVES RETRIEVED; BOTTOM :  RECALL)

| Activity Classes | Number of Active Structures | Pairwise Similarity (Mean) | |
|------------------|-----------------------------|----------------------------|---|
| Renin inhibitors | 1130 | 0.290 | Most Homogenous |
| Angiotensin II ATI antagonists | 943 | 0.229 | |
| HIV Protease inhibitors | 750 | 0.198 | |
| Thrombin inhibitors | 803 | 0.180 | |
| Substance P inhibitors | 1246 | 0.149 | |
| 5HT3 antagonists | 752 | 0.140 | |
| D2 antagonists | 395 | 0.138 | |
| 5HT1A agonists | 827 | 0.133 | |
| 5HT reuptake inhibitors | 359 | 0.122 | |
| Protein Kinase C inhibitors | 453 | 0.120 | |
| Cyclooxygenase Inhibitors | 636 | 0.108 | Most Heterogenous |



Fig. 5.   Process Involved when Performing Similarity Searching using the Proposed Methodology.

## IV. RESULTS AND DISCUSSION

This section is divided into two sub-sections: A and B. Section A basically confirms Vidal's work via replication and compares performance to other fingerprints. Section B discusses the performance of the proposed method regarding time and storage consumption when benchmarked with the conventional method.

### A. Comparing Accuracy between LINGO Profiles and Various different Fingerprints

The performance of the similarity searching process can be evaluated based on its effectiveness. Effectiveness includes the calculation of the recall value in every single search. The recall value, R is calculated by dividing the number of actives retrieved at the end of the process, n, by the number of compounds that available in the activity class, N, as shown in (2). In other words, recall can be defined as the percentage of the active molecules, which is gained from the cut-off point in the ranked list. Some of the cut-off points that have been widely used are at 1% and 5%. In this experiment, we only use 1% cut-off. The recall value gained indicated the probability of structures that are showing positive to the target. Thus, the higher the recall value gained, the higher the number of structures that react positively towards the target, which implies the accuracy of the method. Units

$$R = \frac{n}{N} \qquad (2)$$

The performance of similarity searches using LINGOs was compared with the performance of similarity searching using various fingerprints obtained from the work of Malim [23]. A total of 110 searches were performed using 10 queries from 11 activity classes. These searches were executed in accordance to Fig. 6. Table III presents the average results of the number of actives retrieved and recall values.

From Table III, the superior performance of ECFP4 is evident in comparison with other fingerprints and LINGO Profiles, except for two activity classes where LINGOs outperform ECFP4. However, it was observed that the performance of LINGO was comparable with other fingerprints such as MDL, Daylight, and Unity in general. A closer analysis of the difference in the accuracy between both descriptors (ECFP4 and LINGO) reveals that ECFP4 outperformed LINGO only by a small average difference of 2.975%. Renin recorded the highest difference between both methods at 9.13 %, while the lowest difference value was observed in the Thrombin activity class, which favors LINGO Profiles at 0.41%.

TABLE III. ACCURACY COMPARISON BETWEEN LINGO AND OTHER DESCRIPTORS (TOP: ACTIVES RETRIEVED; BOTTOM: RECALL)

| Activity Classes | Descriptors | | | | |
|---|---|---|---|---|---|
| | *Unity* | *LINGO* | *Daylight* | *ECFP4* | *MDL* |
| 5HT1A | 56 6.79 | 64 7.77 | 59 7.15 | **81 9.79** | 53 6.46 |
| 5HT3 | 59 7.83 | 68 9.10 | 63 8.30 | **89 11.89** | 49 6.55 |
| 5HTReuptake | 21 5.86 | 20 5.82 | 19 5.40 | **24 6.83** | 20 5.58 |
| AT1 | 90 9.49 | 154 16.36 | 99 10.54 | **236 25.02** | 114 12.10 |
| COX | 15 2.41 | 14 2.34 | 21 3.22 | **28 4.45** | 15 2.48 |
| D2 | 19 4.74 | 22 5.70 | 22 5.63 | **27 6.86** | 17 4.33 |
| HIVP | 46 6.19 | 51 6.88 | 37 4.90 | **87 11.57** | 44 5.83 |
| PKC | 21 4.57 | 28 6.23 | 22 4.88 | **35 7.79** | 17 3.75 |
| Renin | 167 14.76 | 316 28.02 | 133 11.76 | **420 37.15** | 126 11.11 |
| SubP | 70 5.61 | **120 9.7** | 57 4.53 | **120 9.7** | 37 2.92 |
| Thrombin | 54 6.69 | **60 7.45** | 33 4.07 | 57 7.04 | **60 7.45** |



Query Compound A
ID 1

CCc1nc(N)nc(N)c1c2ccc3CCCN(CCCOC)c3c2

CCc0, Cc0n, c0nc, 0nc(, nc(N, c(N), (N)n, N)nc, )nc(, nc(N, c(N), (N)c, N)c0, )c0c, c0c0, 0c0c, c0cc, 0ccc, ccc0, cc0C, c0CC, 0CCC, CCCN, CCN(, CN(C, N(CC, (CCC, CCCO, CCOC, COC), OC)c, C)c0, )c0c, c0c0

MDDR Database
ID 3 Compound B

C[C@H](CC(=O)O)C(=O)N1CCC[C@H]1C(=O)O

C[C@, [C@H, C@H], @H](, H)(C, ](CC, (CC(, CC(=, C(=O, (=O), =O)O, O)O), )O)C, O)C(, )C(=, C(=O, (=O), =O)N, O)N0, )N0C, N0CC, 0CCC, CCC[, CC[C, C[C@, [C@H, C@H], @H]0, H]0C, ]0C(, 0C(=, C(=O, (=O), =O)O

$$\frac{\sum_{i=1}^{l} 1 - \frac{N(Ai) - N(Bi)}{N(Ai) + N(Bi)}}{l}$$

Similarity Coefficient: 0.85

Ranked List

Fig. 6. Process Involved when Performing Similarity Searching using Tanimoto Coefficient.

This being the case, according to the work of Hert [38], the nature of the defined activity classes themselves may affect the performance of similarity searches, as more actives may be retrieved in homogenous activity classes, compared to heterogenous ones. Homogenous classes consist of compounds which are less diverse, as opposed to classes with fewer common fragments shared between their compounds, which are described as heterogeneous classes. It can, therefore, be concluded that LINGO works better in homogenous classes as compared to heterogeneous classes. A higher number of active compounds are retrieved in activity classes such as Renin and AT1, in contrast to heterogeneous classes such as COX and 5HT Reuptake. The outcome of this experiment is, then, in agreement with Hert's findings.

Based on the results of this study, it can be summarized that LINGOs may act as an effective alternative to ECFP4 and other fingerprints when performing similarity searching, since this method offers the capability of obtaining a high-accuracy value for a variety of activity classes. It should be noted, however, that the superiority of ECFP4 is widely-known, due to its ability to encode as much structural information as possible when representing the compounds. LINGO profiles, in contrast, only allow for the strings to be observed by shifting one position at a time.

### B. Analyzing the Performance of the Proposed Method in Terms of Time Taken and Storage Consumption

*1) Time complexity:* Measuring the time taken for both methods is a very labour-intensive process, as it depends on the compiler and the type of computer or speed of the processor. For this research, the in-built time libraries in JAVA were used to determine the time taken. The timer was started before importing the input file and ended after the search was completed. The elapsed time was measured in milliseconds and for ease of reading it was then converted to hours.

Performing similarity searching using the proposed method is 782 times faster than the same using the conventional method. Achieving such an increase in speed was due to several reasons. Firstly, the indexed database which was created based on a raw MDDR dataset, contained fewer entries than the raw MDDR dataset itself. There was a total of 4053 unique LINGO Profiles in the indexed database as compared to a total of 4097258 LINGO Profiles which were generated in the raw MDDR dataset. With the reduction of the file size, time taken for a query compound to perform similarity searching using an indexed database would be reduced accordingly as now it only has smaller number of entries to browse through, in contrast to similarity searching performed on a raw MDDR dataset which requires a query compound to scan through the whole file to search for a LINGO Profile. The reduction in the file size, will be described in the next section.

*2) Storage complexity:* Storage complexity is determined by considering the maximum amount of capacity needed by the secondary storage to store the raw MDDR dataset and the indexed database. The measurement unit used in this study

was Megabytes (MB) (1,000,000 bytes in decimal notation). Specifically, there were no tools, libraries or applications used to measure the size of both files, as the sizes of the files were printed automatically by the operating system (OS) after the implementation process. The file size of the indexed database is smaller than the raw MDDR dataset. The reduction by almost half of the file size was achieved through the implementation of the inverted indexing technique, which yields a smaller number of entries in the file. The Raw MDDR dataset contained 4097528 entries, as each entry consisted of a LINGO Profile and its respective index number, as can be seen in Fig. 7.

Each entry here can be referred to as a string and each character within it has a size of a byte (8 bits), as the nature of JAVA language which encodes the strings in UTF-8 format. The '8' in UTF-8 means it uses 8-bit blocks to represent a character. The number of blocks needed to represent a character varies from 1 to 4. Theoretically, one string in a raw MDDR dataset might have a size which falls between 10-11 bytes. Multiplying the size of a string with the number of entries in the raw MDDR dataset and dividing it with the total number of bits in 1 Mb (8000000) will give an approximately similar result. Therefore, having a large number of entries will lead to a larger file size.

As reducing the number of entries is the only way to reduce the size of the file, a compact indexed database comprised of only 4054 entries was constructed for this study. In terms of the number of entries, it can clearly be seen that there is a massive reduction, compared with the raw MDDR dataset. In spite of using the raw MDDR dataset to create the indexed database, all the necessary information was addressed appropriately and the similarity searching process was fully accomplished on this one file.

The underlying process involved in the reduction in the number of entries in the indexed database is explained by the removal of duplicate LINGO Profiles and the mapping of the same index number which belongs to a particular LINGO Profile. The structure of an entry in the indexed database is shown in Fig. 8.

It can be seen that one LINGO Profile "sits" together with its respective index number on a single line. In contrast to raw MDDR dataset, each entry may be duplicating a portion of the same information (the index number or LINGO Profile) from the previous or the next entry of the file. This situation can be observed in the Fig. 9.



Fig. 7. Mapping of LINGO Profile with its Respective Index Number.



Fig. 8. The Structure of an Entry in the Indexed Database.

Fig. 9. Structure of the Entries in the Raw MDDR Dataset.

## V. Conclusions

The inverted indexing scheme has been highlighted in this study as there are several limitations when performing similarity searching using LINGO Profiles. The large raw MDDR dataset which is used in the conventional method to calculate the similarities requires a huge storage capacity, while at the same time increasing the time taken for one query compound to complete the whole process. The proposed method solves this problem by eliminating the redundant LINGO Profiles and multiple occurrences of the same index number. Despite this elimination, the important information associated with the compounds are preserved accordingly. In short, the proposed method makes it possible to process a huge dataset without the help of specialized hardware. In future, this scheme can be used to index a larger chemical database like ChEMBL which consist of more than 1 million compounds data.

### References

[1] R. Dolezal, V. Sobeslav, O. Hornig, L. Balik, J. Korabecny and K. Kuca, "HPC Cloud Technologies for Virtual Screening in Drug Discovery", Intelligent Information and Database Systems, pp. 440-449, 2015.

[2] X. Yan, C. Liao, Z. Liu, A. T. Hagler, Q. Gu and J. Xu, "Chemical Structure Similarity Search for Ligand-based Virtual Screening: Methods and Computational Resources", Current Drug Targets, vol. 17, no. 14, pp. 1580-1585, 2016.

[3] A. Cereto-Massagué, M. Ojeda, C. Valls, M. Mulero, S. Garcia-Vallvé and G. Pujadas, "Molecular fingerprint similarity search in virtual screening", Methods, vol. 71, pp. 58-63, 2015

[4] D. Clark and S. Pickett, "Computational methods for the prediction of 'drug-likeness'", Drug Discovery Today, vol. 5, no. 2, pp. 49-58, 2000.

[5] P. Willett, "Similarity-based virtual screening using 2D fingerprints", Drug Discovery Today, vol. 11, no. 23-24, pp. 1046-1053, 2006.

[6] D. Stumpfe and J. Bajorath, "Similarity searching", Wiley Interdisciplinary Reviews: Computational Molecular Science, vol. 1, no. 2, pp. 260-282, 2011

[7] Y. Hanyf and H. Silkan, "A queries-based structure for similarity searching in static and dynamic metric spaces", Journal of King Saud University - Computer and Information Sciences, 2018.

[8] A. Ciancetta and S. Moro, "Protein–Ligand Docking: Virtual Screening and Applications to Drug Discovery", In Silico Drug Discovery and Design, pp. 189-213, 2015.

[9] "Daylight Theory: SMILES", Daylight.com, 2018. [Online]. Available: http://www.daylight.com/dayhtml/doc/theory/theory.smiles.html. [Accessed: 01- Apr- 2018].

[10] D. Weininger, "SMILES, a chemical language and information system. 1. Introduction to methodology and encoding rules", Journal of Chemical Information and Modeling, vol. 28, no. 1, pp. 31-36, 1988.

[11] C. Steinbeck, Y. Han, S. Kuhn, O. Horlacher, E. Luttmann and E. Willighagen, "The Chemistry Development Kit (CDK): An Open-Source Java Library for Chemo- and Bioinformatics", Journal of Chemical Information and Computer Sciences, vol. 43, no. 2, pp. 493-500, 2003.

[12] D. Vidal, M. Thormann and M. Pons, "A Novel Search Engine for Virtual Screening of Very Large Databases", Journal of Chemical Information and Modeling, vol. 46, no. 2, pp. 836-843, 2006.

[13] P. Thiel, L. Sach-Peltason, C. Ottmann and O. Kohlbacher, "Blocked Inverted Indices for Exact Clustering of Large Chemical Spaces", Journal of Chemical Information and Modeling, vol. 54, no. 9, pp. 2395-2401, 2014.

[14] S. Dandapani, G. Rosse, N. Southall, J. Salvino and C. Thomas, "Selecting, Acquiring, and Using Small Molecule Libraries for High-Throughput Screening", Current Protocols in Chemical Biology, 2012.

[15] Z. Aung and S. Ng, "An Indexing Scheme for Fast and Accurate Chemical Fingerprint Database Searching", Lecture Notes in Computer Science, pp. 288-305, 2010.

[16] "Apache Lucene - Index File Formats", Lucene.apache.org, 2018. [Online]. Available:https://lucene.apache.org/core/3_0_3/fileformats.html#Inverted%20Indexing. [Accessed: 14- May- 2018].

[17] V. Anh and A. Moffat, "Inverted Index Compression Using Word-Aligned Binary Codes", Information Retrieval, vol. 8, no. 1, pp. 151-166, 2005.

[18] H. Yan, S. Ding and T. Suel, "Inverted index compression and query processing with optimized document ordering", Proceedings of the 18th international conference on World wide web - WWW '09, 2009.

[19] F. Hassen and G. Amel, "An efficient synchronous indexing technique for full-text retrieval in distributed databases", Procedia Computer Science, vol. 112, pp. 811-821, 2017.

[20] J. Zobel and A. Moffat, "Inverted files for text search engines", ACM Computing Surveys, vol. 38, no. 2, p. 6-es, 2006.

[21] P. Willett, J. Barnard and G. Downs, "Chemical Similarity Searching", Journal of Chemical Information and Computer Sciences, vol. 38, no. 6, pp. 983-996, 1998.

[22] D. Agrafiotis, J. Myslik and F. Salemme, "Advances in diversity profiling and combinatorial series design", Annual Reports in Combinatorial Chemistry and Molecular Diversity, pp. 71-92, 1999.

[23] N.H.A.H Malim, "Enhancing Similarity Searching," Information School, University of Sheffield, Sheffield, 2011

[24] H. Öztürk, E. Ozkirimli and A. Özgür, "A comparative study of SMILES-based compound similarity functions for drug-target interaction prediction", BMC Bioinformatics, vol. 17, no. 1, 2016.

[25] J. Grant, J. Haigh, B. Pickup, A. Nicholls and R. Sayle, "Lingos, Finite State Machines, and Fast Similarity Searching", Journal of Chemical Information and Modeling, vol. 46, no. 5, pp. 1912-1918, 2006.

[26] G. Hu, G. Kuang, W. Xiao, W. Li, G. Liu and Y. Tang, "Performance Evaluation of 2D Fingerprint and 3D Shape Similarity Methods in Virtual Screening", Journal of Chemical Information and Modeling, vol. 52, no. 5, pp. 1103-1113, 2012.

[27] G. Jayashree and V. Perumal, "Enhancing similarity-based query searching performance using self-organized semantic overlay networks", Proceedings of IEEE International Conference on Computer Communication and Systems ICCCS14, 2014. Available: 10.1109/icccs.2014.7068168 [Accessed 6 February 2019].

[28] D. Vidal, M. Thorman,. & M.Pons, "LINGO, an Efficient Holographic Text Based Method to Calculate Biophysical Properties and Intermolecular Similarities", Journal of Chemical Information and Computer Sciences, vol. 45, pp.386-393, 2014.

[29] H. Hentabli, N. Salim, A. Abdo and F. Saeed, "LINGO-DOSM: LINGO for Descriptors of Outline Shape of Molecules", Intelligent Information and Database Systems, pp. 315-324, 2013.

[30] M. Skinnider, C. Dejong, B. Franczak, P. McNicholas and N. Magarvey, "Comparative analysis of chemical similarity methods for modular natural products with a hypothetical structure enumeration algorithm", Journal of Cheminformatics, vol. 9, no. 1, 2017.

[31] R. Guha, K. Gilbert, G. Fox, M. Pierce, D. Wild and H. Yuan, "Advances in Cheminformatics Methodologies and Infrastructure to Support the Data Mining of Large, Heterogeneous Chemical Datasets", Current Computer Aided-Drug Design, vol. 6, no. 1, pp. 50-67

[32] P. Sharma, S. Salapaka and C. Beck, "A Scalable Approach to Combinatorial Library Design for Drug Discovery", Journal of Chemical Information and Modeling, vol. 48, no. 1, pp. 27-41, 2008.

[33] I. Haque, V. Pande and W. Walters, "SIML: A Fast SIMD Algorithm for Calculating LINGO Chemical Similarities on GPUs and CPUs", Journal of Chemical Information and Modeling, vol. 50, no. 4, pp. 560-564, 2010.

[34] F. Rinaldi, "Text Mining Technologies for Database Curation", Proceedings of the International Conference on Knowledge Discovery and Information Retrieval, 2014.

[35] R. Nasr, R. Vernica, C. Li and P. Baldi, "Speeding Up Chemical Searches Using the Inverted Index: The Convergence of Chemoinformatics and Text Search Methods", Journal of Chemical Information and Modeling, vol. 52, no. 4, pp. 891-900, 2012.

[36] T. Kristensen, J. Nielsen and C. Pedersen, "Using Inverted Indices for Accelerating LINGO Calculations", Journal of Chemical Information and Modeling, vol. 51, no. 3, pp. 597-600, 2011.

[37] E. S. D. Moura, "Text Indexing Techniques", Encyclopedia of Database Systems, 4084–4088,2018.

[38] J. Hert, M. Keiser, J.J. Irwin, T.I. Oprea, and B.K. Shoichet, "Quantifying the Relationship among Drug Classes", Journal of Chemical Information and Modelling, vol 48, pp. 755-765,2008.

# A Framework for Iris Partial Recognition based on Legendre Wavelet Filter

Muktar Danlami[1], Sapiee Jamel[2], Sofia Najwa Ramli[3], Mustafa Mat Deris[4]
Faculty of Computer Science and Information Technology
Universiti Tun Hussein Onn Malaysia
Parit Raja, Malaysia

*Abstract*—**An increasing need for biometrics recognition system has grown substantially to address the issues of recognition and identification especially in highly dense areas such as airport, train stations and for financial transaction. Evidences of these can be seen in some airports and also the implementation of these technologies in our mobile phones. Among the most popular biometric technologies include facial, fingerprints and iris recognition. The iris recognition is considered by many researchers to be the most accurate and reliable form of biometric recognition, because iris can neither be surgically operated with a chance of losing slight nor change due to ageing. However, presently most iris recognition system available can only recognize iris image with frontal-looking and high-quality images. Angular image and partially capture image cannot be authenticated with existing method of iris recognition. This research investigates the possibility of developing a framework for recognition partially captured iris image. The research also adopts the Legendre wavelet filter for the iris feature extraction. Selected iris images from CASIA, UBIRIS and MMU database were used to test the accuracy of the introduced framework. A threshold for the minimum iris image required was established.**

*Keywords*—*Iris recognition; partial recognition; wavelet; Legendre wavelet filter; biometric*

## I. Introduction

The increasing need for a reliable means for an identification and verification system cannot be over emphases [1]. The world population and the need for identifying or verifying people in highly dense areas force the evolution of the use of biometric technologies as alternative and more effective means of access control [2].

The word biometrics is a two combine word of the Greek words bio and metric, which is "life meaning bio and measurement meaning metric". Biometric technology is defined as any technique that can use measurable physiological or behavioral characteristics to discriminate one person from another [3]. Common physiological biometric traits include iris, fingerprints, facial, hand geometry, and retina images. Whereas, common behavioral biometric traits include: voice recording, signature, and keystroke rhythms. It is noted that behavioral biometrics, in general, include a physiological component as well [4].

Although all biometric systems work in the same manner, the first process is enrollment in which each new user is registered into the database. Information about a specific characteristic of the individual is captured. This information is usually passed through an algorithm that turns the information into a template that the database stores. Note that it is the template that is maintained in the system, but not the original biometric measurement as many people may suspect. Compared with the original measurement of the biometric trait, the template has a tiny amount of information; it is no more than a collection of numbers with little meaning except to the biometric system that produced them. When a person needs to be recognized, the system will take the appropriate measurement, translate this information into a template using the same algorithm that the original template was computed with, and then compare the new template with the database to determine if there is a match, and hence, either verification or identification [5].

Today fingerprint and facial recognition system are one of the most used biometric recognition system. Both the fingerprint biometric and the facial recognition system are used in the public domains such as airport, train station and also our financial institution such as banks and Automated Teller Machine (ATM) [6]. However, both the fingerprint and the facial recognition are facing some setbacks. For the fingerprint recognition, the system users need to scan their finger on a fingerprint scanning device, this makes it difficult to authenticate someone with his knowledge and also frequent use of the scanning device often makes the scanning device dirty thus fails during recognition.

Iris recognition has been verified to be one of the most accurate and reliable biometrics authentications, unlike facial recognition, and fingerprint. The facial recognition has great problem due to the fact that the human faces changes over time due to growth development in human nature. The fingerprint unlike the facial recognition does not change for as long as we leave however face setback such as the need for the authenticated individual to scan his or her hand to the scanning device, this make it difficult to authenticate an individual without his or her knowledge, sometimes the scanning device maybe dirty [7]. The identified problems make iris recognition an alternative as the best biometric authentication; iris is neither affected by age nor requires an individual to have a contact with its scanning device

Presently, iris recognition methods can work very well with frontal-looking and high-quality images. Daugman's 2D Gabor wavelet approach has been tested and evaluated using huge databases, such as the CASIA database, UBIRIS database, and MMU database among others, with over 600,000 iris images

with over 200 billion comparisons [8]. However, most existing methods are not designed for non-cooperative users and cannot work with off-angle or partially captured iris images. Recognition can be quite good if canonical poses and simple backgrounds are employed, but changes in illumination and angle create challenges. Recognizing an individual with incomplete or partially captured images in biometric technology continues to be an important challenge today. Despite the advancement made in fingerprint identification techniques, little or not much have been achieved for that of iris recognition. Partially captured image or images with noise or occlusion is a well-known research problem, and many researchers have tried to address the problems in a different capacity.

## II. DATABASES

We selected four different databases to test our method, namely; Chinese Academy of Sciences Institute of Automation (CASIA) [9], University Beira IRIS (UBIRIS) [10], and Multimedia University (MMU) database. The selected database was based on the most frequently used database for the iris recognition algorithm. However, to show the effect of partial recognition, there is a need for the dataset to be carefully selected. We only selected images that are partially captured. However, for registering the iris images to the database, here, we also selected best-captured images. For each subject or eye image, 2-10 images are selected, depending on the availability of the partially captured image of the particular subject or eye.

The iris recognition was implemented with the selected database. The selected databases include CASIA v4 database, UBIRIS v2 database, MMU v2 database, and IITD database. The CASIA v4 database consist of subset namely, CASIA-IRIS-interval, CASIA-IRIS-twins, CASIA-IRIS-distance, CASIA-IRIS-thousand, CASIA-IRIS-syn, however only the CASIA-IRIS-interval and CASIA-IRIS-distance were used. The CASIA-IRIS-interval consists of 249 subjects with a total of 2639 number of iris images, but only 994 images were used from 249 subject. The CASIA-IRIS-distance consists of 142 subjects with a total of 2567 number of images, but only 710 images were used from 142 subject. The UBIRIS v2 database consists of 261subjects and 522 irises with a total of 11102 images, but only 783 were used from 261 subject. The MMU database consists of 100 subjects and 200 irises with a total of 10000 number of image, but only 300 images were selected from 100 subjects as in Table I.

Fig. 1 is a sample of some best-captured eye image from the MMU database; we roundly select then to show how they look. While Fig. 2, is a sample of some partially captured eye image from the MMU database. They are partially captured because either the subject eyes are partially closed or the subject is looking sideway, or the eyelashes of the subject partially closed the eye image.

TABLE I. INFORMATION OF THE SELECTED DATASET

| s/n | Database | Subject | Images |
|---|---|---|---|
| 1 | CASIA-IRIS-interval | 249 | 994 |
| 2 | CASIA-IRIS-distance | 142 | 710 |
| 3 | UBIRIS v2 | 261 | 783 |
| 4 | MMUv2 | 100 | 300 |



Fig 1. Best-Captured Image from MMU Iris Database.



Fig 2. Partially-Captured/Noisy Iris Images from MMU Database.



Fig 3. Some Partially-Captured/Noisy Iris Images from UBIRIS Database.

Fig 4.    Some Best-Captured Image from UBIRIS Iris Database.

Fig. 4 is sample of the best-captured eye image from the UBIRIS database; we roundly select them to show how they look. While Fig. 3 are the partial captured eye image from the UBIRIS database. They are partially captured because either the subject eyes are partially closed or the subject is looking sideway, or the eyelashes of the subject partially closed the eye image.

## III.    METHODOLOGY

The idea here is to find a threshold for which iris can be recognized partially. That is, to find the smallest among of size of iris required to authenticate the subject. The research will consider the normalized iris image at for different percentage; 50 percent, when the normalized iris image is divided into two parts, 25 per cent, when the normalized iris image is divided into four parts, 16.5 per cent, when the normalized iris image is divided into six parts, and 12.5 per cent when the iris image is divided into eight equal parts. With this for different sizes, we find the minimum size of normalized iris required for the recognition process. The processes for the recognition include segmentation, normalization, feature extraction and matching as in Fig. 5.

The first stage of the recognition is the acquisition of the image, for the stage we intend to use the available database online. Some of the databases need some adjusting. Also, the iris images are in different resolution and there is need for a standard size of resolution across the database images. The UBIRIS database images, for example, need to be converted into greyscale image, for others such as CASIA and MMU are all in greyscale. Fig. 6 shows the converted UBIRIS image from colored to a greyscale image.

### A.  Segmentation

For most of the database, the conversion of the image from colored to greyscale is not needed. The process usually starts with segmentation. The iris image is selected from the eye image as in Fig. 7.

### B.  Normalization

Next is to normalized the segmented iris image, here, the rubber sheet mode was used to achieve this function. This is shown in Fig. 8.



Fig 5.    Framework for Partial Iris Recognition.



Fig 6.    Iris Image from UBIRIS Database Converted to Greyscale.



Fig 7.    Iris Segmentation of the Eye Image.

Fig 8.    Iris Normalization.

## C. Feature Extraction

The feature extraction follows after the normalization. Feature that distinguish the iris image are enhance using the Legendre wavelet filter. Following the approach of [11], the Legendre wavelet filter can be define as in Equation (1).

$$
\Psi_{n_1,m_1,n_2,m_2}(x,y) = \Psi_{e,n_1,m_1,n_2,m_2}(x,y)
$$

$$
= \sqrt{\left(m_1 + \frac{1}{2}\right)\left(m_2 + \frac{1}{2}\right)2^{\frac{(k_1+k_2)}{2}}} \quad (1)
$$

$$
\times \ P_{m_1}\left(2^{k_1}x - \hat{n}_1\right)P_{m_2}\left(2^{k_2}y - \hat{n}_2\right)e^{j2\pi(u_1x+v_1y)}
$$

where

$$
e^{j2\pi(u_1x+v_1y)} = \cos[2\pi(u_1x + v_1y)] \quad (2)
$$

$$
+ \ j\sin[2\pi(u_1x + v_1y)]
$$

$u_1$ and $v_1$ are the fundamental frequencies in X and Y direction $m = 0,1,\dots M$ and $n = 0,1,\dots 2^{k-1}$, the coefficient $\sqrt{\left(m_1 + \frac{1}{2}\right)\left(m_2 + \frac{1}{2}\right)}$ is for the orthonomality and the Pm is Legendre polynomial

Generally, image features are pieces of information that describes an image or a part of an image as in Fig. 9. However, in pattern recognition feature is a piece of information which is relevant for solving the computational task related to a certain application. Feature extraction begins from an initial set of measured data and builds derived values feature intended to be informative and non-redundant, helping the subsequent learning and generalizing steps, and in some cases leading to better human interpretations.



Fig 9.    Iris Feature Extraction.

## D. Matching

Lastly the recognition is concluded by the matching, were the unique feature extracted  from the iris image is been compare with the corresponding iris image in the database for verification or the unique feature are searched across the saved feature in the database until a match is found for identification.

## IV.  EVALUATION PARAMETER

False Acceptance Rate (FAR): FAR is the frequency of fraudulent access to imposter claiming identity. This statistic is used to measure biometric performance when operating in the verification mode. A false accepts occurs when the query template of an individual is incorrectly matched to existing biometric template of another individual.

False Rejection Rate (FRR): FRR is the frequency of rejections relative to people who should be correctly verified. This statistics is used to measure biometric performance when operating in the verification mode. A false reject occurs when an individual is not matched correctly to his/her own existing biometric template.

Genuine Acceptance Rate (GAR): GAR is the frequency of genuine access with respect to overall number of attempts.

## V.  RESULT AND DISCUSSION

The Legendre wavelet filter was implemented using Matlab R2015 installed on a Window 7 professional desktop computer, Intel core i7. We considered the Legendre wavelet filter at three different orders. The experimental setting is introduced, including the selected database, parameter setting and performance evaluation. Then, to study the effect of the proposed partial method of the iris code production, comparisons are made between the performances of the iris codes produced by an implementation of traditional iris code generation method.

The iris code generated was tested with the selected images in CASIA-IRIS-interval and the result is as in Table II. The lowest accuracy was achieved at 50% and the highest was achieved at 16.5%. The FAR has its lowest at 50% and its highest at 12.5% while the FRR has its highest at 50% and its lowest at 16.25%. The graphical representation of the accuracy of the CASIA-IRIS-interval is shown in Fig. 10.

TABLE II.    RESULT OF THE PARTIAL RECOGNITION WITH CASIA-IRIS-INTERVAL

| PERCENTAGE OF THE IRIS IMAGE | FAR % | FRR % | GAR % | ACURACY % |
|---|---|---|---|---|
| 50% | 5.48 | 15.89 | 87.05 | 87.05 |
| 25% | 6.45 | 14.26 | 88.26 | 88.26 |
| 16.5% | 6.75 | 13.25 | 92.25 | 92.25 |
| 12.5% | 6.82 | 13.56 | 91.95 | 91.95 |



Fig 10.   Accuracy Result of the Partial Recognition with CASIA-IRIS-Interval.

TABLE III.    RESULT OF THE PARTIAL RECOGNITION WITH CASIA-IRIS-DISTANCE

| PERCENTAGE OF THE IRIS IMAGE | FAR % | FRR % | GAR % | ACURACY % |
|---|---|---|---|---|
| 50% | 7.59 | 16.49 | 84.05 | 84.05 |
| 25% | 8.95 | 15.36 | 85.26 | 85.26 |
| 16.5% | 7.95 | 14.28 | 86.25 | 86.25 |
| 12.5% | 6.42 | 18.59 | 84.95 | 84.95 |



Fig 11.    Accuracy Result of the Partial Recognition with CASIA-IRIS-Distance.

The iris code generated was tested with the selected images in CASIA-IRIS-distance and the result is as in Table III. The lowest accuracy was achieved at 50% and the highest was achieved at 16.5%. The FAR have lowest at 50% and highest at 12.5% while the FRR has highest at 12.5% and lowest at 16.25%. The graphical representation of the recognition accuracy is in Fig. 11.

The iris code generated was tested with the selected images in UBIRISv2 and the result is as in Table IV. The lowest accuracy was achieved at 12.5% and the highest was achieved at 16.5%. The FAR have lowest at 12.5% and highest at 16.5% while the FRR has highest at 12.5% and lowest at 25%. The graphical representation of the recognition accuracy is in Fig. 12.

TABLE IV.    RESULT OF THE PARTIAL RECOGNITION WITH UBIRISV2

| PERCENTAGE OF THE IRIS IMAGE | FAR % | FRR % | GAR % | ACURACY % |
|---|---|---|---|---|
| 50% | 4.59 | 20.69 | 74.05 | 74.05 |
| 25% | 3.95 | 19.86 | 74.36 | 74.36 |
| 16.5% | 4.95 | 24.48 | 74.95 | 74.95 |
| 12.5% | 1.42 | 25.19 | 73.55 | 73.55 |



Fig 12.    Accuracy Result of the Partial Recognition with UBIRISv2.

The iris code generated was tested with the selected images in MMUv2 and the result is as in Table V. The lowest accuracy was achieved at 50% and the highest was achieved at 16.5%. The FAR have lowest at 12.5% and highest at 50% while the FRR has highest at 50% and lowest at 12.5%. The graphical representation of the recognition accuracy is in Fig. 13.

TABLE V.    RESULT OF THE PARTIAL RECOGNITION WITH MMUV2

| PERCENTAGE OF THE IRIS IMAGE | FAR % | FRR % | GAR % | ACURACY % |
|---|---|---|---|---|
| 50% | 3.50 | 13.69 | 92.05 | 92.05 |
| 25% | 3.25 | 12.86 | 92.96 | 92.96 |
| 16.5% | 2.50 | 10.48 | 94.45 | 94.45 |
| 12.5% | 1.50 | 11.19 | 93.55 | 93.55 |



Fig 13.    Accuracy Result of the Partial Recognition with MMUv2.

## VI. Conclusion and Future Work

The main focus of the research was to try the iris recognition with a partially capture image and to also do the recognition partially. So the idea was to find a threshold that can determine the minimum amount of iris region required to identify an individual. Presently the method of partial recognition is applied in fingerprint recognition especially with fingerprint integrated with the mobile hand phone, whereby any part of your fingerprint can be used for the recognition.

Based on the experiment that was carried out, it shows that the partial recognition can also be applied with the iris. Substantially the iris can be recognition with as low as only 12.5% of the iris image. However, best results were achieved with the iris image at 16.5%.

Some of the future work of the research is to create a database that will have only iris images that are partially captured. Providing the database will help standardize the process of the proposed framework evaluation.

Secondly more feature extraction technique can be introduced for better extraction of the iris feature.

### References

[1] K. S. Kumar, K. N. N. S. Ram, K. Kiranmai, and S. S. Harsha, "Denoising of Iris Image Using Stationary Wavelet Transform," in 2018 Second International Conference on Inventive Communication and Computational Technologies (ICICCT), 2018, pp. 1232–1237.

[2] K. Muralidharan, P. Niehaus, and S. Sukhtankar, "Building state capacity: Evidence from biometric smartcards in India," Am. Econ. Rev., vol. 106, no. 10, pp. 2895–2929, 2016.

[3] M. Nabti and A. Bouridane, "An effective and fast iris recognition system based on a combined multiscale feature extraction technique," Pattern Recognit., vol. 41, no. 3, pp. 868–879, 2008.

[4] S. Patil, S. Gudasalamani, and N. C. Iyer, "A survey on Iris recognition system," in International Conference on Electrical, Electronics, and Optimization Techniques, ICEEOT 2016, 2016, pp. 2207–2210.

[5] W. W. Boles and B. Boashash, "A Human Identification Technique Using Images of the Iris and Wavelet Transform," IEEE Trans. Signal Process., vol. 46, no. 4, pp. 1185–1188, 1998.

[6] B. Mary Reni, "Iris recognition based age estimation in security systems using Canny edge detection," Res. J. Pharm. Biol. Chem. Sci., vol. 6, no. 5, pp. 349–357, 2015.

[7] A. Khatun, A. K. M. F. Haque, S. Ahmed, and M. M. Rahman, "Design and implementation of iris recognition based attendance management system," in 2nd International Conference on Electrical Engineering and Information and Communication Technology, iCEEiCT 2015, 2015.

[8] J. Daugman, "New methods in iris recognition," IEEE Trans. Syst. Man, Cybern. Part B Cybern., vol. 37, no. 5, pp. 1167–1175, 2007.

[9] T. Tan, Z. He, and Z. Sun, "Efficient and robust segmentation of noisy iris images for non-cooperative iris recognition," Image Vis. Comput., vol. 28, no. 2, pp. 223–230, 2010.

[10] H. Proenc, "Iris Recognition : On the Segmentation of Degraded Images Acquired in the Visible Wavelength," vol. 32, no. 8, pp. 1502–1516, 2010.

[11] D. Muktar, S. Jamel, S. N. Ramli, and M. M. Deris, "2D legendre wavelet filter for iris recognition feature extraction," in Proceedings of the 3rd International Conference on Cryptography, Security and Privacy, 2019, pp. 174–178.

# The Implementation of Software Anti-Ageing Model towards Green and Sustainable Products

Zuriani Hayati Abdullah[1], Jamaiah Yahaya[2], Siti Rohana Ahmad Ibrahim[3], Sazrol Fadzli[4]
Faculty of Information Science and Technology,
Universiti Kebangsaan Malaysia Bangi, Selangor, Malaysia

Aziz Deraman[5]
School of Informatics & Applied Mathematics,
Universiti Malaysia Terengganu, Kuala Terengganu,
Malaysia

*Abstract*—Software ageing is a phenomenon that normally occurs in a long running software. Progressive degradation of software performance is a symptom that shows software is getting aged and old. Researchers believe that the ageing phenomenon can be delayed by applying anti-ageing techniques towards the software or also known as software rejuvenation. Software ageing factors are classified into two categories: internal and external factors. This study focuses on external factors of software ageing, and are categorized into three main factors: environment, human and functional. These three factors were derived from empirical study that been conducted involving fifty software practitioners in Malaysia. The anti-ageing model (SEANA model) is proposed to support in preventing the software from prematurely aged, thus prolong its usage and sustainable in their environment. SEANA model is implemented in collaboration with a government agency in Malaysia to verify and validate the model in real environment. The prototype of SEANA model was developed and applied in the real case study. Furthermore, the anti-ageing guideline and actions are suggested for ageing factors to delay the ageing phenomenon in application software and further support the greenness and sustainability of software products.

*Keywords*—*Software ageing factor; ageing prevention; software anti-ageing model; SEANA model; SeRIS Prototype System; Green And Sustainable Product; Emprirical Study*

## I. INTRODUCTION

As the increase of dynamic software requirements nowadays from users and stakeholders, software development process is becoming more complex and resulting the degradation of software performance and software quality [1][2]. If this happens to the software which is operating in certain environments, it may get aged prematurely and no longer relevant in their environment. Users may refuse to use the software anymore because it does not fulfil and satisfied the requirements and expectation. Progressive degradation of software performance, such as software crash or hang and accumulation of software errors are reported as the phenomenon of software ageing. The ageing of the software is caused by two factors which are by the changes that have been made throughout its execution and also cause by the failure to adapt with the dynamic environment [3].

Software ageing may occur when there are accumulation of errors or software failure throughout its execution. However, it does not affect or change the functionality of a software, but its effects on the time responsiveness of the software and user

satisfaction over the software [4]. Software failure is closely related to the software downturn during its life cycle. The problem led to a progressive decline in software performance, and caused the software to not function properly. This degradation process is called software ageing [5][23]. Previous studies revealed that we could slowed down software ageing by identifying the influential factors of the ageing phenomenon. There are two types of ageing influential factors which are internal and external factors. However, there are very few studies focuses on external factors [4] [8] that are closely related to software quality in application software. Based on our initial investigation has discovered that some applications get old and aged as early as three years and thus forced the users to not used the application anymore. In this scenario, the application ages prematurely. Currently there is no software anti-ageing mechanism or guideline to assist users or developers to measure and guarantee the software still relevant and young in their operating environment. However, in general, sustainability in software is associated with the ability to operate in a longer time and viable in their environment.

This paper discusses the research background and related works in Section II, the empirical study conducted in this research is presented in Section III. Section IV of this paper presents the development of the anti-ageing model while the validation and the implementation of the model are discussed in Section V. The anti-ageing actions are introduced and presented in Section VI, and lastly concludes with a discussion and conclusion.

## II. BACKGROUND STUDY AND RELATED WORKS

The phenomenon of ageing is applicable in software which is operating in certain environments. By identifying the significant factors and causes of software ageing, it can ageing effect to environment, organization delay and help in preventing the occurrence of ageing phenomenon. This process may be stated as anti-ageing or a rejuvenation process of a software product. Currently in software engineering, the rejuvenation process of a software is considered as one of the mechanism for handling faults tolerant or failure in the long running software [6][21]. Hence, it is essential to find solutions on how to prevent or slow down the ageing process in order to maintain the relevancy of the software and still meet their business requirement. Next section will discuss software ageing issues and software and economics.

## A. *Software Ageing Issues and Effects*

Software ageing has been studied by several researchers as early as in 1990s. Previous studies indicated that issues and challenges related to software ageing in computer software area are normally associated with accumulation of undetermined threads or failure, data corruption, memory related problem such as memory leaking and bloating, data-files fragmentation, residual defects, unreleased files lock, memory lack and overruns [7][8].

Nowadays, software ageing does not only associate to computer software or system software but also being investigated that relates to mobile application such as android and windows mobile application [9][10][11]. Mobile was running for a longer time without rebooted or shut down compared to computer. Therefore, software ageing in mobile devices lead to an extensive challenge to ultimate the user experience and satisfaction. Software ageing issues are not limited to mobile and computer software only but now extended to the ageing phenomenon in cloud computing environment [5][12]. This shows that software ageing is a relevant issue to be explored and investigated further.

Good quality software will delay in the occurrence of ageing and prolong their usage and relevancy. This is supported by the previous researchers who believed that by maintaining good quality of software, it can somehow prevent or minimize the error or failure and thus results in user satisfaction when using the software [13]. Good quality software is referred to the technical behaviour of the software and end user's perspective towards the software, which measures the satisfaction and fulfil expectation of the software. The study done by Yahaya et al. [13] reveals that with these two quality characteristics and measurements will maintain and ensure the software are relevant more longer of time in their operating environment.

In addition to good quality software, Matias, Trivedi and Maciel [14] claimed that software maintenance must be implemented systematically to ensure the optimum quality throughout software life cycle. In software engineering, there are four main maintenance activities which are corrective, adaptive, preventive and perfective. These activities may control and delay the ageing progress of software. For instance, preventive maintenance can help to slow down the occurrence of failures determinable to this cause.

Previous studies revealed that software ageing might gave negatives effects or influences in various aspects [2][4]. For examples its gives drawback to organizational level. Ageing effects on the operating system at resource level such as non-released memory, round-off and data file fragmentations and also debug errors. These delay the work and schedule because of slow time responsiveness on running application.

## B. *Software Anti-Ageing and Rejuvenation*

Software ageing is irresistible manifestation, but there are few studies on how to deal with ageing phenomenon in software. As mentioned earlier, software ageing can be delayed by adopting two approaches which are through software anti-ageing and software rejuvenation. There is a difference between these two approaches where software rejuvenation is used once software is detected ageing, while anti-ageing software can be applied before software becomes old in order to delay the ageing process [15].

According to [12], software rejuvenation process is not a complex task but a very effective technique in increasing the availability of a software by rebooting and refreshing the software. The rejuvenation technique is used to revive ageing software after the ageing status is detected [16][24]. Cotroneo et al. [17] suggests that software maintenance activities are considered as the anti-ageing process where these activities are used to delay the ageing of software. There are four maintenance activities that can be used as software anti-ageing techniques, which are adaptive, corrective, perfective and preventive maintenance [14]. Adaptive maintenance refers to adapting to a new environment or sometimes refers to adapting to new requirements. Corrective maintenance refers to maintenance to repair fault and perfective maintenance refers to perfecting the software by implementing new requirements. In other case it refers to maintaining the functionality of the system, improving its structure and its performance. While preventive maintenance involves performing activities to prevent the occurrence of errors. It tends to reduce the software complexity thereby improving program understandability and increasing software maintainability. In this maintenance activity comprises documentation updating, code optimization, and code restructuring. This is also known as software re-engineering.

It is crucial to prevent software ageing because it not only effects software system, but also effects user and the universe in general. There was a fatal incident that happened about twenty-six years ago where twenty eight soldiers dead and hundreds people were injured because of software failed to detect an Iraqi Scud missile and it strucked the American army barracks [18]. Based on various issues on software ageing discussed in this paper has motivated and led us to explore more on software ageing phenomenon.

## C. *Green and Sustainable Software Product*

Green software as part of information technology (IT) covers environmental sustainability, economic energy efficiency and total cost of ownership, which includes the cost of disposal and recycling. It also refers to the application of IT to create the energy efficient and environment to maintain successful business processes and practices [25]. It incorporates three dimension that are greening IT systems and usage, using IT to support environmental sustainability and also using IT to create green awareness in a way to improve environmental sustainability [26].

Green and sustainable software can be defined as a software with direct or indirect negative impact on economy, society, human being and environment that result from development, deployment and usage of the software. It should have minimal and positive effect towards sustainable development [27]. In addition, Erdelyi [28] defines green software as the development and operation of the software that produce minimal disposal and waste as possible [28].

## III. EMPIRICAL STUDY

This section disscuses on the empirical study that was conducted in Malaysia. The empirical study was conducted through a survey to identify the awareness and acceptance of the software ageing issues and concerns among software practitioners. Second objective of the the survey was to examine and classify the ageing factors identified from literature study that influence the ageing of software. The survey was carried out involving respondents from agencies in public and private sectors. The respondents were chosen using purposive sampling in the group of software practitioners in Malaysia. The background of the respondent may came from diverse types of organization background as shown in Table I.

Questionnaires were distributed by two methods physical questionnaire and online questionnaire (through Google form). Fifty respondents have participated and responded to this survey. Table I shows the respondents's background and their organisations. Majority of the respondent are from service, public administator and ministry department (42%). Respondents are also came from various other backgrounds such as Computer/Security System (8%), Computer Engineering/Design (6%), Software Development (2%), Training/Education/ Consultancy(6%), Internet based organization (20%), Internet based Business e-commerce/web hosting (4%), Healthcare (2%), Intergration system (2%) and others (10%). Fig. 1 illustrates that most of the respondents (48%) have three to ten years working experience, 46% of the respondents have working experience less than three years, while only 6% have working experience from eleven to more than twenty years.

TABLE I.        BACKGROUND OF ORGANISATION FUNCTION

|   | Background of Organisation/ Function | Percentage (%) |
|---|---|---|
| 1 | Service/Public Administrator/Ministry Dept. | 42% |
| 2 | Software Development | 2% |
| 3 | Computer/Security System | 8% |
| 4 | Computer Engineering/Design | 6% |
| 5 | Telecommunicating/Networking | 2% |
| 6 | Internet based Business (ASP)/ e-commerce/ web hosting | 20% |
| 7 | Healthcare | 2% |
| 8 | Integration System | 2% |
| 9 | Training/Education/ Consultancy | 6% |
| 10 | Others | 10% |
| Total : | 100% | |



Fig 1.     Service Period.



Fig 2.     Software Ageing Awareness.

TABLE II.        SOFTWARE AGEING FACTORS

| Factor | Metrics | Mean | % |
|---|---|---|---|
| Functionality | Time responsiveness | 3.62 | 72.4 |
| | Software are unable to meet the user's needs | 3.6 | 72.0 |
| | Failure to function as user's intended | 3.58 | 71.6 |
| | Progressive performance degradation | 3.52 | 70.4 |
| | Software is no longer relevant | 3.52 | 70.4 |
| | High frequency of software error | 3.4 | 68.0 |
| | Failure to upgrading the functionality | 3.34 | 66.8 |
| | Failure to get support | 3.32 | 66.4 |
| | User Interface (UI) | 2.98 | 59.6 |
| Environment | Dynamic to environment changes | 3.42 | 68.4 |
| | Lack of cost for software maintenance | 3.42 | 68.4 |
| | Software not compatible with current hardware technology | 3.4 | 68.0 |
| | Hardware changes | 3.36 | 67.2 |
| | Business process changes | 3.22 | 64.4 |
| | Business need | 3.12 | 62.4 |
| | Changes of software technology | 3.06 | 61.2 |
| Human | Lack of expertise in upgrading and maintaining software | 3.4 | 68.0 |
| | Weak software quality practices among practitioners | 3.32 | 66.4 |
| | Inefficient software management by management team | 3.14 | 62.8 |
| | Software is not user friendly | 3.08 | 61.6 |

Based on the first objective of the empirical study, the analysis shows that most of the respondents (64%) are not aware and realise the presence of the phenomenon of software ageing in their operational software. Nevertheless, they agree that they have experienced in the scenario of ageing occurrence in their daily software operation (refer to Fig. 2). Majority of respondents reveal that there are no standard mechanism or policy that can be referred to measure the quality states of the software. This finding discovers the inadequate of awareness among software practitioners and software developers in software quality and related aspects. Inadequate of awareness in software quality practices can be one of possible factors that influence to the occurrence of software ageing.

This study also determines other possible factors that might influence and contribute to the software ageing phenomenon for application software. The initial discover and indentification of metrics were done through literature study and brainstorming approach among research team and experts. They include software engineering experts, software engineers and academicians in Malaysia. Twenty initial metrics have been identified and verifies through this empirical study. Table II shows the findings which verified the ageing factors. From the findings, the software ageing factors are further mapped and classified into three categories which are functionality, environment, and human.

The percentages shown in Table II are obtained from the score given by the respondents of this survey. The higher percentage means that the factor has high influence and crucial toward the ageing process. The study reveals that functionality factor contributes the highest percentage (76.4%), environment (65.72%) and human factor (64.7%) (refer Table II). This indicates the importance of these factors based on respondents perceive and perspective toward the association of ageing occurrence and phenomenon in software environment. Time responsiveness shows the highest score proving that slow response of the software will contribute to the ageing of application software. Software ageing may result in slowing the system performance in performing tasks and therefore, contributes to user dissatisfaction towards the software. Based on respondent's feedback, user interface metric obtains the lowest percentage which is 59.6%. However, software that have dull and unattractive user interface and not user friendly can also lead to software ageing.

## IV. SEANA: THE SOFTWARE ANTI-AGEING MODEL

The development of anti-ageing model (or SEANA model) is based on the findings from the previous empirical study discussed in previous section as well as literature findings. The SEANA model comprises of the ageing factors and metrics (the ageing instrument), software ageing assessment process and reporting process. The last two components include the measurement algorithm, ageing levels, anti-ageing guideline and actions. SEANA model is demonstrated in Fig. 3. The following sub sections explain each of the components in this model.



Fig 3. Software Anti-Ageing Model (SEANA).

*A. Ageing Factor*

The ageing factors comprises in SEANA model are designed to be used and applied by different users. Therefore, a proper and systematic instrument is needed. The objective of the instrument is to measure the ageing level of the targeted software. It can be used by the organisation or company, the stakeholder or the owner of the application software who wants to investigate the ageing status of the software. The instrument is designed in three categories of software ageing factors (which are functional, environment and human) and twenty metrics that have been identified and verified in the empirical study discussed in previous section.

All metrics in the instrument are formulated and designed to be answered by the respondents by giving scores between 1 to 4 (Likert scale of four).

*B. Assessment Process*

The second component of this model is the process of assessment. The assessment process component aims to do the following tasks:

- To measure the score for each categories of ageing factors that have been indicated in the instrument.

- To measure the ageing score of the overall factors.

- To map the score obtained in the assessment with the ageing level.

*C. Assessment Report*

The third component is the report. This component is the report generated after the completion of the assessment process. The report consists of the ageing index level, which are defined as young, semi-old and very old. If the assessment result shows that the ageing score is very old, the suggestion of anti-ageing action will be included in the report. The anti-ageing action will be suggested to the tested software according to highest scores of ageing factors obtained by the respondent during software assessment. Table III shows the ageing classification level based on the score obtained in the assessment. The ageing levels or classifications are adopted from [19] by defining ageing level as big ageing and little ageing. However, for this study purposes and compatibility, we classify software ageing into three level which are young, semi-old and very old.

TABLE III. CLASSIFICATION OF SOFTWARE AGEING

| Class | Score (%) |
|---|---|
| Young | 68- 100 |
| Semi-old | 35 - 67 |
| Very Old | 1-34 |

## V. VALIDATION AND IMPLEMENTATION IN CASE STUDY

This section discussed on the validation and implementation of SEANA model in real case study.

*A. The Case Study*

The case study has been conducted in order to validate and implement the proposed model. We conduct a case study in one of the semi-government agency in east coast of Malaysia. We choose this agency (referred as KET) because they develop their own application in-house and have their own maintenance team to monitor their system application state. The aim of this case study was to assess three application systems that operated in their environment (referred as App1, App2 and App3 System).

*1) App1 system*: App1 System is an information system that manages the geographical data for a particular region. The system function similar to Google Maps but only stores and keeps the data in the east coast area of Malaysia in order to assist KET to find the rural area for ITC development purposes. Table IV shows the result obtained from the case study for App1 system.

TABLE IV. ASSESSMENT RESULTS FOR APP1 SYSTEM

| Factors | Average Score (1-4) | Percentage |
|---|---|---|
| Functionality | 1.78 | 44.5% |
| Environment | 1.86 | 46.5% |
| Human | 2.00 | 50.0% |

Based on the result, the average cumulative score for this software product is 47% which is mapped into the ageing index level of Semi-old. A post assessment meeting, and review with the owner of the system discovered that the App1 System performed normally slow to retrieve geographical information such as images or maps. However, the functionality of the software is still in good condition because they practice the software maintenance activities and continuously upgrading the system regularly in order to achieve the user satisfaction.

*2) App2 system*: The second selected system to be tested in this case study was App2 System. It is a document management system that allows KET staff to manage document online such as letters, memos and filing. Table V shows the result obtained from this assessment based on the three factors.

TABLE V. ASSESSMENT RESULTS FOR APP2 SYSTEM

| Factor | Average Score (1-4) | Percentage |
|---|---|---|
| Functional | 2.56 | 64% |
| Environment | 2.71 | 67.8% |
| Human | 2.75 | 68.8% |

*3) App3 system*: The third assessment of application software is App3 system. App3 System helps KET to monitor and manage the development of Rural Transformation Centre (RTC). The result of this assessment is shown in Table VI.

TABLE VI.     ASSESSMENT RESULTS FOR APP3 SYSTEM

| Factor | Average Score (1-4) | Percentage |
|---|---|---|
| Functional | 3.22 | 80.5% |
| Environment | 3.00 | 75.0% |
| Human | 3.00 | 75.0% |

Based on the computed result, App3 System ageing index is 76.8% which is equivalent to Young in ageing index level. We conducted a review and meeting session to verify the scored obtained for the system with the system owner and the owner agreed with the result. They claimed that App3 System is still in excellent condition and has been used actively by KET staff. Feedback from system owner is consistent with the finding that we gained from the case study.

Based on the result obtained for App2 shows that App2 System scored is 66.9% which is mapped into ageing index level of Semi-old. Feedback from the software owner claims that App2 System does not have any major problem but the application system has difficulties to upgrade some of the new functions that require by the staff.

### B. The Prototype of SEANA Model

This section presents the prototype which developed in this research in order to validate and automate the ageing process as defined in SEANA model. The prototype is called the Software Anti-Ageing and Rejuvenation Index System (or SeRIS). The development of SeRIS was based on prototyping approach. The system was undergone through alpha and beta testing to validate the correctness and verify based on actual user's requirements. It was also being validated and applied in specific software product in real environment. This confirmation study was carried out collaboratively with industry. Feedbacks from the testing and validation activity were used in refining the SEANA model and SeRIS prototype system.



Fig 4.     SeRIS Main Page.

Fig. 4 to Fig. 6 illustrate the interfaces of SeRIS. The system was implemented to validate the propose ageing measurements and automate the ageing process to ensure the correctness of the computational that involve in the model. The SeRIS prototype system assists users in applying SEANA model in the real environment. SeRIS provides interface to input data of the targeted application software to be evaluated, computes the ageing scores for each factors and produces the ageing level and report.



Fig 5.     SeRIS Main Page - After Assessment.



Fig 6.     Assessment Report Generated by SeRIS System.

## VI. THE ANTI-AGEING GUIDELINES AND ACTIONS

### A. The Anti-Ageing Guideline

The next step in the anti-ageing process is to identify the necessary anti-ageing actions associated with each of the ageing factors. Based on the ageing index level shown in Table III, the anti-ageing actions are proposed to be carried out and applied on the software that has the lowest score during the ageing assessment described in Section IV. The anti-ageing actions are derived from software maintenance activities. Previous researchers revealed that one way to manage software ageing was through systematic maintenance of the software [20]. Software maintenance is a vast activity process of modifying and upgrading the software or part of the software to repair software error or fault, to add new functions or modification, or adaptation to a new environment [22][23].

As discussed in Section II, software maintenance activity can be categorized into four different types: adaptive, corrective, perfective and preventive. This research adopts the maintenance approaches as defined by Matias et al [14] as the baseline of the anti-ageing actions. The proposed anti-ageing actions are derived from literature and case study findings. The actions are recommended to ensure the software stays relevant and fulfil users' need and expectation in the dynamic environment today for longer time of usage. Table VII-IX show the anti-ageing guideline and actions for functional, environment and human factor related to the ageing factors and measurements.

TABLE VII.    ANTI-AGEING FOR FUNCTIONALITY FACTORS

|   | Metrics | Anti-Ageing Action |
|---|---------|---------------------|
| 1 | Time responsiveness | *Perfective*<br> - Monitor the memory usage<br> - Improve the quality aspect of the software.<br> - Check the software structure (optimisation) |
| 2 | Users' requirement and expectation | *Adaptive*<br> - Check quality assessment based on user's perspective and approach.<br> - Enhance or modify software based on user requirements and expectations regularly. |
| 3 | Functionality | *Corrective*<br> - Correct the error and fault of the software accordingly and systematically. |
| 4 | Degradation of Performance | *Perfective*<br> - Improve the functional of software service to increase performance<br>*Corrective*<br> - Check and correct faults tolerant and failure regularly |
| 5 | Software Relevancy | *Perfective*<br> - Improve and enhance the software function to ensure the up-to-date service/functions are available. |
| 6 | Software Faults and Failures | *Corrective*<br> - Correct the fault and error accordingly and systematically.<br> - Improve the change request and process. |

TABLE VIII.    ANTI-AGEING FOR ENVIRONMENTAL FACTORS

|   | Metrics | Anti-Ageing Action |
|---|---------|---------------------|
| 1 | Dynamic to environment change | *Corrective*<br> - Easy maintenance for environment change<br> - Easy maintenance for business change. |
| 2 | Cost for software maintenance and software upgrading | *Training*<br> - Focus on in-house training for staff<br> - Minimize outsourcing<br> - In-house maintenance |
| 3 | Technology demand and compatibility | *Adaptive & Perfective*<br> - Improve and enhance the software according to new/current technology demand and compatibility. |
| 4 | Hardware Changes | *Adaptive*<br> - Improve and enhance the software for meeting new/current hardware demand and compatibility. |
| 5 | Business process change and demand | *Adaptive*<br> - Improve and enhance the software function/services to ensure the demands in business processes are maintained and achieved. |
| 6 | Software technology change | *Adaptive*<br> - Enhance the software for meeting new/current software technology demand and compatibility. |

TABLE IX.    ANTI-AGEING FOR HUMAN FACTORS

|   | Metrics | Anti-Ageing Action |
|---|---------|---------------------|
| 1 | Upgrading and maintenance expert | *Training*<br> - Focus on training for staff in software maintenance and related. |
| 2 | Software quality practice | *Training*<br> - Awareness to the staff for quality assurance<br> - Train staff for software quality practices and implementation. |
| 3 | Software management capability | *Training*<br> - Train for software management practices. |
| 4 | User Interface | *Adaptive*<br> - Improve and enhance the software for user friendliness and usability aspect<br> - Improve and enhance software functionality. |

### B. The Anti-Ageing Action Implementation

The proposed anti-ageing actions are suggested to the systems that have been tested in the case study discussed in previous section. Based on the assessment results, we choose the lowest score among the three systems. In this case the App1 system has been selected with the lowest percentage (47%) and ageing index level of Semi-old. The implementation of the anti-aging actions has been carried out on the three metrics that obtained the lowest score during the assessment. Table X shows the metrics and the proposed anti-ageing actions.

TABLE X.        IMPLEMENTATION OF ANTI-AGEING FOR ACTION

| Metrics | Score (1..4) | Anti-Ageing Action |
|---|---|---|
| Time responsiveness | 1 | Refer to Table VII to apply the anti-ageing action for time responsiveness. Suggested action to be applied is perfective maintenance. |
| User Interface | 1 | Refer to Table IX to apply the anti-ageing action for User Interface. Suggested action to be applied is adaptive maintenance. |
| Cost for software maintenance and software Upgrading | 1 | Refer Table VIII to apply the anti-ageing action for reducing the cost for software maintenance and software upgrading. Suggested action to be applied is by conducting in-house training for manintenance and try to minimize outsourcing. |

## VII. DISCUSSION

This research focuses on identification of software ageing factors and measurements as the fundamental of further related research which in our scope is the anti-ageing model, guideline and actions. This paper starts the discussion with presenting the background and related work of software ageing. The underlying issues and works related to software ageing, rejuvenation, anti-ageing and green and sustainability have been investigated and studied. Later we conducted the empirical study to explore more from real industrial perspective on these related issues and topics.

The findings from empirical study were used as the input to the development of the anti-ageing model. This model is called **S**oftwar**e An**ti-**A**geing or SEANA model. The main objective of the survey was to validate and verify the ageing factors among software practitioners. The empirical study which was conducted in Malaysia revealed three main ageing factors and associated measurements. The ageing factors are categorized as functional, environment and human. Based on the analysis, twenty metrics have been recognized to measure software ageing. The metrics were assigned with numerical scales for further quantifying of the software ageing score and level.

The SEANA model was developed as shown in Fig. 3 and could be used to measure the ageing status/level of any application software operating in certain environments. After the ageing level has been identified based on the assessment, the anti-ageing actions can be generated further aligned with the results of the ageing index. The anti-ageing actions are proposed in order to counteract and minimize the occurrence of ageing phenomenon in the specific targeted software. This is believed will prolong the relevancy of the software operating in their environments. The anti-ageing guideline and actions defined in this model will assist and ease the software owner or the stakeholder to make decision on the solution to be taken in order to deal with software ageing phenomenon if it occurs in their applications.

The SEANA model has been validated and applied in real case study collaborated with local industry in Malaysia. In this agency, three application software were used as case study as described in this paper. Furthermore, the prototype system, SeRIS was developed to validate and automate the process. The case study and SeRIS prototype system prove the effectiveness and practicality of SEANA model.

## VIII. CONCLUSIONS

As a conclusion, even though software ageing is inevitable it can be delayed  by applying anti-ageing techniques that has been presented in this paper. This study focuses on external factors of software ageing and identifies three main and essential ageing factors which are environment, human and functional. The anti-ageing model for application software was developed and tested in real industrial environment, and further recommended an anti-ageing guideline and actions associated with ageing phenomenon in software. For future work, it is suggested that the anti-ageing model proposed in this research to be applied and aligned with the green and sustainability context of software product. Sustainability dimensions which are social, economy and environment can be embedded in the new enhance anti-ageing model. Furthermore, with the prolong usage of software and delaying the aged of the software will reduce the waste and maintain minimum waste disposal of software product and development process.

## ACKNOWLEDGMENT

### REFERENCES

[1] Li, Y. Qi, and L. Cai, "A Hybrid Approach for Predicting Aging-Related Failures of Software Systems," 2018 IEEE Symp. Serv. Syst. Eng., pp. 96–105, 2018.

[2] J. H. Yahaya and A. Deraman, "Towards the Anti-Ageing Model for Application Software," Proc. World Congr. Eng., vol. II, 2012.

[3] L. Parnas, "Software Aging Invited," ICSE '94 Proc. 16th Int. Conf. Softw. Eng., pp. 279–287, 1994.

[4] S. Ahamad, "Study of Software Aging Issues and Prevention Solutions," Int. J. Comput. Sci. Inf. Secur., vol. 14, no. 08, pp. 307–313, 2016.

[5] Melo, J. Araujo, V. Alves, and P. Maciel, "Investigation of software aging effects on the OpenStack cloud computing platform," J. Softw., vol. 12, no. 2, pp. 125–138, 2017.

[6] Cotroneo, A. K. Iannillo, R. Natella, R. Pietrantuono, and S. Russo, "The software aging and rejuvenation repository: Http://openscience.us/repo/software-Aging/," 2015 IEEE Int. Symp. Softw. Reliab. Eng. Work. ISSREW 2015, pp. 108–113, 2016.

[7] R. Mohan and G. Ram Mohana Reddy, "Software aging trend analysis of server virtualized system," Int. Conf. Inf. Netw., pp. 260–263, 2014.

[8] J. H. Yahaya, Z. N. Zainal Abidin, and A. Deraman, "Software Ageing Measurement and Classification Using Goal Question Metric (GQM) Approach," Sci. Inf. Conf. 2013, pp. 160–165, 2013.

[9] Y. Zhao, J. Xiang, S. Xiong, Y. Wu, J. An, S. Wang, and X. Yu, "An Experimental Study on Software Aging in Android Operating System," 2015 2nd Int. Symp. Dependable Comput. Internet Things, pp. 148–150, 2015.

[10] J. Araujo, V. Alves, D. Oliveira, P. Dias, B. Silva, and P. Maciel, "An Investigative Approach to Software Aging in Android Applications," 2013 IEEE Int. Conf. Syst. Man, Cybern., pp. 1229–1234, Oct. 2013.

[11] S. Huo, D. Zhao, X. Liu, J. Xiang, Y. Zhong, and H. Yu, "Using Machine Learning for Software Aging Detection in Android System," 2018 Tenth Int. Conf. Adv. Comput. Intell., pp. 741–746, 2018.

[12] J. Araujo, R. Matos, V. Alves, P. Maciel, F. V. de Souza, R. M. Jr., and K. S. Trivedi, "Software aging in the eucalyptus cloud computing infrastructure," ACM J. Emerg. Technol. Comput. Syst., vol. 10, no. 1, pp. 1–22, 2014.

[13] J. H. Yahaya, A. Deraman, S. R. A. Ibrahim, and Y. Y. Jusoh, "Software Certification Modeling: From Technical to User Centric Approach" Aust. J. Basic Appl. Sci., vol. 7, no. 8, pp. 9–18, 2013.

[14] R. Matias Jr., K. S. Trivedi, and P. R. M. Maciel, "Using Accelerated Life Tests to Estimate Time to Software Aging Failure," 2010 IEEE 21st Int. Symp. Softw. Reliab. Eng., pp. 211–219, Nov. 2010.

[15] Z. H. Abdullah, J. Yahaya & A. Deraman "The Anti-Ageing Model for Assessment of Application Software," Postgrad. Res. Work. , SoftTech Asia 2018, 2018.

[16] F. Machida, J. Xiang, K. Tadano, and Y. Maeno, "Lifetime extension of software execution subject to aging," IEEE Trans. Reliab., vol. 66, no. 1, pp. 123–134, 2017.

[17] D. Cotroneo, R. Natella, R. Pietrantuono, "A Survey of Software Aging and Rejuvenation Studies," ACM J. Emerg. Technol. Comput. Syst. - Spec. Issue Reliab. Device Degrad. Emerg. Technol. Spec. Issue WoSAR 2011, vol. V, no. 212, p. 30, 2010.

[18] D. N. Arnold, "The Patriot Missile Failure," 2000.

[19] D. Cotroneo, R. Natella, R. Pietrantuono, and S. Russo, "Software Aging and Rejuvenation: Where We Are and Where We Are Going," 2011 IEEE Third Int. Work. Softw. Aging Rejuvenation, no. 30, pp. 1–6, Nov. 2011.

[20] J. Zhao, K. S. Trivedi, Y. Wang, and X. Chen, "Evaluation of software performance affected by aging," 2010 IEEE Second Int. Work. Softw. Aging Rejuvenation, vol. 3, pp. 1–6, Nov. 2010.

[21] I. Sommerville, Software Engineering (Tenth Edition). 2016.

[22] Z. H. Abdullah, J.H.Yahaya, Z. Mansor & A.Deraman. "Software Ageing Prevention from Software Maintenance Perspective – A Review," Journal of Telecommunication, Electronic and Computer Engineering, vol. 9, no. 3-4, pp. 93-96, 2017.

[23] J. H. Yahaya, A. Deraman & Z. H. Abdullah. "Evergreen Software Preservation: The Conceptual Framework of Anti-Ageing Model," Information Science and Applications, Lecture Notes in Electrical Engineering, vol. 339, pp 899-906, 2015.

[24] Z. N. Zainal Abidin, J.H. Yahaya, A. Deraman & Z. H. Abdullah. "Rejuvenation Action Model for Application Software," The 6th International Conference on Information and Communication Technology (ICoICT 2018), 3-5 May 2018.

[25] S. Murugesan & P.A. Laplante. "IT for a greener planet," IT Pro January/February, pp. 16–20, 2011.

[26] Murugesan, S. "Harnessing green IT: Principles and practices," IEEE IT Professional, vol. 10, no. 1, pp.5-6, 2008.

[27] M. Dik, J. Drangmeister, E. Kern & S. Naumann. "Green software engineering with agile methods in green and sustainable software (GREENS)," Proceedings of 2013 2nd International Workshop on Green and Sustainable Software , pp 78–85, 2013.

[28] K. Erdelyi. "Special factors of development of green software supporting eco sus-tainability," Proceedings of EEE 11th international symposium on intelligent systems and informat-ics (SISY), pp 337–340, 2013.

# Understanding Customer Voice of Project Portfolio Management Software

Maruthi Rohit Ayyagari[1]

College of Business, University of Dallas
Irving, Texas, USA

Issa Atoum[2]

Department of Software Engineering
The World Islamic Sciences and Education, Jordan

*Abstract*—**Project Portfolio Management (PPM) has gained success in many projects due to its large number of features that covers effective scheduling, risk management, collaboration, and third-party software integrations to mention a few. A broad range of PPM software is available; however, it is essential to select the PPM with minimum usage issues over time. While many companies use surveys and market research to get users feedback, the PPM product software reviews carry the voice of users; the positive and negative sentiments of the PPM software reviews. This paper collected 4,775 reviews of ten PPM software from Capttera.com. Our approach has these phases- text preprocessing, sentiment analysis, summarization, and categorizations. The software reviews are filtered and cleaned, then negative sentiments of user reviews are summarized into a set of factors that identify issues of adopted PPM software. We report the most important issues of PPM software which were related to missing technological features and lack of training. Results using Latent Dirichlet Allocation (LDA) model showed that the top ten common issues are related to software complexity and lack of required features.**

*Keywords—Project Portfolio Management (PPM); software reviews; sentiment analytics; text summarization; LDA*

## I. Introduction

Every organization strive to achieve its strategic goals by executing a set of cornerstone projects[1]. Managing and controlling diverse, interrelated projects as a portfolio is nontrivial. The projects face problems related to change management, scoping [2], complexity [3], timelines, and tracking. Project and portfolio managers must harvest the features of project portfolio management (PPM) software to ensure proper control. Common issues of PPM include prioritization [4], inaccurate reporting, resource utilization, and software development lifecycle. The primary objective of PPM is to execute projects that support strategic organization goals under constraints of scope, time and resources.

Recently, an increase in PPM software is noted. According to Gartner [5], the estimated $2.5 billion project portfolio management market demonstrates stability, as well as an increasing level of change. The driving forces to PPM software tools are traced back to requirements of PPM practitioners and stakeholders, organization configuration management, demand of collaboration between users, and increased the complexity of enterprise projects [3]. If the PPM software tools are implemented according to business needs, they have the potential to improve organizational business benefits aligned with business strategy at the portfolio level.

The organizations that use PPM tools are 44% more likely to complete projects on budget, and 52% more likely to get the anticipated Return on Investment [6]. Implementors and adopters of PPM got a decrease of failure by 59%, spent 37% less per project, reduced the wide variety of redundant projects by 78%, and increased resource constructiveness by 14% [7].

Although PPM tools implementation is recognized in practice, current understanding issues of PPM tools are limited [8], [9]. As organizations strive to turn out to be globally aggressive while increasing shareholder's value, they are always compelled to reduce infrastructure costs to get products into the market cheaper, quicker and with better high-quality. Therefore, the PPM software vendors attempt to continue software evolution based on user needs. While users' feedbacks can be gained by surveys[8] where target PPM software is showing an increase in the market place, globalization causes extremely high competition between PPM vendors. Moreover, enterprise software adopts the roles of different levels of practitioners; therefore, studying all these stakeholders could increase vendors' revenue by planning for the next product release. Therefore, the need for a systematic study on the issues behind the diffusion of PPM software tools in organizations is decisive.

This paper develops a taxonomy of PPM software tools issues of a set of selected of 10 favorite PPM software tools—Microsoft Project, Wrike, Atlassian, Basecamp, Trello, Asana, teamwork projects, Podio, Smartsheet, JIRA. The proposed approach shown in Fig. 1 is based on software reviews of the selected PPM software. We collected 4,775 PPM reviews and yielded 4,397 reviews after removing empty and invalid details. First, we do a preprocessing step by removing stop words and changing words to lowercase. Then, we do sentiment analysis of the cons part of the reviews, as cons are supposed to have negative comments—the output of the sentiment analysis which is the list of negative reviews. Next, we do a summarization based on genism TextRank algorithm [10]. Finally, we allow two coders to read the summarized reviews and analyze potential issues manually. The coders end up with categories and subcategories of PPM software related issues.

A key difference in this paper compared to previous work, is the application of text mining techniques [11] to reduce the overwhelming number of software reviews. Instead of depending of market research that might not give instant output of the PPM issues trend, the proposed approach provides instant results of key success/failure factors of PPM.

Fig. 1.    Context Diagram of the Proposed Model.

There have been many studies in the era of PPM success factors [12]; however, they are limited to identify failure factors of PPM tools [13]. The objective of this paper is to identify the negative voice of the PPM software practitioners so that software vendors can enhance their products accordingly. Therefore, the major contribution here is to both vendors and users. Users have the trend and of PPM software and vendors get mined categorized and instant results of the customers' voice.

The remainder of the paper is structured as follows. Section two summarizes a background of PPM and its issues, and mining software reviews. Section three summarizes related work. Section four illustrates the proposed model while Section 5 evaluates the proposed model. Section 6 provides conclusions, with implications and future research.

## II. BACKGROUND

### A. Project Portfolio Management Software

The complexity of nowadays enterprise projects makes project management success challengeable for many companies. Many studies showed that poor project management means projects finish late, over-budget, and with the lower return of investment. Therefore, businesses often rely on project management software, which controls the project, manages change, and mitigates risk by identifying potential project issues.

A study from Forrester Research Group concluded that implementing a PPM solution produced an expected ROI of 255% [14]. The Aberdeen Group, organizations that the use a PPM tool, are 44 percent more likely to complete projects on time [6]. Consequently, an acceptable PPM software should provide features that draw project success with maximum benefits.

Project management software can help streamline business processes by applying state of the art PPM software development frameworks [15]. The process gets increased in difficulty as the number of projects increase. The problem will be how to maximize the success of various interleaved projects given restrictions on time, budget, and resources. Therefore, PPM software gives a hit of decision making by providing visibility, and oversight with dashboard tools to prioritize and manage current projects.

The PPM software tools provide tools to manage and control projects, provide proper communication, collaboration with resources, progress and management reports, risk management, integration with other enterprise application such HR and CRM, and predicting project future.

### B. Issues of PPM

Enterprise PPM software requires extensive configuration to support business goals. The PPM software tools especially those installed and maintained by the organization are difficult to set up due to the complexity of software that needs a large scale set of experience of hardware and software. The next issue that comes in is the learning curve that can take time for practitioners to understand and find needed information [16]. Generally, a project or portfolio manager must take training on specific software to control and monitor gain the maximum from the PPM tool.

As projects are temporary endeavor undertaken to create a unique product or service [17], they are bound to the environment and change management complexities; therefore, customizing the PPM tool to a specific environment, project, size, and user needs and constraints are essential. Moreover, the PPM tool should be easy to use and user-friendly, where practitioners can easily do the job they are seeking.

PPM tools that run on the cloud suffer from poor efficiency and low availability factors that could affect running critical projects [18]. Moreover, without high-level dashboard decision makers could not provide a creative solution [19]. Therefore, most PPM tools provide reporting, dashboards, and integration with enterprise applications.

### C. Mining Software Reviews

One value-powerful technique to unveil the underlying elements influencing PPM adoption is to research the online software evaluations, furnished with the aid of practitioners on their reviews and practices of PPM systems. Software reviews are an essential source of information for evaluating as it captures the essence of the user's voice. Therefore, the emergent of several software reviews sites has gained attention to may researchers who are most often referred to as option mining or sentiment analysis.

The sentiment analysis is the technique of computationally figuring out and categorizing reviews expressed in a bit of text, mainly to determine whether the author's mindset towards a specific topic, product, and many others are fantastic, sad, or impartial. The valuable information contained in software reviews is critical for users and decision makers.

However, the immense volume of online reviews makes it difficult for users to comprehend text. The Worldwide data will grow 61% to 175 zettabytes by 2025 according to IDC [20]. Most often software reviews are accompanied with software rating which is a number that shows the user satisfaction of a product. However, the details of user satisfaction factors are latent in the text. Therefore, text mining could reveal hidden success factors or joint issues.

Text mining provides a set of tasks that can be used to mine interesting patterns of data. Sentiment analysis of PPM software reviews enables users to identify the pros and cons of products. The massive amount of text could be reduced using text summarization techniques. Textual content summarization refers to the technique of shortening lengthy pieces of textual content. The goal is to create a coherent and fluent precis having only the main points mentioned within the text. In our

paper, summarization is based on ranks of text sentences using a variation of the TextRank algorithm [10].

## III. RELATED WORK

There have been many works that study project portfolio management systems. However, most of them are propriety without details of how results are systematically generated [21]. The problem is that ranking PPM software is done by marketing research leaving the details of why there is a region for improvement on PPM software.

Gartner, a leading research and advisory company provides reports of PPM software tools over the years. According to Gartner, the tools are categorized into four groups known as the magic quadrant group- leaders, challengers, visionaries, niche players. However, the criteria for the Magic Quadrant is leaned towards large vendors than towards buyers [8]. Similarity, Softwareadvice.com, capttera.com provides tools to compare PPM software tools.

Many researches identify reasons for project failures. Failures are due to a list of issues related to organization, people, technical or the process [13]. Failures are often related to lack of executive support, project management skills, communication and risk management. However, these issues are not related to PPM tools. Therefore, without evidence, we cannot adopt or leave a PPM software tool.

The work of [9] proposes a mathematical model prioritize projects to gain maximum resource utilization. Several other works provide best practice techniques to gain project success [22][23]. A review by [24] showed an increased focus on human actors, project management alternative acceptance, adapt to change. The work uses a fuzzy approach to order key success factors of projects [4].

## IV. PROPOSED APPROACH

The PPM Software reviews provide essential feedback about the products usage by practitioners. However, the number of available reviews makes the decision uncomprehensive as a user must read long reviews which are the time-consuming process. Although reviews provide a rating of the product (usually between 1 and 5), the rating does not provide information about why the user was happy or satisfied. Therefore, reviews must be cleaned and then analyzed for possible exciting patterns.

Fig. 2 depicts the context diagram of the proposed approach. The data collection process comprises a set of steps to accumulate software reviews from the website. We choose to scrape PPM reviews from Capterra.com, as the website has a long list of high-quality reviews. The selected PPM software tools are the most rated and used PPM products—Microsoft Project, Wrike, Atlassian, Basecamp, Trello, Asana, Team work Projects, Podio, Smartsheet, JIRA. Our target was to scrape the most recent 500 reviews from each PPM software. We collected 4,775 reviews as Podio had 225 reviews during project execution (step.1). as we were concentrating on the negative sentiment of PPM reviews- which is the cons element in a review- all reviews that have empty cons or invalid details were removed (step 2). Therefore, the total numbers of valid reviews were 4,732 reviews.



Fig. 2. Proposed Approach.

In the text preprocessing step (step 3), the applied transformations on cons element of the reviews are – lowercase, remove accent, parse HTML, remove URLs. Then the resulting text was tokenized using word tokenize pattern next; the text was filtered by removing stop words and all special characters. Consequently, we apply the sentiment analysis algorithm of [25] to identify the negative reviews in the cons element of the reviews. As a result, 63 reviews were omitted from the reviews resulted in 4,669 reviews (step 4).

To reduce the number of reviews, we applied a summarization implementation of Gensim [26] that is based on TextRank Algorithm [10], one of the summarization algorithms that output useful results (step 5). The summarization reduced reviews by 95% (234 reviews).

The last step (step 6) was completed by two project management experts to ensure consistency and validity. The resultant reviews (cons element) along with other details were passed to the experts, and they were asked to look at the cons part and identify the issue of the reviews. They did the job alone in the first phase, and they could discuss and find a solution when disagreements were found. Before they commenced, we discussed with them and described to them that the objective of their job; the objective to see the failure factors that make the PPM software practitioners unhappy with the software.

## V. EVALUATION AND DISCUSSION

### A. Validating Proposed Approach with Experts

The main reported issues of the PPM software where categorized by experts as shown in Table I, ordered based on how issues are mostly covered by reviews (frequency). The most common issue that practitioners face is the complexity of the PPM tools. To support large scale projects vendors had added many features to the software that added extra complexity to users; they had difficulties in finding information, managing cards, set permission, and manage complex projects. Therefore, the finding draws attention to vendors that they should find a way to configure the tool as per project size by probably hiding advanced features unless requested. A few junior practitioners reported that the used language in the PPM interface was not easy to understand; mainly the help material was comprehensive but not will be categorized in many cases. An overall view of the word is shown in Fig. 3. According to the figure, users are worried about time, task, use, features as shown in the middle of the figure.

Furthermore, we discovered that the PPM suffers from issues related to ease of use. Because studied PPM software is large-scale software, the users reported that the software was

not easy to use which is an indirect result of software complexity. They reported that they were not able to navigate quickly and some software was lacking needed project templates for perceived ease of use. This finding direct vendors to customize the navigation and reduce options per needs, perhaps by project size or based on user experience. Given a complicated PPM software that is not always user-friendly, the learning curve was increased as reported by users.

There were issues related to the unavailability of needed features. Many users reported that there were not able to do require reporting unless they integrate with other third-party tools. They also reported that the most common missing feature was handling risk management as part of the PPM software. Therefore, vendors could take this missing feature in the upcoming versions of the software.



Fig. 3.   Word Cloud of PPM Software Issues.

Furthermore, the PPM software for critical applications needs a cloud system that is available and efficient in situations where data is exchanged in complex projects. However, non-cloud-based PPM tools need extended configuration and initial setup before startup.

Also, one common suggestion by users is that the PPM tool should support the automation of some tasks-workflows and triggering task completions automatically.

### B. Validating Proposed Approach with Topic Modeling

We are aware of other related works that discuss the issues of PPM software adoption[12], however, most of studied works have different comparison parameters [13] that does not map with our approach. Moreover, the intent of this research is the *negative* voice of the customer not the technique used to get this voice; therefore, we opt out comparisons with market researchers or comparing datamining techniques performance. Therefore, another additional technique was used to validate the proposed model.

We run another experiment to show that the summarization algorithm does not omit the necessary PPM issues. We run topic modeling, the latent Dirichlet allocation (LDA) of [27], [28], an unsupervised method that can identify top topics discussed by users. The LDA results of the set of reviews are shown in Table II. Using the terms of the topic models, we assign the topic relative to our findings in column (3). However, we were not able to map the last two topics, due to the accuracy and completeness of the summarization carried out on cons element of the review only. Consequently, our previous findings are 80% accurate compared to topic modeling; therefore, our findings are complete and more accurate than automated topic modeling results.

TABLE I.    TAXONOMY OF PPM SOFTWARE ISSUES (NOTE THAT REVIEWS HAVE SPELLING ERRORS, WE KEEP THEM AS IS)

| Issue Category | Issue Subcategory | Sample PPM Review |
|---|---|---|
| **Complexity** | resource card management | *Given the robust features and functions that Jira supports, there's a lot of depth and breadth to the software, so it can be somewhat complex and confusing to newer users, especially those who haven't worked much with project management tools previously. For project managers who are leading multiple teams simultaneously, things can get a bit overwhelming because the email notifications and project alerts can quickly start to become excessive, and there's no way to consolidate notifications across projects/teams.* |
| | too many data update | |
| | extra plugins and options | |
| | task management: task for managers, filter tasks, dependency, and change | |
| | manage complex projects | |
| | setting permissions | |
| | match project size | |
| | tracking complexity | |
| | plan execution | |
| | control resources | |
| | setup addon | |
| | Comprehensive help material | |
| | difficult interface language | |
| **Ease of use** | templates | *This isn't really a con, but there is a bit of a learning curve in terms of navigating and understanding how to utilize all of the features. Navigation is not completely innate so that took a minute to figure out. It's a big tool though so that shouldn't be a surprise. They have great resources and support to help new users get adjusted, which is awesome. It only took a day or so of messing around with the different features to feel fully comfortable in the capacity that I needed.* |
| | navigation | |
| | User experience (UX) | |
| | find information | |
| **Learnability** | - | *There's a very steep learning curve. I'm writing from the perspective of a product manager so it could be that developers* |

| | | |
|---|---|---|
| | | *have a much easier time grasping the concepts here. For me, it took awhile to get an understanding of the relationships for setting up the sprint logic and mapping the progression of a story through the board. Once I started to immerse myself in it, I started to grasp it more and more. That said, I could see a more in-depth sample project being incredibly helpful for tutorializing within the tools themselves.* |
| **Lack of features** | built in reports | *This software does not have some of the features like time sheets and a risk management section to help tract risks and keep everyone up to speed with risks.*<br><br>*We trialed Microsoft PPM in one of our subsidiaries, which is a multi-platform/mixed operating system entity. Microsoft PPM caused several issues for our workflows.* |
| | risk management | |
| | manage old tickets | |
| | export file types | |
| | comparative's features | |
| | type of addon | |
| **Customization** | workflow customization | *Even if TRELLO have great options and features, the customization of the screen is a little tricky, I wish to reorder some social alerts on my own hierarchy and the most time TRELLO over my desires and rearrange everything and it's really frustrating to rearrange everything each time that I login to my account.* |
| | admin involvement | |
| | lack of boards | |
| | edit board | |
| **Configuration** | software setup required | *When we scaled, the License was not as flexible in higher user tiers as it was in the lower ones. Also once upon a time an update went wrong and we had to invest a bucketload of time and brain power to get everything back up to speed. We did run jira on premise and hat some custom plugins so there might be some of the tripwires: In case you use plugins, they must be available in the version you wish to upgrade to and this might require some research prior to updates. Also, you can make a science out of the workflow and notification configuration. In case it is not documented well when you do it it might get you into trouble.* |
| | configuring PPM software | |
| | integration | |
| **Collaboration** | log meetings | *I have been using this software for the last year and I still can't find a way to assign a task to more than one person. This seems like a great software if each project is solely owned by one person. The collaboration with this program is not the greatest, but for individual projects it is great.* |
| | overwhelming emails | |
| | allow user-user communication | |
| | log surveys | |
| **Service Quality** | efficiency | *Because it is a web application, if you present problems with the internet connection you will not be able to access the information of your project, it is always a negative point since the data should be available.* |
| | availability | |
| **Automation** | drag and drop | *The product requires a high investment of time to not only understand yourself but also train your team, when I became a mid-level user I found it tedious to set up projects in the customizable way we required.* |
| | automate workflows | |
| | trigger cards | |
| | reminders | |
| **General** | cross-technology compatibility | *There have been times when I've thought that the lingo, they use on the app might be confusing for new users, or that other wording would have better explained what a certain tool was.* |
| | cost | |

TABLE II. TOP 10 TOPICS OF ISSUES OF PPM SOFTWARE TOOLS

| # | Terms | Topic |
|---|---|---|
| 1 | *tool, found, option, small, integration, interface, isn, cost, due, improved* | complexity |
| 2 | *features, bit, app, people, difficult, confusing, love, platform, understand, slow* | lack of features |
| 3 | *software, users, learning, management, issues, user, good, nice, support, team* | learnability |
| 4 | *setup, tickets, button, ticket, navigate, connection, awesome, screen, bugs, phone* | configuration |
| 5 | *make, platform, email, visually, back, navigate, loved, moving, function, order* | ease of use |
| 6 | *management, create, add, information, made, nice, update, items, fast, drop* | service quality |
| 7 | *support, reporting, tools, easier, integrations, workflows, customer, number, files, unlike* | configuration |
| 8 | *tool, sheet, support, collaboration, dislike, samepage, share, high, move, task* | collaboration |
| 9 | *time, complex, drag, ability, manage, found, smartsheet, file, program, personally* | mixed |
| 10 | *work, create, view, clunky, excel, bit tricky, apps, flow, ons, products* | mixed |

## VI. IMPLICATIONS

The implication of this work affects the PPM software vendors who should find the current issues that are degrading their software. Therefore, they could plan their maintenance activities better to fix bugs available bugs. Moreover, the results show that vendors should take care more of user needs of usability and User Experience (UX) issues. Furthermore, they should change their training and software design strategies to reduce software complexity. On the other hand, users could get a glance of what they could choose based on previous analyzed knowledge and experience of previous practitioners.

## VII. LIMITATIONS

Although the proposed approach has reported interesting issues of PPM software readers has to consider a few limitations. The list of reviews was collected from one site; therefore, researchers should consider this before generalizing results. Also, we assume that the latest 500 reviews have the issues of the current version of the PPM; however, it might always be right all the time especially in cases where the practitioner is admired by the software. Moreover, all the ten most popular PPM software is included in as coherent contiguous documents, which limits the identity of the PPM that raises the issue.

## VIII. CONCLUSION

This paper proposes an approach to identify the most common issues of PPM software. The approach is based on sentiment analysis to choose negative reviews and on summarization techniques to reduce the number of reviews. The study was carried out over ten popular PPM software— Microsoft Project, Wrike, Atlassian, Basecamp, Trello, Asana, teamwork projects, Podio, Smartsheet, and JIRA. We reported the issues fine-grained with the issue category. Results reported showed that the issues are related to complexity, ease of use, learnability, lack of features, customization, collaboration, service quality, and configuration. Our findings reported that users were not worried about compatibility and cost of PPM software. In the future, we will show issues per PPM software.

### REFERENCES

[1] O. Zwikael, Y.-Y. Chih, and J. R. Meredith, "Project benefit management: Setting effective target benefits," Int. J. Proj. Manag., vol. 36, no. 4, pp. 650–658, 2018.

[2] S. Shafiee, K. Kristjansdottir, L. Hvam, and C. Forza, "How to scope configuration projects and manage the knowledge they require," J. Knowl. Manag., vol. 22, no. 5, pp. 982–1014, 2018.

[3] W. Vogel and R. Lasch, "Single approaches for complexity management in product development: An empirical research," in Supply Management Research, Springer, 2019, pp. 151–215.

[4] T. Yaghoobi, "Prioritizing key success factors of software projects using fuzzy AHP," J. Softw. Evol. Process, vol. 30, no. 1, p. e1891, 2018.

[5] D. B. Stang and M. Light, "Gartner, Magic Quadrant for Project Portfolio Management, Worldwide," 2018.

[6] The Aberdeen Group, "Managing the Project Portfolio to Improve Profits," 2014.

[7] IDC, "Measuring the ROI of PPM (Project and Portfolio Management)," 2018.

[8] T. Byrne, "Looking beyond the magic quadrant to find the nitty-gritty," 2009. [Online]. Available: https://www.realstorygroup.com/Blog/1660-Looking-beyond-the-magic-quadrant-to-find-the-nittygritty.

[9] B. Titarenko, A. Hasnaoui, R. Titarenko, and L. Buzuk, "Robust data analysis in innovation project portfolio management," in MATEC Web of Conferences, 2018, vol. 170, p. 1017.

[10] F. Barrios, F. López, L. Argerich, and R. Wachenchauzer, "Variations of the similarity function of textrank for automated summarization," arXiv Prepr. arXiv1602.03606, 2016.

[11] G. Di Fabbrizio, A. Aker, and R. Gaizauskas, "Summarizing Online Reviews Using Aspect Rating Distributions and Language Modeling," IEEE Intell. Syst., vol. 28, no. 3, pp. 28–37, May 2013.

[12] T. J. Kloppenborg, D. Tesch, and R. R. Chinta, "Demographic determinants of project success behaviors," Pulse, vol. 18, 2019.

[13] H. Taherdoost and A. Keshavarzsaleh, "A theoretical review on IT project success/failure factors and evaluating the associated risks," 2018.

[14] C. Symons, "The ROI of project portfolio management tools," Forrester, May, vol. 8, 2009.

[15] D. Lock and R. Wagner, The Handbook of Project Portfolio Management. Routledge, 2018.

[16] R. R. Nelson, "IT project management: Infamous failures, classic mistakes, and best practices.," MIS Q. Exec., vol. 6, no. 2, 2007.

[17] K. Heldman, PMP: project management professional exam study guide. John Wiley & Sons, 2018.

[18] J. A. Araúzo, J. Pajares, and A. Lopez-Paredes, "Simulating the dynamic scheduling of project portfolios," Simul. Model. Pract. Theory, vol. 18, no. 10, pp. 1428–1441, 2010.

[19] S. Rajegopal, P. McGuin, and J. Waller, Project portfolio management: Leading the corporate vision. Springer, 2007.

[20] P. Andy, "IDC: Expect 175 zettabytes of data worldwide by 2025," 2018. [Online]. Available: https://www.networkworld.com/article/3325397/idc-expect-175-zettabytes-of-data-worldwide-by-2025.html. [Accessed: 01-Apr-2019].

[21] Gartner® Inc, "FrontRunners Methodology," 2018. [Online]. Available: https://www.saimgs.com/imglib/other_pages/FrontRunners/MethodologyOverview.pdf. [Accessed: 28-Feb-2019].

[22] H. Kerzner, Project management best practices: Achieving global excellence. John Wiley & Sons, 2018.

[23] C. N. Enoch, Project portfolio management: a model for improved decision making. Business Expert Press, 2019.

[24] L. K. Hansen and P. Svejvig, "Towards rethinking Project portfolio management," 2018.

[25] B. Liu, M. Hu, and J. Cheng, "Opinion observer: analyzing and comparing opinions on the Web," in Proceedings of the 14th international conference on World Wide Web, 2005, pp. 342–351.

[26] R. Řehůřek and P. Sojka, "Software Framework for Topic Modelling with Large Corpora," in Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks, 2010, pp. 45–50.

[27] Q. Mei, X. Ling, M. Wondra, H. Su, and C. Zhai, "Topic sentiment mixture: modeling facets and opinions in weblogs," in Proceedings of the 16th international conference on World Wide Web, 2007, pp. 171–180.

[28] D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent dirichlet allocation," J. Mach. Learn. Res., vol. 3, pp. 993–1022, Mar. 2003.

# Analysis of the Emotions' Brainwaves

Witman Alvarado-Díaz[1], Brian Meneses-Claudio[2], Avid Roman-Gonzalez[3]

Image Processing Research Laboratory (INTI-Lab), Universidad de Ciencias y Humanidades

Lima, Peru

*Abstract*—Currently in Peru, patients with degenerative diseases, such as Amyotrophic Lateral Sclerosis (ALS) have lost of communication ability. Many researchers' papers that establish basic communication system for these patients. It is also essential to know their feelings or state of mind through their emotions, in this study, we present an analysis of electroencephalographic signals (EEG) applied to emotions such as fear, tenderness, happiness and surprised; it was used linear discriminant analysis (LDA) to get the identification and classification of the 4 emotions with a success rate of 63.36% on average.

*Keywords*—*Electroencephalogram; EEG; emotions; amyotrophic lateral sclerosis; degenerative diseases; classification learner*

## I. INTRODUCTION

Approximately 210 cases of ALS of which more than 50 are located in Lima, more cases in people aged 30 to 59 years. It should be mentioned that the total cases of degenerative diseases are approximately 5% of the national population, of which a significant part has already lost their motor and communication capacity. In this sense, it has been working in many applications and EEG studies such as: [1], [2], [3], [4] in which systems are created in order to support patients suffering from degenerative diseases, however it is also essential to study the emotions in this type of patients, in order to achieve improvements in communication as we can understand their emotions and feelings.

In the present work, the classification of emotions is carried out, [5] which are intense affective events that arise in the face of the perception of transcendental situations for a person. Emotions are states in which a combination of sensations and feelings determines the behavior of people; we will focus on studying fear, tenderness, happiness, and surprise. Since, they are essential to be able to know the primary emotional state of the patient.

Fear is the emotion that communicates that the person is in danger; with this, we can know the potential danger in which the patient or the people that surround them can be.

Happiness is produced as a result of something favorable, that likes or benefits; we could identify if the patient is comfortable with the care provided.

The surprise guides us in unexpected situations; thanks to this emotion, we can identify that the patient could be in an uncomfortable situation.

The tenderness reflects the affection, love for a person or animal or thing, with this emotion, the patient will be able to express his affection towards his loved ones as well as the people that surround him.

In [6] they analyzed the emotions sadness, anger and surprise, mention that there are investigations that show that there is a link between the frontal lobe of the brain and emotions, finding that positive emotions are shown on the left and negative ones on the right.

In [7] they decompose the signals generated by emotions into ten sub-bands and use neural networks with a Feedforward Backpropagation algorithm to achieve the classification of the signals.

In [8] propose a method of spectral asymmetric index (SASI), for data divided into 5 bands, in positions fp1 and fp2, in order to detect emotions in the beta (13 - 25Hz) and theta (4 - 8 Hz) frequencies also use a gradient boosting decision tree (GBDT) algorithm for the classification of positive and negative signals.

In [9] they show that the most active area of the brain, concerning emotions is the right side, they achieved it by observing the power of the alpha band (8 - 12 Hz), for the condition of disgust in the temporal regions, front and back.

In [10] they worked with emotions happiness, sadness, angry and relaxed. Dividing the signals into the gamma, beta, alpha, and theta bands, demonstrating the beta and gamma bands are reliable bands for recognizing emotion with the EEG signals; further mentioning that the alpha band can also be considered for the recognition of emotions while theta band can be ignored.

In [11] they study happy, calm, sadness and frightening emotions, showing that the alpha and beta bands contain useful characteristics, also mentions that when a subject sees emotional stimuli, the power decreases in the alpha band but increases in the beta band, In other words, the distribution of the power spectrum in the brainwave patterns changes.

Section II presents the methodology that has been followed for the research work. In Section III, the reader will find the preliminary results obtained. Section IV shows the respective discussions and conclusions for this research work.

## II. METHODOLOGY

For the present work, the method to follow is outlined in the block diagram shown in Fig. 1.

### A. Acquisition of Data

For the data acquisition stage, the OpenBCI system has been used, [12] which consists of the main board that is the Cyton Board, which is an Open Source tool, used for the

acquisition of EMG, EEG and ECG signals; together with its complement (Daisy Board), they have 16 channels, from which data can be obtained at a sampling frequency of 125Hz. The OpenBCI system communicates wirelessly via Bluetooth to the computer through OpenBCI Dongle (Programmable USB).

For the reception of data, the pc uses Python, as a bridge for the transmission of data via Lab Streaming Layer, then the data is received in Matlab, with the use of a graphical interface which is in charge of collecting the data and save them in text files (Fig. 2).

The EEG signals are obtained through the 16 channels of the OpenBCI system, which corresponds to the positions Fp1, Fp2, C3, C4, T5, T6, O1, O2, F7, F8, F3, F4, T3, T4, P3, P4, of the international system 10/20 (Fig. 3).



Fig. 1.   Work Scheme.



Fig. 2.   System for Data Acquisition.



Fig. 3.   Positioning of Electrodes According to the 10/20 System.

### B. Processing and Classification

As part of the data processing, we write a code in Matlab which links the data, in Fig. 4, the reader can see the flow chart which we describe below.

First, we request the entry of the name of the data, and create a directory then the following steps are followed:

- We read the first file and delete column 17th since it only contains data of the time elapsed during the acquisition of data.

- We create a vector with the class to which the data belongs.

- Next, we enter a loop that is responsible for reading the files one by one and joins them to the first file mentioned above; also, it does the same for the classes to which the data belongs.

- Finally, the whole matrix is transformed into a table of 192 x 7501 in which the last column corresponds to the classes.



Fig. 4.   Flow Chart for Grouping the Data.

## III. RESULTS

Samples were taken to five people, to which they are shown 3 videos for each emotion of 1 minute each, for the feelings, surprise, happiness, tenderness and fear as the reader can see in Fig. 5, the 12 videos are shown in a way randomly, and also the collected data is saved in different files.



Fig. 5. Captures of the Videos Shown During Data Acquisition.

As mentioned above, the Matlab Classification Learner application was used to classify the data into four different classes corresponding to each emotion studied, using linear discriminant analysis (LDA) with a full covariance structure and five cross-validation cuts.

The linear discriminant analysis (LDA) also known as Fisher's linear discriminant, named for its inventor Sir RA Fisher; this analysis has been applied for decades, according to [13] and [14], the LDA has a statistical approach to reduce the dimensionality of the data, calculating an optimal projection, in order to minimize the distance within the classes and maximize the distance between classes. The classical LDA requires that the total dispersion matrix will not be singular however according to [14] and [15] in many problems such as information retrieval, facial recognition, machine learning, bioinformatics, data analysis, etc.; the total dispersion matrix mention previously can be singular, since the dimensional data is usually very high and the dimensions generally exceed the number of data points generating a sampling problem or singularity.

In [15], [16], [17], [18] many mathematical formulas and LDA algorithms are described, it is mention for example that the LDA can be used to classify an X observation of a q-dimensional vector which it obtains by observing one of several classes that may be unknown; one of the several ways to describe the LDA is through the use of probability models; assuming that the class $i$th has a density$f_i(X)$ and a probability $\pi_i$, taking into account Bayes formula, which tells.

$$P(class = i|X) = \frac{f_i(X)\pi_i}{\sum_j f_i(X)\pi_j} \qquad (1)$$

To demonstrate the Bayes rule or classifier, which says that the largest conditional probability classification will obtain the least expected number of errors, it assumes that class $i$ has a Gaussian distribution with mean $\mu i$ and covariance $\sum$ , it is

demonstrating that classifying the maximum conditional probability is equivalent to classifying

$$\arg\ max_i\ (L_i) \qquad (2)$$

Where $L_i$ is the discriminate function.

$$L_i = X^T \sum{}^{-1} \mu i - \mu_i^T \sum{}^{-1} \mu_i/2\ + \log(\pi_i) \qquad (3)$$

When using the maximum similarity estimates for $\mu i$ take into account that $L_i$ is a linear function of X since it goes at the LDA procedure.

The LDA is a fast and accurate algorithm, which assumes that the different classes generate data based on Gaussian distributions; in the MATLAB tool, the LDA uses the "fitcdiscr" (Fit discriminant analysis classifier) function, which returns a discriminant analysis model or a classifier, based on the input variables contained in a table and the output, responses or labels.

When the application started and later adding the data already prepared, we can see its graph Fig. 6; then we choose the analysis that we want to do and start the application.



Fig. 6. Data Plotted by the Classification Learner Application.



Fig. 7. Confusion Matrix Generated by the Classification Learner Application.

As a result, the application performs the classification of the data and provides us with a confusion matrix Fig. 7 in which we can see the performance of the chosen algorithm.

Having in mind that the data entered into the application have only been linked, so it is "unprocessed" with which efficiency of 63.36% is obtained on average in Table I, we can see respective efficiencies for each test subject.

TABLE I.    TABLE OF ACCURACY

| Name Data | Accuracy (%) |
|-----------|--------------|
| Subject-1 | 31.8 |
| Subject-2 | 67.2 |
| Subject-3 | 94.8 |
| Subject-4 | 43.8 |
| Subject-5 | 79.2 |

## IV. DISCUSSION AND CONCLUSIONS

For the classification of emotions, different algorithms can be used as in the works [19], [8]; we obtain better results with linear discriminant analysis, and the use of the Classification Learner application.

In work [20] they obtained 55.3% applying linear discriminant analysis studying the emotions anger, fear, and surprise, evoking the emotions through videos.

Also in the consultation [6], they obtained an efficiency between 48.78% and 57.04% using the linear discriminant analysis and support vector machine, working only in the identification of positive and negative emotions, in the right and left part of the frontal lobe.

Finally, in [10] an efficiency of 91.01% is obtained, in the study of happiness, sadness, anger and relaxed, it was achieved by decomposing their data in the signals corresponding to the frequency bands of the brain.

As we can see in other works, it has been possible to classify the brain signals produced by emotions; in our case, a preliminary efficiency of 63.36% was obtained on average.

For future work, we will take into account the accuracy obtained to process the data and perform the analysis on more people, which will allow us to improve the methodology and achieve better results.

### REFERENCES

[1] W. Alvarado-Díaz, B. Meneses-Claudio, and A. Roman-Gonzalez, "Implementation of a brain-machine interface for controlling a wheelchair," in 2017 CHILEAN Conference on Electrical, Electronics Engineering, Information and Communication Technologies, CHILECON 2017 - Proceedings, 2017.

[2] B. Meneses-Claudio, W. Alvarado-Diaz, A. Roman-Gonzalez, and E. Z. Villaorduña, "Implementation of a Wireless system for the processing and comparison of cerebral waves of patients with Amyotrophic Lateral Sclerosis through matlab identifying their basic needs," in 2017 CHILEAN Conference on Electrical, Electronics Engineering, Information and Communication Technologies, CHILECON 2017 - Proceedings, 2017.

[3] A. Roman-Gonzalez, N. I. Vargas-Cuentas, M. Hoyos, J. Diaz, and M. Zimic, "Brain Computer Interface to Answer Yes-No Questions," 2017.

[4] B. Meneses-Claudio and A. Roman-Gonzalez, "Study of the Brain Waves for the identification of the Basic Needs of Patients with Amyotrophic Lateral Sclerosis," Congr. Argentino Ciencias la Inform. y Desarro. Investig. CACIDI 2018, pp. 1–6, 2018.

[5] M. Pallarés, Emociones y sentimientos : dónde se forman y cómo se transforman. Marge Books, 2010.

[6] M. Cao, G. Fang, and F. Ren, "EEG-BASED Emotion Recognition In Chinese Emotional WordS," pp. 3–7, 2011.

[7] S. G. Mangalagowri and P. C. P. Raj, "EEG feature extraction and classification using feed forward backpropogation algorithm for emotion detection," 2016 Int. Conf. Electr. Electron. Commun. Comput. Optim. Tech. ICEECCOT 2016, pp. 183–187, 2017.

[8] S. Wu, X. Xu, L. Shu, and B. Hu, "Estimation of valence of emotion using two frontal EEG channels," Proc. - 2017 IEEE Int. Conf. Bioinforma. Biomed. BIBM 2017, vol. 2017–Janua, pp. 1127–1130, 2017.

[9] P. C. Petrantonakis and L. J. Hadjileontiadis, "A novel emotion elicitation index using frontal brain asymmetry for enhanced EEG-based emotion recognition," IEEE Trans. Inf. Technol. Biomed., vol. 15, no. 5, pp. 737–746, 2011.

[10] E. S. Pane, M. A. Hendrawan, A. D. Wibawa, and M. H. Purnomo, "Identifying Rules for Electroencephalograph (EEG) Emotion Recognition and Classification," Proc. 2017 5th Int. Conf. Instrumentation, Commun. Inf. Technol. Biomed. Eng. ICICI-BME 2017, no. November, pp. 167–172, 2018.

[11] R. M. Mehmood, R. Du, and H. J. Lee, "Optimal feature selection and deep learning ensembles method for emotion recognition from human brain EEG sensors," IEEE Access, vol. 5, no. c, pp. 14797–14806, 2017.

[12] OpenBCI, "Herramientas Open Source Biosensores (EEG, EMG, EKG, y más) - OpenBCI," 2016. [Online]. Available: http://openbci.com/. [Accessed: 11-Apr-2017].

[13] Shuiwang Ji and Jieping Ye, "Generalized Linear Discriminant Analysis: A Unified Framework and Efficient Model Selection," IEEE Trans. Neural Networks, vol. 19, no. 10, pp. 1768–1782, 2008.

[14] W. Y. Yang, S. X. Liu, T. S. Jin, and X. M. Xu, "An optimization criterion for generalized marginal Fisher analysis on undersampled problems," Int. J. Autom. Comput., vol. 8, no. 2, pp. 193–200, 2011.

[15] D. Chu, L. Z. Liao, M. K. P. Ng, and X. Wang, "Incremental Linear Discriminant Analysis: A Fast Algorithm and Comparisons," IEEE Trans. Neural Networks Learn. Syst., vol. 26, no. 11, pp. 2716–2735, 2015.

[16] N. A. A. Shashoa, N. A. Salem, I. N. Jleta, and O. Abusaeeda, "Classification depend on linear discriminant analysis using desired outputs," 2016 17th Int. Conf. Sci. Tech. Autom. Control Comput. Eng. STA 2016 - Proc., no. 1, pp. 328–332, 2017.

[17] P. Markopoulos, "LINEAR DISCRIMINANT ANALYSIS WITH FEW TRAINING DATA Panos P . Markopoulos Department of Electrical and Mircoelectronic Engineering Rochester Institute of Technology," IEEE Int. Conf. Acoust. Speech, Signal Process. 2017, pp. 4626–4630, 2017.

[18] G. M. James and T. J. Hastie, "Functional linear discriminant analysis for irregularly sampled curves," J. R. Stat. Soc. Ser. B Stat. Methodol., vol. 63, no. 3, pp. 533–550, 2001.

[19] P. S. Ghare and A. N. Paithane, "Human emotion recognition using non linear and non stationary EEG signal," Int. Conf. Autom. Control Dyn. Optim. Tech. ICACDOT 2016, pp. 1013–1016, 2017.

[20] M. S. Park, H. S. Oh, H. Jeong, and J. H. Sohn, "EEG-based emotion recogntion during emotionally evocative films," 2013 Int. Winter Work. Brain-Computer Interface, BCI 2013, pp. 56–57, 2013.

# Multi-Agent Architecture of Intelligent and Distributed Platform of Governance, Risk and Compliance of Information Systems

## Design and Implementation of a Communication System Dedicated to this Platform

Soukaina Elhasnaoui[1], Saadia Drissi[2], Hajar Iguer[3], Hicham Medromi[4]

(EAS- LRI) Systems Architecture Team, ENSEM, Hassan II University
BP.8118, Oasis Casablanca, (LPRI) EMSI
Casablanca, Morocco

*Abstract*—**Governance, risk management and compliance of information technologies (IT GRC) is the responsibility of the company's executives. The IT GRC responds to the important concerns of information systems managers, to ensure the necessary changes in the Information System (IS) over time, and enable it to meet the needs of risk mitigation, regulatory compliance, value creation and strategic alignment. Like a large number of organizations' activities, the IT GRC has to find a solution that is equipped through IS applications. Although these tools do exist, they are never developed by considering the IT GRC processes as a whole. We respond to this lack of consideration by proposing an intelligent and distributed platform of risk, governance and compliance of information systems that deploys a variety of IT GRC best practices and frameworks and makes an intelligent choice under constraints and parameters of the best framework to evaluate the objectives and processes in question. EAS-COM (communication system dedicated to the IT GRC platform) is our second proposal in this work: it ensures end-to-end communication between the different layers of the proposed IT GRC platform. This approach is based on Multi-Agent System (MAS) intelligence to manage the interactions between the distributed systems of the IT GRC platform.**

*Keywords—IT Governance risk; and compliance; information system; multi agent systems*

## I. INTRODUCTION

### A. Governance, Risk Management and Compliance of Entreprise

GRC is the acronym for "Governance, Risk and Compliance". It is a concept that describes the integration of activities to improve the efficiency and effectiveness of many internal functions of organizations. In other words, a comprehensive and systematic approach to governance, risk management and compliance leads to a deeper understanding of the management of what is happening in a business. This approach improves strategy definition, decision-making, risk monitoring and oversight, improved performance, improved internal processes and controls, and so on [1]. As Enterprise Resource Planning (ERP), GRC is becoming one of the most important business requirements of an organization [2], mainly because of globalization [3][4]. We now present a brief description of governance, risk management and compliance definitions.

*1) Governance:* Corporate governance refers to the processes, systems and controls by which organizations operate. A more concrete definition states that "governance is the culture, values, mission, structure, policies, processes and measures through which organizations are directed and controlled". ISO / IEC 38500 subdivides IT governance into three main tasks: To evaluate, direct and monitor the implementation of plans and policies in order to meet the objectives of the company.

*2) Risk:* Risk definitions generally refer to the possibility of loss or harm created by an activity or by a person [4]. Risk management aims to identify, assess and measure risks and develop counter measures to address them, while communicating risk decisions to stakeholders. Typically, this does not mean eliminating risk, but rather seeking to mitigate and minimize impacts. From the point of view of the GRC, the most appropriate concept of Enterprise Risk Management (ERM): "Enterprise risk management is a process implemented by a consulting entity, Management, and other personnel used to establish the company-wide strategy to identify potential events that may affect the organization and manage the risk to provide reasonable assurance regarding the achievement of the organization's objectives " [5]. A well-structured risk management should be aligned and linked to both governance and compliance activities in order to achieve benefits such as better decision-making and increased confidence between the parties Regulatory compliance.

*3) Compliance:* Compliance means not only the establishment of laws, regulations and standards, but also contractual obligations and internal policies [4]. Compliance must ensure that the organization meets all its obligations, and therefore operates within defined prescribed and voluntary limits. The diversity of activities, processes and behaviors that are related to compliance can be very large. But if organizations can manage all these activities, they will operate more efficiently, compete more effectively, and build their

brands in the market. Governance, risk management and compliance as separate concepts are not new [6], but the activities of each discipline share a common set of information, knowledge, methodology, processes and technology. The ultimate goal is to identify, integrate and optimize the processes and activities that are common across the GRC.

### B. Governance, Risk Management and Compliance of Information Technology

For information systems research, a subcategory of the GRC is of particular interest: GRC processes that support the information technology operations of an organization, Called IT-GRC (see Fig. 1).

Research in the field of information systems considers that the integration of governance, risk and compliance is interesting from two main perspectives [8].

First, the IT GRC is seen as a mechanism: How information systems can support the integration of the GRC (business) into the activities of an organization, and how the integration of The GRC can be applied to the information technology of an organization? IT GRC is better understood as a subset of the GRC that supports IT operations in the same way that the GRC as a whole supports business activities. It is aligned with the IT activities and the overall GRC strategy of an organization. Integration of IT governance, IT risk management and compliance has not yet been adequately addressed. Since more than half of GRC publications deal primarily with software technology [8], it can be assumed that there is great potential for integration in technology.

The review of the literature reveals that research priorities in the IT GRC field have not emerged so far, and that a wide variety of aspects ranging from a powerful technical consideration involving the development of a IT GRC application.

Pedro Vicente [7] proposes a business architecture that describes the integration of the main IT governance processes, IT risk management and IT compliance based on a process model for IT GRC. The latter is considered the first process model for IT GRC, it was proposed by the analysis and combination of three references that treat GRC as a separate subject: a process model of ISO / IEC 38500: 2008 for IT Governance; The COSO ERM framework for risk management; and a generic model for IT compliance. Although the process model is directed at IT, it takes into account only three frameworks of good practices, dropping the benefits of standard multitudes and existing standards in this area [9].

Puspasari has created a tool that combines governance, risk management and compliance of information technology [10]. This application consists of managing the following processes: policy management process, risk management, compliance, audit management, business continuity, disaster recovery planning and incident management. Each domain represents a module in the proposed application. The architecture proposed by Puspasari responds to a specific need that is the Bank XYZ who hopes to manage the risks by complying with the regulations associated with this process. Therefore, this

architecture cannot meet all companies and SI environments. In addition, it supports only process management in relation to risk management. Moreover, it does not follow the recommendations of any good practice guidelines.

### C. Positioning of Good Practice Guidelines

As noted by Johannsen and Croeken [11] (see Fig. 2), several frameworks are interdependent and some of their aspects overlap. It is important, however, to identify the appropriate standard to support the appropriate level of IT GRC requirements, for example:

- Help IT managers make the right decisions.
- Define and regulate service management processes.
- Deploy these processes and required procedures, job instructions and monitoring functions.

From an academic point of view, these benchmarks of good practice can be considered as an interesting subject of scientific research, not only because these models are widespread, but also because they integrate enormous consolidated knowledge.

The approaches we have cited (frameworks, standards and best practices) are incomplete with respect to the management of all IS activities of the GRC. Some processes are not covered by certain approaches, and no approach covers all processes related to the IS management of the GRC. This means that the approaches are not complete but fragmented.



Fig. 1.   IT GRC is a Sub-unit of the Entreprise's GRC.



Fig. 2.   Classification of Best Practice Standards.

This is probably due to the fact that the approaches have been produced with the objective of meeting a specific need for governance, risk management of compliance without taking into account all aspects of these three disciplines. The most comprehensive approach is COBIT. However, the functionalities are partial for the IT Governance because the COBIT approach remains generalist. COBIT can be used at the highest level of IT governance, providing a global control framework, based on a computer process model that is generically tailored to each company. There is also a need for detailed, standardized practitioner processes. Specific practices and standards, such as ITIL and ISO / IEC 27002, cover specific areas and can be mapped to the COBIT framework, thus providing a hierarchy of guidance documents.

It should be noted that today there is no IT GRC approach covering the entire IT GRC needs. The objective of this work is motivated by this evidence. Our intention is to address the lack of a comprehensive and structured vision of the underlying concepts of the IT GRC on the one hand and of the IT GRC processes on the other.

In recent years, an array of ITIL (IT Infrastructure Library) or COBIT: Control Objectives for Information and related Technology, as well as internal frameworks, Microsoft operations framework (MOF), ITSM Hewlett-Packard) and ITPM (IT Process model of IBM) were developed. These frameworks, which are also summarized under the theme of Information Technology Governance, describe the objectives, processes and organizational aspects of IT management and control. These best practice models were developed based on practical experiences in IT organizations.

These numerous frameworks that exist on the market make it possible to optimize the functioning of the information system. They offer considerable inputs, but also a very large number of elements not applicable to certain scenarios some organization some systems.

### D. Problematic

We are addressing a twofold challenge to respond to the needs of companies on the adaptation of the IT GRC and the choice of the best framework of good practices to implement the IT processes and generate the action plan, and on the other hand, the management of information workflows (communication) in order to meet business needs, namely from the expression of the need to the implementation of the action plans of the associated processes. We propose in this work to study the tools that help good governance, risk management and compliance of information systems. The lack of effective solutions of this kind (adaptable to any business and environment) is a fundamental problem that deserves further study and raises several research questions:

- Steps in setting up the IT GRC are not clarified

What are the steps that structure the implementation of the IT GRC? What is the nature of these steps?

- IT governance faces changing objectives. Despite this, the maintenance of good governance over time is little taken into account.

Decision-making is often cited as a key element guiding evolutionary actions. How can we then grasp the concept of decision-making in order to maintain good governance over time? What are the impacts of decisions on Information System objectives and on Information System in general?

- The adoption of best practice guidelines until now does not take into account the parameters and constraints of each company

What are the criteria for choosing the best framework that should enable to support processes activities and processes related to the core business of the company?

- Implementation of end-to-end IS activities cannot be considered without effective interactions management. In spite of this, the consideration of a communication system that manages the workflows is little considered.

What is the nature of the interactions that a communication system must support in managing GRC-related processes from the expression of needs to the generation of action plans? And what are the technologies to be used to achieve this result?

We propose to deal with the following problem:

How can IT GRC processes be managed effectively to meet the strategic needs of information system? What is the best framework of good practices to implement the activities of these processes? How can interactions and information workflows be managed to build a platform to support good governance, risk management and information systems compliance?

### E. Research Methodology

Our proposal for the construction of an IT GRC platform is based on the following:

- An understanding of the nature of the IT GRC implementation process

- An IT GRC implementation model, or modeling the architecture of the IT GRC platform, the proposed platform is a smart, distributed, multi-frameworks solution that provides good governance, risk management and compliance of information technology within a company, including a set of distributed systems that:

  o Assures and intelligently evaluates the alignment of the company's business objectives with the IS objectives and strategy,

  o Manages IT processes,

  o Prioritize IT investments in line with business value.

  o Manages and evaluates IT risks,

  o Ensures compliance with the legal framework,

  o Choose the best reference system for the Governance, Risk and Compliance of Information Systems to perform the tasks mentioned above

o Update frameworks according to the latest versions available on the market,

The platform is based on the most widely used standards and methods of Governance, Risk and Compliance of Information Systems (COBIT, ITIL, PMBOK, ISO27001, ISO27002, ISO27002, ISO27005, MEHARI, EBIOS)

- The modeling of a communication architecture, which manages the interactions between the distributed systems of the IT GRC platform, ensuring end-to-end communication between the different layers of the solution. It comprises a communication block per layer for the particularity of the workflows of each layer and the specificity of the processing to be launched before redirecting the information flow to the following layer.

In this way, we wish to respond to the needs and current failures of IS engineering research on the formalization of the IS concepts and the need observed on the adaptation of the frameworks and references of the IT GRC.

The next section presents the global IT GRC solution proposed to address the problematic. Recall that the latter refers to the observation of a lack of adaptation of the processes of the Governance, the management of the risks and the conformity of the Information Systems to the needs of the companies.

## II. PROPOSED GRC IT PLATFORM ARCHITECTURE

The analysis of the literature shows the weakness of research investigations in this field. We address this problematic by proposing an intelligent and distributed Platform of Governance, Risk and Compliance based on Multi-frameworks of Information Systems consisting of:

- A strategic system whose objective is to ensure and evaluate the alignment of the company's business objectives with the IS objectives and strategy.

- A decision-making system whose objective is to choose the best reference system for Governance, Risk and Compliance of Information Systems.

- A communication system that manages the communication (flow of information) between the different systems of the GRC IT platform in a smart way.

- Processing systems whose objective is to manage the IT processes according to the reference system chosen by the decision-making system.

- An updating system which serves to update the frameworks of good practices considered by each processing system.

Fig. 3 illustrates the overall architecture of the Distributed IT GRC platform based on multi-agent systems.

The proposed architecture consists of the following layers:

### A. Strategic Layer

The strategic layer allowing persisting the dynamic and static configurations of the company, to encapsulate the objectives related to the Information Technology of the various departments of the company and to correspond them with the IT objectives and the adequate computer processes, Edit the matrix of responsibilities, the maturity model and the control objectives of the strategy in question. To ensure these functionalities, the strategic layer is based on inter-organization workflows based on multi-agent systems to ensure the orchestration of workflows from different independent and non-pre-packaged business departments for a common final objective for one or more initial business objectives. Moreover, it puts at the service of its users a semantic engine allowing translating their business objectives into a query that can be interpreted by all IT GRC frameworks. Requests are archived for the enrichment of the framework set to initial state.

The strategic layer of the platform is based on the EAS-STRATEGIC system making it possible to make the static configuration of the company necessary for all the components, namely the general information, the resources, the departments, the certifications obtained or prepared, the constraints, strategies implemented ... etc., in addition to the dynamic configuration consisting of expressing the current specific business objectives of a given department. The persistence of the configuration, the translation of the business objectives expressed in language comprehensible by all IT GRC frameworks and the intelligent correspondence between business objectives, IT objectives and IT processes. In order to serve as a reference to the objectives expressed by the business manager, the choice is based on the COBIT framework for which a multi-agent decomposition has been made, which will constantly feed the semantic engine, plus requests already processed that are stored at the level of the knowledge base (learning aspect).

At the end of its processing, this layer sends the synthesis of the results to the communication layer for a possible redirection to the processing components for the purpose of specialization.

### B. Decision Making Layer

The decision layer is capable of selecting the best framework of IT governance, risk management and compliance for a request from the strategic layer, capable of detailing the activities and measures to be executed for an IT process according to its category (Governance, Risk, Compliance) based on the company configuration and the IT process evaluation criteria per framework.



Fig. 3. Global Architecture of the EAS-IT GRC Platform.

To respond to these functionalities, the decision-making layer is based on a multi-criteria intelligent choice capable of designing frameworks to be mobilized in order to respond effectively to the user's demand. It offers two decision-making modes: an IT-oriented mode and an activity oriented mode; According to the needs of the company. Each mode is supported by intelligent agents running two algorithms of choice first by criteria and the second by framework.

The decision-making layer of the IT GRC platform is based on the EAS-DECISION system making it possible to make an intelligent choice of the best framework to process a request from EAS-Strategic. A decisional categorization of the IT processes is made at the level of the communication layer and then the two algorithms of choice of the framework are executed by the agents responsible. Any other decision of the processing layer must be redirected to the communication layer, example: choice of the best risk management strategy.

## C. Processing Layer

The processing layer encapsulates each IT GRC framework in an intelligent, stand-alone system that deploys actions and implements all of the framework's recommendations in an interactive way. The interaction is done by sending a specification request to the strategic layer to request static information to be configured or to open an exchange form with a potential user whose answers are redirected to the knowledge base of the System in question.

The processing layer of the proposed IT GRC platform is based on several EAS-Processing processing systems (EAS-ITIL, EAS-PMP, EAS-ISO 27001 ...) which are notified by

EAS-C OM after recovering the decision from EAS-Decision. Each EAS-Processing system encapsulates a specific IT GRC framework and puts it into production through Intelligent Agents that communicate with each other in order to detail the process acquired in input. For example EAS-ITIL represents the ITIL framework, so once one or more IT processes have to be dealt with this framework the agents of the latter choose the process of an appropriate ITIL cycle with the associated recommendations. A communication with a potential user is possible to detail the request.

## D. Communication Layer

The communication layer provides end-to-end communication between the different layers of the solution in two different modes synchronous by message sending and asynchronous by information sharing, each mode is triggered according to the specificities of the organization and the strategy in question. It comprises a communication block per layer for the particularity of the flows of each layer and the specificity of the processing to be launched before redirecting the information flow to the following layer.

The communication layer of our IT GRC platform is based on the EAS-COM system, which is responsible for exchanging flows and messages between EAS-Strategic, EAS-Decision and EAS-Processing. Two communication modes are involved: communication mode by sharing information and the second by message sending.

This system constitutes the second scientific contribution in this work which we will present in the following section.

## E. Update Layer

The update layer supports updating versions of frameworks of best practices used to periodically upgrade the entire platform. This upgrade is ensured from a correspondence between the processes of the old and the new version, injecting the necessary information to the knowledge bases of the different blocks of the platform. The updating layer of the IT GRC platform is based on the EAS-Updater system which upgrades the versions of all the frameworks deployed to the platform: a correspondence is made from the official documentation between the old and the new version in flat files, an intelligent agent at the level of EAS-Updater loads the received files into the knowledge bases of the different layers.

The IT-GRC platform is a solution based on the concept of distributed systems, based on multi-agent systems (MAS) in its various parts namely user interface, static and dynamic configuration of the organization management profiles, choice of the best framework and processing of processes, it takes advantage of the autonomy and learning aspect of the MAS as well as their communication and coordination of high level.

However, these technological components are difficult to manipulate, or, users lack the skills necessary to use them properly. In this situation, the modeling of communication architecture is necessary, with the aim of adapting the functionalities of the platform to the needs of the users. To help achieve these objectives, it is necessary to develop a functional and intelligent communication architecture that is adaptable and capable of providing a support framework, thus allowing access to the functionalities of the systems independently of the physical and temporal constraints.

A functional architecture defines the logical and physical structure of the components that make up a system and the interactions between them [12][13][14]. If we focus on intelligent and distributed architectures, the main paradigm to consider is the multi-agent system.

EAS-COM is a new architecture focused on product development based on multi-agent systems. It integrates this technology to facilitate the development of a flexible distributed system by taking advantage of the characteristics of interaction between agents to model functional system.

## III. EAS-COM

EAS-COM (see Fig. 4: EAS-COM is represented by the transverse layer of the platform) is a communication system that facilitates the integration of distributed systems of the IT GRC platform. This system must be dynamic, flexible, robust, adaptable to each user's request, scalable and easy to use and maintain. However, this architecture is extensible to integrate the desired processing system, without dependence on a specific programming language. The systems integrated into the IT GRC platform must follow a communication protocol that must integrate. Another important feature is that, thanks to the capabilities of the agents, the developed systems can make use of the learning techniques to manage the decisions previously taken and which are recorded in knowledge bases.

EAS -COM offers a new perspective, where multi-agent systems and Web services are integrated to provide communication needs and leverage their strengths and avoid their weakness.

In two previous works we proposed two architectures of the EAS-COM system based on the two modes of communication between the agents: communication by sharing information and communication by message sending [16][17][18][19][20].

The use of the information sharing mode for the modeling of a communication system within a distributed platform has several advantages that [15] also point out in their work, namely:

- There is no need to treat communication participants directly.

- There is no need for synchronous links between communication participants.

- There is little loss of information.

On the other hand, the use of information sharing, especially for communication between agents of the same ADM, risks accumulating unnecessary data and less communication flexibility. These disadvantages appear during collaboration between the agents of the same subsystem (agents of the Strategic-Com, agents of Decision-Com and agents of processing-Com). These agents need to have freedom of expression to achieve the desired goal of each subsystem.

On the other hand, the use of the information sharing mode to establish the communication between EAS-COM and the other distributed systems of the IT GRC platform raises the problem of the synchronization of execution of the requests of these systems by our system Communication.

Concerning the second proposal was the use of the message sending mode for the modeling of a communication system within a distributed platform. This proposal has several advantages, namely:

- Freedom of expression

- Flexible communication

- Parallelization

On the other hand, the use of message sending, especially for messages containing the most relevant information (IT service requested, categorized IT service, IT service decided and result of processing), risks losing this information and therefore the Workflow of the communication layer will be interrupted. Therefore, focusing only on sending a message is likely to saturate communication, especially between the three EAS-COM (Strategic-Com, Decision-Com and processing-Com) subsystems. These three multi-agent systems need to have a permanent backup of the data that we deemed most relevant to achieve the desired goal of each subsystem.

The architecture of the hybrid communication system that we are going to propose in this section combines the two modes of communication: information sharing and message sending. This solution will overcome the shortcomings

encountered in the two previous architectures (see the evaluations of the two proposals).

The exclusive use of one of these two modes of communication does not provide a persistence of the data to be exchanged. However, in view of their complementarities in this context, their association provides relevant results in the coordination and control of the interactions between distributed systems of the IT GRC platform, between EAS-COM subsystems and these latter. Therefore, a high level of interaction is achieved in a smart way.

In this third version of the architecture of the EAS-COM (see Fig. 5), we have combined the two modes of communication.



Fig. 4. Basic Architecture of the EAS-COM Communication System.



Fig. 5. Hybrid Architecture of the EAS-COM Communication System.

In order to solve the problem of managing communication workflows within the IT GRC platform, we break down the EAS-COM system into subsystems. Each subsystem is concerned with the execution of a specific task of the whole communication problem.

There is a close link between the choice of agents and the objectives for which they are designed. Since we intend to manage workflows between components of the IT GRC platform based on the importance of their content to users, we need to perform the following main tasks:

*1)* The categorization of IT services received from the strategic layer.

*2)* Request and receive the processing of the decision (interaction with the decision layer) in relation to the best references.

*3)* Management of processing systems (sending of IT services to be processed and reception of processing results) taking into account the quality of their processing and their performance. Each task can be assigned to an agent or group of agents.

- We call the multi-agent system dedicated to the categorization of IT services (interaction with strategic layer) "Strategic-com". It contains task agents (1).

- We call the assigned multi-agent system to communicate with the Decision-Com decision-making layer. It contains agents responsible for executing task (2).

- We call the "processing-com" multi-agent system for managing the processing of IT services (interaction with the processing layer). The agents of this multi-agent system are responsible for task (3).

### A.  Strategic-COM

The Strategic-COM subsystem ensures communication with the strategic layer represented by the EAS-STRATEGIC system. This one translate strategic needs of the user in terms of IT service. The deduced IT services are redirected to the Strategic-COM subsystem which categorize the IT processes included in the IT service requested. Categorization consists of associating each IT process into one or more good practices/ frameworks to manage activities of the IT process. Here after the diagram explaining the procedure for categorizing an IT service received by the strategic layer (see Fig. 6).

The IT Service categorization procedure is as follows the IT service received must be divided according to the IT processes that contain.

Each IT process is associated with one or more good practice references according to the discipline to which it belongs (IT Governance, IT Risk Management, IT compliance)

The elements of the matrix are constructed as the following form: {Proc i, (Ref 1, Ref 2, ... Ref n)}.

These ones are then grouped together in order to construct the final matrix as: {{Proc a, (Ref i, Ref j, ... Ref n)}, {Proc b, (Ref i, Ref j, ... (Ref i, Ref i, ... Ref n)}}. This matrix represents the categorized IT service ready to be processed by the second EAS-COM subsystem.

We defined three types of agents: Collector Agent, Manager Agent, Constructor Agent (see Fig. 7).

*1)* *Collector agent:* Collector Agent performs an organizational task. It checks the structure of the web services received, it classifies them according to the date of their creation by the user (date of creation is specified in all IT service). At the end of its processing, it transfers the IT Services to the Manager Agent.

*2)* *Agent manager:* Manager Agent is the heart of Strategic-COM. It categorizes IT services by associating each IT process with one or more appropriate frameworks for its implementation. At the end of the processing, it merges the elements of the matrix which will constitute the IT service categorized as {IT process, {ref1, ref2, ..., refn}}. This result will be transferred to the builder agent.

The Agent Manager has a knowledge base, this one depends of the mapping of the COBIT processes with the other frameworks. This mapping list will be fed from the IT GRC platform.

*3)* *Constructor agent:* The objective of this agent is to provide a comprehensible representation of the IT service, while preserving as much as possible the IT service setting data (the user creating the IT service, the date of its creation, Priority of IT processes ...). To achieve this goal, it retrieves the result of categorizing the IT processes provided by the Manager Agent and constructs the final matrix that represents the categorized IT service that will be sent to the decision layer (EAS-Decision) as a web service.



Fig. 6.   Procedure for Categorizing an IT Service by Strategic-COM.



Fig. 7.   Strategic-COM Agents.

Fig. 8.    Distribution of Strategic-COM Agents According to their Tasks.

In the following figure (Fig. 8), we present the distribution of Strategic-COM Agents according to their tasks.

### B. Decision-COM

The Decision-COM ensures communication with the decision layer represented by the EAS-Decision system (see Fig. 9). This communication consists of sending the categorized IT service to the decision-making layer represented by the EAS-DECISION system. Once the decision is taken in relation to the best frameworks to be associated with each of the IT processes included in the IT service, Decision-COM receives the result of the decision, represented by the IT service decided. The latter must have the following format: {(Proc a, ref i), (Proc b, ref j), ..., (Proc z, ref n)}.

*1) DD agent:* This agent ensure the communication of the IT service with the decision layer. It receives the categorized IT service from the Constructor Agent and translates it into a web service so that it can be sent to the decision layer via network (knowing the IP address of the server in which EAS-Decision runs) and it remains Listening to receive the result of the decision. Once it is received, it is transferred to the processing-Com subsystem for processing.

### C. Processing-Com

The Processing-COM ensures communication with the processing layer. Processing systems of the EAS-PROCESSING layer manage the IT processes following the recommendations dictated by the framework chosen by the EAS-COM decision-making system in order to generate the action plans to be implemented to meet the needs users of the IT GRC platform. We defined four types of agents: Agent ComIn, Agent Admin, Agent Directory, ComOut Agent.

*1) Comin agent:* Agent ComIn is a communicating agent. It receives the decision-making IT service from the Decision-com and transfers it to the Admin agent to determine the processing systems capable of managing the IT processes.

*2) Admin agent:* The Admin agent invokes the processing system that is best placed.

If there are several systems that can solve the requested task, the Admin agent has the ability to select the optimal choice. This decision-making capacity in relation to the choice of the processing system depends on the performance of the latter, its execution number, its availability.... This information is stored in its knowledge base which it uses during the resolution of conflicting situations. With each choice made, it communicates with the agent ComOut and determines the best system to trigger.

*3) Agent directory:* The Directory Agent records system processing reports, as well as the information about them (system performance, number of execution...).

*4) Comout agent:*  Notifying and triggering processing systems that can handle all the processes of an IT service is a complex task that can lead to additional processing time, and therefore can slow down this task. In this step, we propose a new approach whereby process triggering of IT service processes can be partitioned. Our idea is to trigger the set of processing systems chosen to implement the processes of the same IT service. During this trigger, the ComOut agent receives the list of processing systems to be notified. This list must contain the information of these systems, namely the name of the system, the description, the IP address of the server in which the processing system is running.

This method provides simultaneous processing of all processes included in the IT service. However, there may be situations where multiple processing requests are not permitted, including requests to process multiple processes through the same processing system, which could significantly reduce the processor's performance. In these cases, the Admin agent instructs the ComOut agent to check the status of the affected system and notify it that it is busy and cannot accept other requests until it finishes.

In the following figure (Fig. 10), we present the distribution of the Processing-COM Agents according to their tasks.



Fig. 9.    Decision-COM' Agent.



Fig. 10.  Distribution of the Processing-COM Agents According to their Tasks.

We have defined three subsystems that make up our EAS-COM system: Strategic-Com, Decision-Com and Processing-Com: they are multi-agent systems made up of several agents that interact to guarantee the achievement of the goals to which they are Affected. During this interaction, agents intervene to manage possible workflows. To achieve their objectives, our agents act according to their knowledge and skills. In Table II and Table III, we summarize the main characteristics of our agents in the Annexure B.

## IV. IMPLEMENTATION AND EXPERIMENTATION

The AUML modeling of the EAS-COM system and the realization of the simulation platform were followed by the implementation of the communication and management system for the interactions between the distributed components of the IT GRC platform that run on networked machines. The IT GRC platform was tested on a local network and an internet network. This platform is based on the hardware architecture (see Fig. 11) composed of:

- 1 PCs representing the EAS-STRATEGIC application server;

- 1 PCs representing the EAS-DECISION application server;

- 1 PCs representing the EAS-COM application server;

- 3 PCs on which the software components of the processing systems of the EAS-PROCESSING layer are installed respectively;

A router is through which these PCs are connected. The figure below illustrates this architecture:

The EAS-STRATEGIC system is in direct contact with the user of the IT GRC platform. It makes it possible to make the static configuration of the company necessary for all the components: general information, resources, departments, certifications obtained or prepared, constraints, strategies implemented ... etc, in addition to the dynamic configuration of expressing the current specific business objectives of a given department. This system allows the users of the platform to translate the business objectives expressed in language comprehensible by all GRC IT frameworks and the intelligent correspondence between business objectives, IT objectives and IT processes. Once performed, it sends the business requirement expressed in IT processes (IT Service) through a RESTful service by specifying the IP address of the PC on which the EAS-COM system is running and using JSON as the format of data. Here is the request sent by the EAS-STRATEGIC system in the execution of this version of experimentation:

http://10.10.19.147:8080/EAS-COM/?query{"idService":50,"user":utilisateur1,"date":"Jul 15, 2015 6:45:26 PM","seviceHasItprocesses":[{"idService":50,"idItprocess":1," itprocessname":"PO1"},{"idService":50,"idItprocess":2,"itproc essname":"PO2"},{"idService":50,"idItprocess":4,"itprocessna me":"PO4"}, {"idService":50,"idItprocess":8,"itprocessname":"PO8"}], "priorite":[["PO8",4],["PO4",3],["PO2",2],["PO1",1]]}



Fig. 11. Experimental Platform Architecture.



Fig. 12. Launching the EAS-COM GUI.

This query starts our EAS-COM system (see Fig. 12). It retrieves the requested IT service, and displays its data in a table. Then, it proceeds to the categorization of the IT processes included in the IT service by consulting the knowledge base. The latter follows the mapping between the COBIT processes and the ITIL, PMBOK, ISO 27001 and ISO 27002 frameworks (see Table I in Annexure A). In our case:

- PO1 is associated with ITIL

- PO2 is associated with ITIL, ISO 27001 and ISO 27002

- PO4 is associated with ITIL, ISO 27001 and ISO 27002

- PO8 is associated with ITIL, PMBOK, ISO 27001 and ISO 27002

The categorization result is then displayed in the second table (see Fig. 13).



Fig. 13. Categorization of Requested IT Service.

The EAS-COM system prepares the request to send to the decision-making system with respect to the frameworks associated with the IT processes, EAS-DECISION, this request translates the categorized IT service, the latter has the following format:

http://10.10.19.167:8080/EAS-decision/?query {"IdService":50,"user":utilisateur1,"date":" Jul 15, 2015 6:45:26 PM"serviceIT":[["PO1","ITIL"],["PO2","ITIL","ISO27002","ISO27001"],["PO4","ITIL",          ,"ISO27002","ISO27001"],["PO8","ITIL",          "PMBOK","ISO27002","ISO27001"] "Priorite":[["PO8",4],["PO4",3],["PO2",2],["PO1",1]]}

The EAS-DECISION system makes it possible to make an intelligent choice of the best framework of the four IT processes included in the categorized IT service (PO1, PO2, PO4, PO8). The result of the decision is sent via a RESTful service specifying the IP address of the PC running EAS-COM (10.10.19.147), and http as the transport protocol. Here is the query from EAS-DECSION to EAS-COM:

http://10.10.19.147:8080/EAS-COM/?queryChoice={"choice":[["PO1","ITIL"],["PO2","ITIL"],["PO4","ISO          27002"],["PO8","ISO PMBOK"]],"IdService":1,"user":utilisateur1,"date":"Jul Jul 15, 2015 6:45:26 PM"}

The reception of this request by the EAS-COM communication system allows the EAS-COM communication system to dissect selected frameworks, the result is displayed in the decision results table (see Fig. 14)

EAS-COM then sends the processing requests. To do this, it associates to each IT process an appropriate processing system according to the reference system chosen by the EAS-DECSION system. In our case:

- PO1 will be managed by the EAS-ITIL processing system

- PO2 will be managed by the EAS-ITIL processing system

- PO4 will be managed by the EAS-ISO 27002 processing system

- PO8 will be managed by the EAS-PMBOK processing system

The requests (notifications) to send to the processing systems are as follows:

The request sent to EAS-ITIL:

http://10.10.19.110:8080/EAS-ITIL/?query       {"process": "PO1", "IdService":50,"user":utilisateur1,"date":"Jul 15, 2015 6:45:26 PM" }

The request sent to EAS-ITIL:

http://10.10.19.110:8080/EAS-ITIL/?query       {"process": "PO2", "IdService":50,"user":utilisateur1,"date":"Jul 15, 2015 6:45:26 PM" }

The request sent to EAS-ISO27002:

http://10.10.19.111:8080/EAS-ISO27002/?query {"process":                               "PO4", "IdService":50,"user":utilisateur1,"date":"Jul 15, 2015 6:45:26 PM" }

The request sent to EAS-PMBOK: http://10.10.19.112:8080/EAS-PMBOK/?query       {"process": "PO8", "IdService":50,"user":utilisateur1,"date":"Jul 15, 2015 6:45:26 PM" }

These queries will allows to launch the interfaces of the processing systems in order to follow the execution of the execution of the four IT processes. (Note: the EAS-ITIL processing system is executed twice but each execution concerns a different IT process: PO1 for the first execution and PO2 for the second).

Each processing system deploys the actions and implements all the recommendations of the framework in an interactive way. Once it completes its processing, it sends the processing report of the requested IT process to EAS-COM by specifying the download link (the report to the format of a PDF file stored in the server in which the treatment) (see Fig. 15).



Fig. 14. Receiving the Result of the Decision.



Fig. 15. Receiving Processing Results from Processing Systems.

The processing reports are received by EAS-COM, which stores them in the database of the platform. EAS-COM calculates the parameters of each of the four processing systems: performance, quality, and number of execution.

## V. ANALYSIS OF RESULTS

The simulation and experimentation described in this section makes it possible to highlight the interest of the support provided by EAS-COM to face the design problems of the applications of the IT GRC platform. EAS-COM is designed to deal with the different problems encountered by distributed systems implemented in our platform.

### A. Interest in Decoupling Functionality

The problem of decoupling functionality appears in the proposed IT solution GRC. EAS-COM addresses these issues, in particular through the use of a service-oriented (RESTful) approach and the use of agent technology:

*1) Distribution:* Distribution appears in the architecture of the IT GRC platform. The associated problems are largely handled at the service infrastructure level. The interest of EAS-COM in this case is therefore the possibility to build on these existing infrastructures and thus benefit from the solutions they provide to manage the decentralization, security and reliability of communications.

*2) Reusability:* The problem of reusability also appears in the IT GRC platform. On the one hand, applications have been developed primarily from existing functionalities. On the other hand, certain functionalities such as those of processing systems can be used in several applications. This problem is partly addressed by the use of an approach-oriented service, but EAS-COM increases the reusability by integrating an explicit representation of the context in the descriptions of the functionalities.

*3) Heterogeneity:* Two types of heterogeneity appear in the applications presented: the heterogeneity of the functionalities and the heterogeneity of the infrastructures of these applications. EAS-COM addresses the heterogeneity of functionalities through the use of the service-oriented approach, the heterogeneity of application infrastructures by making it possible to integrate these systems without taking into account its programming language.

### B. Interest in Robust IT Platform GRC

The problem of application robustness is present in the IT GRC platform.

*1) Deployment:* All applications presented in the IT GRC platform are defined in an abstract way and dynamically deployed in a given environment. EAS-COM exploits in particular the mechanisms of assembly of functionalities proposed by the applications integrated in the IT platform GRC.

*2) Breakdown:* Fault tolerance is not specifically detailed, but it appears in the case of the unavailability of one of the processing systems. In particular, we mentioned that when EAS-COM chooses a processing system to manage an IT process and that system disappears or fails, it is possible to use an alternate functionality (choose another processing system that Can take over the management of the same IT process according to the recommendations of the same reference system decided by EAS-DECISION).

*3) Evolution:* The evolution appears in the case of the EAS-Processing layer, in which new processing systems appear gradually. These systems are supported by EAS-COM and integrated without modification of the general architecture of EAS-COM. EAS-COM can thus take care of the evolution of an attentive environment without requiring internal modification. This capability is based on the presence of Admin agents capable of interpreting the descriptions of the new processing systems in its knowledge base.

## VI. CONCLUSION

We proposed a new intelligent distributed platform of Governance, Risk Management and Compliance of Information Systems based on the multi-agent system. In order to adapt the functionalities of the platform to the needs of the users and to help achieve its objectives, it is necessary to develop a functional and intelligent communication architecture that is adaptable and capable of providing a support framework, Accessing the functionality of the IT GRC platform's distributed systems regardless of physical and time constraints. The architecture of the proposed intrusion detection system is based on a new detection model consisting of two independent analyzers using a new functional approach. EAS-COM is a communication architecture dedicated to managing the interactions and information flows between the distributed systems of the IT GRC platform, focusing on the development of products based on multi-agent systems. It integrates this technology to facilitate the development of a flexible distributed system by taking advantage of the characteristics of interaction between agents to model functional system. This approach is based on the intelligence of Multi-Agent Systems (SMA). Intelligent agents, distributed across the three subsystems that make up EAS-COM, cooperate and communicate to effectively manage the IT needs of IT users. To manage this communication, we have established three versions of the architecture: the first is based on the information sharing paradigm, the second is based on the mode of sending messages, and the last one we opted for the implementation, is based on the combination of these two communication modes (hybrid communication architecture).

We subsequently realized an experimentation of the IT GRC platform, implementing our communication system. This system was concretized and validated by the actual tests. It uses web services (RESTful) to interact with components of the general platform that are connected to a local network or an internet network. As for the execution of internal functionalities, it relies on the technology of multi-agent systems by deploying different types of agents who communicate and interact with one another in order to achieve the intended objectives.

In perspective, we continue our work to finalize the experimental platform, adding other processing systems and ensuring their implementations in the platform through the

communication system, and then submitting to real tests. Then we will expand the IT platform GRC in such a way to set up a layer of change management and performance that will set up the action plans generated by the processing systems. To do this, we will adapt our communication system to connect this layer to the existing components of the IT GRC platform. Finally evolve into a marketing platform.

REFERENCES

[1] Llanaj, G. (2010). Meeting the Challenges of Governance, Risk and Compliance.

[2] Gill, S., & Purushottam, U. (2008). Integrated GRC-is your organization ready to move. Governance, risk and compliance. SETLabs Briefings, 37-46.

[3] Frigo, M. L., & Anderson, R. J. (2009). A strategic framework for governance, risk, and compliance. Strategic Finance, 90(8), 20.

[4] Tarantino, A. (2008). Governance, risk, and compliance handbook: technology, finance, environmental, and international guidance and best practices. John Wiley & Sons.

[5] COSO Enterprise Risk Management--Integrated Framework: Application Techniques. Committee of sponsoring organizations of the treadway commission, 2004.

[6] PricewaterhouseCoopers. (2004, April). IAB Internet Advertising Revenue Report, 2003 Full-Year Results [PDF]. Internet Advertising Bureau. Retrieved 9 May, 2004, from the World Wide Web:http://iab.net/resources/ad_revenue.asp

[7] Vicente, P., & da Silva, M. M. (2011, June). A conceptual model for integrated governance, risk and compliance. In International Conference on Advanced Information Systems Engineering (pp. 199-213). Springer, Berlin, Heidelberg.

[8] Teubner, A., & Feller, D. W. I. T. (2008). Informationstechnologie, Governance und Compliance. Wirtschaftsinformatik, 50(5), 400-407.

[9] Racz, N., Weippl, E., & Seufert, A. (2010, May ). A frame of reference for research of integrated governance, risk and compliance (GRC). In Communications and multimedia security (pp. 106-117). Springer Berlin Heidelberg.

[10] Puspasari, D., Hammi, M. K., Sattar, M., & Nusa, R. (2011, December). Designing a tool for IT Governance Risk Compliance: A case study. In Advanced Computer Science and Information System (ICACSIS), 2011 International Conference on (pp. 311-316). IEEE.

[11] Johannsen, W., & Goeken, M. (2011). Referenzmodelle für IT-Governance: Methodische Unterstützung der Unternehmens-IT mit COBIT, ITIL & Co. Dpunkt. verlag.

[12] Franklin, S., & Graesser, A. (1996). Is it an Agent, or just a Program?: A Taxonomy for Autonomous Agents. In Intelligent agents III agent theories, architectures, and languages (pp. 21-35). Springer Berlin Heidelberg.

[13] Bedia, M. G., & Corchado, J. M. (2002). A planning strategy based on variational calculus for deliberative agents. Computing and Information Systems, 9(1), 2-13.

[14] Anastasopoulos, M., Niebuhr, D., Bartelt, C., Koch, J., & Rausch, A. (2005, October). Towards a reference middleware architecture for ambient intelligence systems. In ACM conference on object-oriented programming, systems, languages, and applications.

[15] Krummenacher, R., Kopecký, J., & Strang, T. (2005, October). Sharing context information in semantic spaces. In On the Move to Meaningful Internet Systems 2005: OTM 2005 Workshops (pp. 229-232). Springer Berlin Heidelberg.

[16] Soukaina Elhasnaoui, Aziza Chakir, Meriyem Chergui, Hicham Medromi, Adil Sayouti, "A multi agent architecture for an IT GRC Platform using an intelligent communication system", JDSI' 16 EMSI CASA.

[17] S. Elhasnaoui,H.Iguer,H. Medromi,L.Moussaid, A multi agent architecture for communication workflow management system integrated within an IT GRC Platform using sharing information mode International conference on information technology for organizations Development (IT4OD-2016) Fez, Morocco, March 30th-April 1st, 2016.

[18] S.Elhasnaoui, A.Chakir, M.Chergui, H.Iguer, S.Faris and H.Medromi, "Communication system architecture based on sharing information within an SMA for IT GRC Platform" International Journal of Engineering and Innovative Technology (IJEIT) Volume 4, Issue 3, September 2015 .

[19] S. Elhasnaoui, L. Moussaid, H. Medromi, A. Sayouti , "A Communication System Architecture Based on Sharing Information to Integrate Components of Distributed Multi-Agent Systems within an IT GRC Platform" International Journal of Advanced Engineering Research and Science Vol-3,Issue-12,December 2016.

[20] S.Elhasnaoui1, A.Chakir, M.Chergui, H.Medromi, L. Moussaid, Communication System architecture to integrate distributed systems of an IT GRC platform based on agent technology and web services International journal of soft computing, volume 11 issue 5.

ANNEXEURE A

The grid in Table I is as follows: each cell mentions the specific value (s) of the IT Processes for the framework under consideration. When supported by a best practice we include a star (*) in the corresponding cell. When a COBIT IT process is not supported by a best practice we mention a dash (-).

Each framework provides processes and best practices for the implementation of IS GRC activities. Supporting tools and applications exist to support governance activities, but they are fragmented, dedicated to a specific framework.

TABLE I.     COBIT ALIGNMENT WITH ITIL, PMBOK, ISO 27001, ISO 27002 FRAMEWORKS

| COBIT IT process | ITIL | PMBOK | ISO 27001 | ISO 27002 |
|---|---|---|---|---|
| **PO 1** : Define a strategic IT plan | * | - | - | - |
| **PO 2** : Define the information architecture | * | - | * | * |
| **PO 3** : Determine the technological orientation | * | - | * | * |
| **PO 4** : Define processes, organization and working relationships | * | - | * | * |
| **PO5** : Manage IT investments | * | * | * | * |
| **PO 6** : Communicate management goals and directions | * | - | * | * |
| **PO 7** : Manage IT human resources | * | * | * | * |
| **PO 8** : Manage quality | * | * | * | * |
| **PO 9** : Evaluate and manage risks | * | * | * | * |

| | | | | |
|---|---|---|---|---|
| **PO 10** : Manage projects | * | * | - | - |
| **AI 1** : Find IT solutions | * | * | * | * |
| **AI 2** : Acquire and maintain applications | * | * | * | * |
| **AI 3** : Acquire and maintain applications | * | * | * | * |
| **AI 4** : Facilitate operation and use | * | * | * | * |
| **AI 5** : Acquire IT resources | * | * | * | * |
| **AI 6** : Manage changes | * | - | * | * |
| **AI 7** : Install and validate solutions and modifications | * | * | * | * |
| **DS1** : Define and manage service levels | * | - | * | * |
| **DS2** : Manage third-party services | * | * | * | * |
| **DS3** : Manage performance and capacity | * | - | * | * |
| **DS4** : Provide continuous service | * | - | * | * |
| **DS5** : Ensuring the security of systems | * | - | * | * |
| **DS6** : Identify and charge costs | * | - | - | - |
| **DS7** : Educate and train users | - | - | * | * |
| **DS8** : Manage customer support and incidents | * | - | * | * |
| **DS9** : Manage configuration | * | - | * | * |
| **DS10** : Manage issues | * | * | * | * |
| **DS11** : Manage data | * | - | * | * |
| **DS12** : Manage the physical environment | * | - | * | * |
| **DS13** : Manage the operation | * | - | * | * |
| **SE 1** : Monitor and evaluate IS performance | * | * | * | * |
| **SE 2** : Monitor and evaluate internal control | - | - | * | * |
| **SE 3** : Ensure compliance with external obligations | - | - | * | * |
| **SE 4** : Put in place IS governance | * | - | * | * |

ANNEXURE B

We have defined three subsystems that make up our EAS-COM system: Strategic-Com, Decision-Com and Processing-Com: they are multi-agent systems made up of several agents that interact to guarantee the achievement of the goals to which they are Affected. During this interaction, agents intervene to manage possible workflows. To achieve their objectives, our agents act according to their knowledge and skills. In the following Table I, we summarize the main characteristics of our agents:

TABLE II.      PRINCIPAL CHARACTERISTICS OF EAS-COM AGENTS

| **Goals** | |
|---|---|
| Collector Agent | Receives and transfers IT service coming from Strategic Layer to Manager Agent according to their arrivals. |
| Manager Agent | Categorizes processes included into IT service requested and generates matrix elements {IT Process/frameworks} |
| Constructor Agent | Constructs final matrix which represents IT service categorized, and transfers it to DD agent |
| DD Agent | Receives IT service categorized and transfers it to decision layer. Then it waits to receive the result of decision (IT service decided). |
| Com-In Agent | Receives IT service decided from DD agent and transfers it to admin agent |
| Admin Agent | Associates every IT process to processing system according to decision made and generates {IT process/system processing} |
| Directory Agent | Intervenes to calculate performance of processing system and increment his number of execution after each implementation. |
| Com-Out Agent | Sends requests to processing systems and  receives its responses (processing reports) |
| **Knowledge** | |

| | |
|---|---|
| Collector Agent | Data of IT service requested<br>Vérifier l'identité du Sender |
| Manager Agent | Mapping between IT processes and all existing frameworks |
| Constructor Agent | The result of categorizing IT processes<br>The parameters of the IT service |
| DD Agent | Current state of IT service categorized<br>Identity of EAS-Decision<br>Format of service web to send and to receive |
| Com-In Agent | Current state of IT service decided |
| Admin Agent | List of existing processing systems (mapping between frameworks and systems) |
| Directory Agent | Function to calculate performance<br>Increment execution number |
| Com-Out Agent | Format of web service to send and receive from processing systems<br>IP address of computers where systems are running |
| **Competences** | |
| Collector Agent | Receives IT service from EAS-Strategic<br>Checks IT service requested structure<br>Ranks IT services requested according their arrivals date<br>Sends IT service to Manager Agent |
| Manager Agent | Categorizes every IT process of the IT service requested<br>Associates every IT process to the appropriate IT GRC discipline (IT governance, IT Risk and IT compliance)<br>Associate every IT process to one or more frameworks of best practices<br>Generates matrix elements {IT process/frameworks}<br>Transfers categorization results to Constructor Agent |
| Constructor Agent | Checks matrix elements structure that are received from Manager Agent<br>Constructs final matrix (IT service categorized)<br>Transfers IT service categorized to Decision-COM |
| DD Agent | Verifies IT service categorized structure<br>Transfers IT service categorized to EAS-Decision<br>Receives IT service decided (result of decision)<br>Checks IT service decided structure<br>Transfers IT service decided to Processing-COM |
| Com-In Agent | Verifies IT service decided structure<br>Transfers it to Admin Agent |
| Admin Agent | Identifies IT processes and their associated best framework<br>Consults processing system performance/execution number<br>Associates every IT process to adequate processing system<br>Generates processing system notification {IT process/system Processing}<br>Transfers system processing choice to COM-OUT agent |
| Directory Agent | Stores processing reports into database<br>Identifies run time of each processing system<br>Increments execution number of each implemented processing system |
| Com-Out Agent | Checks notifications structure<br>Sends notification system to processing systems chosen (number of notification=number of processes included into the IT service)<br>Supervises the progression of processing of the IT service<br>Receives the response of each processing system invoked<br>Checks processing reports structure<br>Transfers all reports to Directory Agent |

Fig. 16 summarizes the operation of our multi-agent system. It presents the messages exchanged between the agents when receiving the IT services, the decision on the best framework to apply and the implementation of the IT processes by the processing systems. These messages are summarized in the following table (Table II):

TABLE III.        MESSAGES EXCHANGES BETWEEN EAS-COM AGENTS

| N° | Message | Description |
|---|---|---|
| 1 | Inform (IT service requested data) | Message sent by the collector agent to manager agent. It contains the initial data of IT service requested |
| 2 | Aggregation (demand-categorization) | Message sent by a Manager agent to its instances created according to the number of IT processes included in the IT service requested. |
| 3 | Aggegation (response-categorization) | Message sent by Manager Agent instances to the Manager agent that asks categorization of IT processes. This last one synthetizes all responses. |
| 4 | Inform (final-categorization-data) | Message sent by the Manager Agent to constructor Agent. It contains final data of all IT processes categorized handled by Manager Agent. |
| 5 | Inform (IT service Categorized) | Message sent by Constructor Agent to DD agent. |
| 6 | Inform (IT service decided) | Message sent by DD agent to Com-In Agent. This message contains data of IT service decided. |
| 7 | Inform (IT service decided) | Message sent by Com-In agent to Admin Agent. This message contains data of IT service decided. |
| 8 | Aggregation (demand-processing) | Message sent by Admin Agent to its instances which are created according to the number of IT processes included into the IT service decided. |
| 9 | Help (demand-info-system) | Message sent by the admin agent instances to directory Agent in order to get information about concerned processing system to perform action of processing. |
| 10 | Help (response-info-system) | Message sent by Directory agent to Admin agent' instances to tell it whether there is information about the asked processing system. |
| 11 | Aggregation (response-processing) | Message sent by Admin Agent' instances to the Admin agent that asks association of IT processes to the appropriate processing systems. This last one synthetizes all responses. |
| 12 | Inform (processing IT service-Demand) | Message sent by Admin agent to Com-Out agent. It contains final data of IT service demands of processing: every IT process is associated with the adequate processing system. |
| 13 | Notify (conflicts) | Message sent by Com-Out Agent to Admin agent in order to notify it if a processing system is "busy". |
| 14 | Notify (end processing) | Message sent by Com-Out agent to Directory agent in order to notify it that a processing system has finished its processing. |

Fig. 16. Summary of Communication between EAS-COM Agents (The Numbers Indicate the Messages in the Previous Table).

# Comparison of Accuracy between Convolutional Neural Networks and Naïve Bayes Classifiers in Sentiment Analysis on Twitter

P.O. Abas Sunarya[1], Rina Refianti[2], Achmad Benny Mutiara[3], Wiranti Octaviani[4]

Dept. of Informatics Engineering, STMIK Raharja Jl. Jenderal Sudirman No. 40, Tangerang 15117, Indonesia[1]
Faculty of Computer Science and Information Technology
Gunadarma University, Jl. Margonda Raya No. 100, Depok 16424, Indonesia[2, 3, 4]

*Abstract*—The needs and demands of the community for the ease of accessing information encourage the increasing use of social media tools such as Twitter to share, deliver and search for information needed. The number of large tweets shared by Twitter users every second, making the collection of tweets can be processed into useful information using sentiment analysis. The need for a large number of tweets to produce information encourages the need for a classifier model that can perform the analysis process quickly and provide accurate results. One algorithm that is currently popular and is widely used today to build classifier models is Deep Learning. Sentiment analysis in this research was conducted on English-language tweets on the topic "Turkey Crisis 2018" by using one of the Deep Learning algorithms, Convolutional Neural Network (CNN). The resulting of CNN classifier model will then be compared with the Naïve Bayes Classifier (NBC) classifier model to find out which classifier model can provide better accuracy in sentiment analysis. The research methods that will be carried out in this research are data retrieval, pre-processing, model design and training, model testing and visualization. The results obtained from this research indicate that the CNN classifier model produces an accuracy of 0.88 or 88% while the NBC classifier model produces an accuracy of 0.78 or 78% in the testing phase of the data test. Based on these results it can be concluded that the classifier model with Deep Learning algorithm produces better accuracy in sentiment analysis compared to the Naïve Bayes classifier model.

*Keywords*—*Sentiment-analysis; convolutional neural network; deep learning; Naïve Bayes classifier*

## I.  INTRODUCTION

The needs and demands of the community for the ease of accessing information encourage the increasing use of social media facilities to share, deliver, and search for information needed. One of the popular social media that is widely used by people from various backgrounds is Twitter. Twitter provides facilities with features that are easy to understand for users to publish daily activities, inform a news or fact, and express opinions. This makes Twitter still popular today.

Twitter receives tweets from users' as many as 55 million messages every day [1]. The number of large tweets shared by Twitter users every second, making a collection of tweets can be processed into useful information such as to find out a review or public opinion about a particular product, service, or topic.

The process of processing tweet data to get information requires a method that can find patterns of linkages and classify these tweets, one of which uses sentiment analysis. Sentiment analysis is done to classify data into positive, negative and neutral classes.

The need for a large number of tweets to produce information encourages the need for a classifier model that can perform the analysis process quickly and provide accurate results. One algorithm that is currently popular and widely used today to build classifier models is Deep Learning. Deep Learning Algorithm, one of them is Convolutional Neural Network (CNN) which utilizes the Neural Network concept to carry out many learning processes applied in analyzing and predicting processes. The CNN algorithm is inspired by the workings of human brain neurons which consist of several layers. Each neuron is interconnected and will forward information between layers. Information will go through the iteration and distribution process to each subsequent layer to produce the final output as needed. This iteration process helps the machine to learn and identify information so that it will produce a classifier model that can do the classification process of new data with a good level of accuracy.

The CNN algorithm is generally more implemented to analyze and predict two-dimensional objects (images) but there are several studies that apply the CNN algorithm to one-dimensional objects such as text. One example of research that applies the CNN algorithm in text classification is the research of Yoom Kim (2014) [2]. Based on the research, it was found that the classifier model with CNN algorithm showed good classification performance in text classification (such as sentiment analysis) and since it became the basic standard in text classification.

Based on the above background, in this research will use Convolutional Neural Network (CNN) and Naïve Bayes Classifier (NBC) algorithms in the sentiment analysis process using Twitter data which is expected to produce classifier models with good accuracy. Accuracy results from the CNN classifier model will then be compared with the results of the accuracy of the NBC classifier model; so that it can be seen which algorithm is capable of producing classifier models with better accuracy values.

The limitations of the problem in this research can be formulated as follows:

*1)* The sentiment analysis process was carried out related to the topic "Turkey Crisis 2018" with a tweet obtained from Twitter totalling 45,443 data based on a hash tag (#) relating to the topic taken.

*2)* The tweet used in this research was only an English tweet.

*3)* The process of sentiment analysis and the making of the classifier model in this research use the Python programming language version 3.6.

*4)* Classification of tweet data obtained into positive, negative, and neutral classes using the Text Blob library in Python.

*5)* Using Deep Learning algorithm, Convolutional Neural Network (CNN) and Machine Learning algorithm, Naïve Bayes Classifier (NBC) to build classifier models that can classify sentiments of new data.

*6)* Compare the results of the accuracy values produced by the CNN classifier model with the results of the accuracy of the NBC classifier model.

*7)* Visualize the comparison results of the accuracy from CNN and NBC models into tables and graphs.

The aim of this research is to use Deep Learning algorithm, namely Convolutional Neural Network (CNN) in the sentiment analysis process on English tweets related to the topic "Turkey Crisis 2018" on Twitter data and compare the results of the accuracy values obtained from the CNN classifier model with the results of accuracy values from the Naïve Bayes Classifier model to find out which classifier models produce better accuracy values in text classification.

In the rest of paper, we show briefly the literature review and related work in Section II. In Section III the research methodology is presented. The implementation and results related to our research are also shown in Section IV. The last section is conlusion and future work of our research.

## II. Literature Review

### A. Sentiment Analysis

According to B. Liu (2010) [3], sentiment analysis or opinion mining is a process of understanding, extracting and processing textual data automatically to get information on sentiments contained in an opinion sentence. Sentiment analysis is done to see opinions or trends of opinion on a problem or object by someone, whether they tend to have a negative or positive opinion or opinion.

As in [4], the basic task in sentiment analysis is to classify the polarity of the text in documents, sentences, or features, namely whether the opinions expressed in the document, sentence or feature are positive, negative or neutral.

### B. Twitter

Twitter is a website that is a service from microblog, which is a form of blog that limits the size of each post, which provides facilities for users to be able to write messages in Twitter updates containing only 140 characters. Twitter was founded by three people, namely Jack Dorsey, Biz Stone, and Evan William in March 2006 and was launched in July of the same year.

All users can send and receive tweets via Twitter sites, compatible external applications (cell phones), or with short messages (SMS) available in certain countries. Users can write messages by topic using the # (hashtag). Whereas to mention or reply to messages from other users can use the @ (et) sign.

The characteristics of a microblogging or Twitter, which has a status update commonly referred to as tweet totaling 140 characters shorter than other media; can comment on the tweet made by following by using reply, then it can be written using the RT @ username function; have their own way of sharing photos and videos commonly referred to as tweetpic as in [5].

### C. Naive Bayes Classifier

Naïve Bayes Classifier (NBC) is a text mining method that can be used to solve opinion mining problems. NBC can be used to classify opinions into positive and negative opinions. NBC can function properly as a method of text classifiers.

The Naïve Bayes classification algorithm utilizes the probability theory proposed by British scientist Thomas Bayes, which predicts future probabilities based on past experience. The simple NBC algorithm and its high speed in the training and classification process make this algorithm interesting to use as a classification method. The classification process is usually divided into two phases, namely, learning and test. In the learning phase, some of the data that has been known for the data class is fed to form an approximate model. Then in the test phase the model that has been formed is tested with some other data to determine the accuracy of the model.

In the Naïve Bayes Classiffier algorithm each tweet is represented by a pair of attributes "$x_1, x_2, x_3, \ldots x_n$" where $x_1$ is the first word, $x_2$ is the second word and so on. Whereas $V$ is a set of tweet categories. At the time of classification the algorithm will look for the highest probability of all categories of tweets tested ($v_{MAP}$), where the equation is as follows:

$$v_{MAP} = arg\max_{v_j \in V} \frac{P(x_1, x_2, x_3 \cdots, x_n | v_j) P(v_j)}{P(x_1, x_2, x_3 \cdots, x_n)} \quad (1)$$

For $P(x_1, x_2, x_3, \ldots, x_n)$ the value is constant for all categories ($v_j$) so that the equation can be written as follows:

$$v_{MAP} = arg\max_{v_j \in V} P(x_1, x_2, x_3 \cdots, x_n | v_j) P(v_j) \quad (2)$$

The above equation can be simplified as follows:

$$v_{MAP} = arg\max_{v_j \in V} \prod_{i=1}^{n} P(x_i | v_j) P(v_j) \quad (3)$$

where $v_j$ = category of tweets $j$ = 1, 2, 3, ... $n$, in this study $j$ = 1 indicates a category of negative sentiment tweets, $j$ = 2 indicates a category of positive sentiment tweets and $j$ = 3 indicates a category of neutral sentiment tweets:

$P(x_i | v_j)$ = probability xi in the category $v_j$;

$P(v_j)$ = probability of $v_j$.

For $P(v_j)$ and $P(x_i | v_j)$ it is calculated during the training where the equation is as follows:

$$P(v_j) = \frac{|docs_j|}{|example|} \tag{4}$$

$$P(x_i|v_j) = \frac{n_k+1}{n+|vocabulary|} \tag{5}$$

where

$P(v_j)$ = The probability of each document against a set of documents;

$P(x_i | v_j)$ = The probability of the occurrence of the word $x_i$ in a document with the class category $v_j$ ;

| docs | = number of documents in each category $j$ ;

|example| = number of documents from all categories;

$n_k$ = number of times the frequency of occurrence of each word;

$n$ = number of frequency of occurrence of words from each category

There are several forms of representation of the Naïve Bayes Classifier method, including:

*1) Gaussian naive bayes*: Gaussian Bayes are usually used to represent the conditional probability of the continue feature in a class $(x_i | y)$, and are characterized by two parameters: mean and variant.

*2) Bernaulli naïve bayes*: In Naïve Bayes Bernaulli, weighting is carried out using binaries (0 and 1) in weighting each term, this is different from the calculation of frequency terms that do weighting on each term.

*3) Multinomial naïve bayes*: Multinomial Naïve Bayes assumes independence between the appearance of words in a document, without taking into account the order of words and context of information in sentences or documents in general. Besides this method takes into account the number of occurrences of words in the document.

The Naive Bayes algorithm that is often used for text mining is Multinomial Naive Bayes. Multinomial Naïve Bayes is one of the specific methods of the Naïve Bayes method. Multinomial Naïve Bayes is also a supervised learning machine in the process of classifying text by using the probability value of a class in a document.

### D. Deep Learning

Deep Learning is a branch of science learning based on artificial neural networks (ANN) or it can be said that the development of ANN teaches computers to be able to take actions that are considered natural by humans. For example, to learn from examples. In deep learning, a computer can learn to classify directly from images, sounds, texts or even videos. A computer is trained using data sets labeled and the numbers are very large which can then change the pixel value of an image into an internal representation or feature vector where classification can be used to detect or classify patterns at input input [6][9][10][11].

Deep learning method is a method of learning representation with several levels of representation, where representation forms a neural network architecture field that contains many layers (layers). The deep learning layer consists of three parts, namely the input layer, hidden layer, and output layer. In the hidden layer can be made in layers to find the right algorithm composition to minimize errors in output [6].



Fig 1.    Deep Learning Layers.

Fig. 1 illustrates deep learning layers that have p + 2 layers (p hidden layer, 1 input and 1 output layer). Blue circles represent neurons. There are one or more neurons in each layer. These neurons will be connected directly to other neurons in the next layer [6].

### E. Convolutional Neural Network

Convolutional Neural Network (CNN / ConNet) is one of the deep learning algorithms which is the development of the Multilayer Perceptron (MLP) which is designed to do data into two dimensions, for example: images or sound. Convolutional Neural Network is used to classify the labeled data by using supervised learning method, the way it works is that there is training data and there are variables that are targeted so that the purpose of this method is to group data into existing data.

In general, the CNN layer type is divided into two parts, namely:

*1) Feature extraction layer (feature extraction layer)*: The image that is located at the beginning of the architecture is composed of several layers and in each layer arrangement of the neurons connected to the local region (local region) of the previous layer. The first type of layer is the convolutional layer and the second layer is the pooling layer. At each layer the activation function is applied with its intermittent position between the first and second types. This layer accepts image input directly and processes it until it produces a vector output to be processed in the next layer.

*2) Classification layer*: This layer is composed of several layers which in each layer are composed of fully connected neurons with other layers. This layer receives input from the output of the image feature extraction layer in the form of a vector which is then transformed as in the Multi Neural Network with the addition of several hidden layers. Output results in the form of class accuracy for classification.

Fig 2. Examples of Convolutional Layer Diagrams [7].



Fig 3. Examples of MAX Pooling Layer Diagrams [7].

As in [7], Convolutional Layer is a main core of CNN, where this layer has a collection of filters that can be used to study input images. Through this layer, the feature will be extracted and then proceed to the next layer in order to extract more complex features. Examples of Convolutional Layer diagrams can be seen in Fig. 2 where the input image size given is 28x28 and a 4x4 filter or kernel.

Pooling Layer is a resizing process that is a process to change the size of different input images, one of them is using the MAX operation. This aims to help reduce the number of parameters and calculation times needed when training the network as well as the work of Bui and Chang in [7]. An example of a Pooling Layer diagram can be seen in Fig. 3.

In Fig. 3, the entered image is 4x4 in size and then resized into 2x2-sized image with a depth of 16. Each value is at Max Pooling, for each 4 pixels a maximum value is taken. Seen in Fig. 3 at 4 pixels in blue, the maximum value to be taken is 5. At 4 pixels in red, the maximum value that will be taken is 9. In pixels in green, the maximum value that will be taken is 8. On pixels in orange, the maximum value to be taken is 7. So as to produce a reduced image.

And the third layer on CNN is Fully Connected Layer, where this layer takes all the neurons in the previous layer (Convolutional Layer and MAX Pooling Layer) and connects them to each single neuron that exists, as we can see in [8].

## III. RESEARCH METHOD

The process of designing a classifier model for sentiment analysis in this research consists of five stages:

### A. Data Retrieval

The first stage of the process of designing this classifier model is data retrieval using the Twitter API service. The Python programming language has provided tweepy library that can facilitate retrieval of data from Twitter. Data is then saved in .csv or .txt format.

### B. Pre-Processing

The second stage is pre-processing, namely the stage where the tweets that have been obtained will be extracted and cleaned from noise, namely random or variant errors in measured variables consisting of RT components, hashtag, digits, user (@), punctuation, url, and others components that are considered to interfere with the tweets classification process. Removing noise object is an important goal of cleaning data because noise inhibits most types of data analysis. The flow of the cleaning process can be seen in Fig. 4 below.

Tweets that have been through the cleaning process will then be classified into three categories of sentiment class namely positive, negative and neutral using the TextBlob library. The flow of the data classification process can be seen in Fig. 5 below.



Fig 4. Flow of Data Cleaning.

Fig 5.    Flow of Data Classification with TextBlob.

### C. Model Design and Training

The third stage is the design and training of classifier models, the tweets that have been classified as sentiment classes will be divided into three parts, namely, data train, validation data, and test data. Data train is used to train new classifier models using the Convolutional Neural Network algorithm and the Naïve Bayes Classifier.

### D. Testing Model

The fourth stage is the testing phase of the classifier model that has been trained using test data by looking at the value of the accuracy produced. Calculation of accuracy values is done using Confusion Matrix to see how much accuracy is produced by the two classifier models in the training and testing process so that it can be known which model produces better accuracy in sentiment analysis.

### E. Visualization

In the final stage, it displays the results in the form of diagrams, graphs and tables.

## IV.  IMPLEMENTATION AND RESULTS

### A. Data Retrieval

Retrieving data from Twitter is first done by making a Twitter API connection. The first step that must be done is to create an application on Twitter by visiting the https://apps.twitter.com/ site to get the keys and access tokens used to access the Twitter API.

After getting the key and access token, it takes a Python library that can implement the Twitter API call, one of which is the tweepy library. The next step is to open the Spyder software and install the tweepy library to be able to pull data from Twitter using the Twitter API.

The next step is to retrieve tweets from Twitter based on the hashtag (#) or predefined keywords. Tweets taken are only English-language tweets, taken randomly from ordinary users or Twitter's official media accounts. The topics discussed in this sentiment analysis were "Turkey Crisis 2018" and several hashtags used to search tweets including #TurkeyCrisis, #TurkeyLira, #Turkey, #Erdogan, and #Trump. The tweets that were successfully retrieved were English tweets totalling 45,443 data.

### B. Pre-Processing

Data from Twitter that has been taken next will go through the pre-processing stage which consists of the cleaning process of tweets and the classification process of tweets based on positive, negative, or neutral sentiment classes using the TextBlob library. The purpose of pre-processing data is to transform raw data into a format suitable for analysis.

### C. Cleaning Data

At this stage the cleaning process of tweet data from noise is carried out, namely random or variant errors in measured variables consisting of RT components, hash tag, digits, user (@), punctuation, url, and other components that are considered to interfere with the tweets classification process.

Tweet data obtained from Twitter often contains components that are not needed and can interfere with the classification process of tweets so that the need for deletion of these components. In the Python programming language, the data cleaning process can use the Beautiful Soup library. After going through the cleaning process, the tweet initially amounted to 45,443 to 33,107 clean tweets.

### D. Data Classification using TextBlob

The next stage after cleaning the data tweet is the data classification stage. Tweets that have been cleared from noise components will then be classified to be divided into three sentiment classes, namely positive (1), negative (2) and neutral (0) classes. Data classification at this stage utilizes the TextBlob library. TextBlob classifies tweets into three sentiment classes based on their polarity.

Tweets will be classified into positive sentiment class if the polarity sum of each word in the sentence produces a value greater than 0 it will be labeled 1. The Tweet will be classified into the negative sentiment class if the polarity sum of each word in the sentence produces a value less than 0 it will be labeled 2. Tweets will be classified into neutral sentiment class if the polarity sum of each word in the sentence produces a value equal to 0 it will be labeled 0.

From the 33,107 tweeted and classified tweets, a neutral category of 14,443 tweets, a positive category of 12,142 tweets, and a negative category of 6525 tweets were obtained. The tweet data that has been classified will be equalized for each class of sentiment because the tweet data is uneven and tends to be neutral. Alignment of the number of tweets will follow the amount of data in the sentiment class with the least data, namely the negative class with the number of data 6525. After leveling, the number of tweets for each sentiment class is 6525.

### E. Model Design and Training

The tweet data that has been through the cleaning process and the classification process using the TextBlob library will then be used to build a classifier model using the Convolutional Neural Network (CNN) algorithm and the Naïve Bayes Classifier (NBC) algorithm.

### F. Split Dataset

Tweets that have been classified as sentiment class will be divided into training data, validation data, and test data which will later be used in designing classifier models using the CNN and NBC algorithms.

The data split in this research was done using the Python library, the Scikit-learn library with the split_train_test method. Data will be divided into three parts including:

*1)* Data Train: the data set used for the learning process by the classifier model.

*2)* Data Validation: the data set used to set the parameters of the classifier and provide an unbiased evaluation of a model.

*3)* Data Test: the data set used to assess the performance of the final model.

### G. Designing the CNN Classifier Model

The process of constructing a classifier model for analysis sentiments using the Convolutional Neural Network algorithm consists of several stages, namely importing datasets, dividing datasets, feature extractions using word2vec, tokenization and padding sequences, designing layers in models, model training and evaluation, model testing and visualization. Fig. 6 shows the flow of the design of the classifier model with the CNN algorithm.

### H. CNN Model Training and Evaluation

The training phase is carried out as a process to find the patterns of linkages between input variables and output variables from the data studied so that later this model can be used to analyze sentiment on new data. Based on the data splitting at the beginning, the data train amounted to 15,655 data with a 33.05% negative share, 33.45% positive, 33.50% neutral. The training process will be carried out with 10 epochs and the results of the training model will be stored whenever an increase in the accuracy value is generated at each epoch.

From the eight (8) classifier models generated from the training and validation process, the best model is the model produced at the 3rd epoch because it produces the best accuracy value of 0.89 and a loss value of 0.33 before the classifier model become overfitting.

Table I will display the accuracy value and loss value generated by the CNN classifier model at each epoch during the validation process.

The best classifier model that is produced, namely the model in the 3rd epoch will then be used to test the test data.

### I. CNN Model Trial of Data Test

The CNN Classifier model that has been trained and evaluated in the previous stage will be tested with test data to see whether the resulting accuracy value will be as good as the accuracy value at the training and validation stages. The test data used in this research is the data obtained based on the split process in the initial part consisting of 1957 data with a 34.18% negative, 33.67% positive, and 33.50% neutral. Calculation of accuracy values for the test data is done using Confusion Matrix to determine the value of precision, recall and f1-score generated by the model.

From the result of testing using test data, the CNN classifier model give an accuracy value of 0.88 or 88% with a loss value of 0.33 or 33%. Fig. 7 will display the classification report from the CNN classifier model test in the test data using the Confusion Matrix calculation.



Fig 6.    Design Flow of CNN Model Classifier.

TABLE I.    ACCURACY AND LOSS VALUE OF DATA VALIDATION IN EACH EPOCH

| Epoch | Validation Accuracy | Validation Loss |
|---|---|---|
| 1 | 0.785 | 0.569 |
| 2 | 0.869 | 0.390 |
| 3 | 0.893 | 0.332 |
| 4 | 0.907 | 0.370 |
| 5 | 0.893 | 0.424 |
| 6 | 0.916 | 0.429 |
| 7 | 0.923 | 0.465 |
| 8 | 0.926 | 0.495 |
| 9 | 0.929 | 0.502 |
| 10 | 0.929 | 0.525 |

```
Confusion Matrix

         predicted_netral  predicted_positive  predicted_negative
netral              558                   34                   37
positive             22                  572                   65
negative             40                   36                  593
-------------------------------------------------------------------
Classification Report

              precision    recall  f1-score   support

      netral       0.90      0.89      0.89       629
    positive       0.89      0.87      0.88       659
    negative       0.85      0.89      0.87       669

 avg / total       0.88      0.88      0.88      1957
```

Fig 7.    CNN Model Classification Report For Data Test.

### J.  Visualization of CNN Model Accuracy Results

Visualization is made to make it easier to understand the results obtained from the training process to the trial. Fig. 8 shows a comparison graph of the validation accuracy value of the train accuracy at each epoch.

Based on Fig. 8, it can be concluded that the accuracy value produced by the classifier model during the training process increases at each epoch. While the accuracy value generated during the validation process has increased in the first to fourth epoch, but at the 5th epoch accuracy has decreased. At the 6th epoch, the accuracy value increases again and starts from the 7th to 10th epoch, the accuracy value is stable.



Fig 8.    Comparison Graph of Validation Accuracy against Train Accuracy.



Fig 9.    Comparison Graph of Validation Loss against Train Loss.

Based on Fig. 9 it can be concluded that the loss value generated during the training process has decreased in each epoch while the loss value generated during the validation process has decreased to the 3rd epoch but starting from the 4th to 10th epoch loss values are increasingly experiencing which increase indicates that the classifier model is overfitting. The best model is the model produced at the 3rd epoch because it produces the best accuracy value of 0.89 and a loss value of 0.33 before the classifier model become overfitting.

### K.  Design of NBC Classifier Model

The process of building a classifier model with the Naïve Bayes Classifier algorithm in this research consisted of several stages, namely importing datasets, dividing datasets, feature extractions, conducting model training on several gram n-features, testing the accuracy of each gram n-feature model, validating, testing model and visualization. In Fig. 10 is shown the process flow design of the Naïve Bayes classifier model.



Fig 10.   Design Flow of the Naïve Bayes Classifier Model.

| | negative | positive | netral | total |
|---|---|---|---|---|
| turkey | 4717 | 4566 | 4835 | 14118 |
| the | 3321 | 3562 | 2848 | 9731 |
| to | 2452 | 2441 | 2037 | 6930 |
| in | 1973 | 1765 | 1739 | 5477 |
| of | 1653 | 1884 | 1519 | 5056 |
| and | 1764 | 1791 | 1478 | 5033 |
| is | 1788 | 1778 | 1252 | 4818 |
| you | 1113 | 976 | 789 | 2878 |
| on | 1128 | 891 | 842 | 2861 |
| not | 1217 | 878 | 753 | 2848 |

Fig 11.  Custom Stop Words List.

*L.  NBC Model Training and Evaluation*

The training phase is carried out on several N-gram models to get the model with the best accuracy value which will then be evaluated with the Naïve Bayes Classifier algorithm. N-gram is a method for retrieving bits of letter characters of n from a word. N-gram has three types of processing models in a sentence; the type of processing includes Unigram for separating one word in a sentence, Bigram for separating two words in a sentence, and Trigram for separating three words in a sentence.

Classifier training models will be carried out on the N-Gram model with several conditions, among others, the unigram with stop words, unigram without stop words, and unigram without custom stop words. Custom stop words are stop words derived from 10 words that most often appear on the corpus. In Fig. 11, it shows the custom stop words list in this study.

Based on the training and validation process carried out on the three unigram models with the conditions mentioned, namely with stop words, without stop words, and without custom stop words, the highest accuracy was generated by the unigram with stop words model with an accuracy value of 77.82% with the number feature 3000.

After getting the results that the highest accuracy value is generated from the unigram with stop words model, then an experiment will be conducted to conduct training and accuracy testing on the bigram and trigram with stop words models to see whether there will be an increase in accuracy.

Based on the results of the accuracy obtained from the unigram, bigram, and trigram with stop words training and validation processes, the best accuracy values for each n-gram model are as follows:

*1)* Unigram: on 3,000 features with validation accuracy of 77.82%

*2)* Bigram: on 5,000 features with validation accuracy of 75.27%

*3)* Trigram: on 5,000 features with validation accuracy of 74.71%

The Unigram with stop words classifier model that produces the best accuracy values will then be used to test on the test data.

*M.  NBC Model Trial of Data Test*

Based on the training process and the validation of the NBC classifier model in the previous stage, it was known that

the unigram with stop words model produced the highest accuracy in 3000 features. At this stage, the accuracy of the classifier model will be tested for the test data. The test data used in this research is the data obtained based on the split process in the initial part consisting of 1957 data with a 34.18% negative, 33.67% positive, and 33.50% neutral. Accuracy testing of the data test was done using Confusion Matrix to determine the precision, recall and f1-score values produced by each model.

The unigram NBC classifier model which was tested with the data test gave an accuracy value of 0.78 or 78%. Fig. 12 will display the classification report from the NBC classifier model test in the test data using the Confusion Matrix calculation.

*N.  Visualization of NBC Model Accuracy Results*

Fig. 14 shows a comparison graph of the accuracy results obtained from the training and validation process carried out on the three unigram models with the conditions mentioned, namely with stop words, without stop words, and without custom stop words.

Based on Fig. 13 can be concluded as follows:

*1)* The best accuracy of unigram without stop words is the 13,000 feature with an accuracy value of 74.55%

*2)* The best accuracy is unigram with stop words, namely the 3000th feature with an accuracy value of 77.82%

*3)* The best unigram accuracy without custom stop words is the 3000th feature with an accuracy value of 77.72%

```
null accuracy: 66.79%
accuracy score: 77.82%
model is 11.04% more accurate than null accuracy
-------------------------------------------------------------
Confusion Matrix

          predicted_negative  predicted_positive  predicted_netral
negative          557                 72                 51
positive           70                520                 37
netral            120                 84                446
-------------------------------------------------------------
Classification Report

             precision   recall  f1-score   support

    negative      0.84     0.69      0.75       650
    positive      0.77     0.83      0.80       627
      netral      0.75     0.82      0.78       680

 avg / total      0.78     0.78      0.78      1957
```

Fig 12.  Unigram Classification Report in 3000 Feature against Data Test.



Fig 13.  Comparison of Unigram Model Accuracy with Conditions.

Fig 14.   Comparison of Model Accuracy of Unigram, Bigram, and Trigram.

After getting the results that the highest accuracy value is generated from the unigram with stop words model, then an experiment is conducted to test the accuracy of the bigram and trigram with stop words models to see whether there will be an increase in accuracy. In Fig. 14 can be seen the comparison of the results of accuracy produced by unigram, bigram, and trigram with stop words models.

Based on Fig. 14 can be concluded as follows:

*1)* The best accuracy of Unigram on the 3000 feature with an accuracy value of 77.82 %.

*2)* Best accuracy of Bigram on the 5000 feature with an accuracy value of 75.27 %.

*3)* The best accuracy of Trigram on the 5000 feature with an accuracy value of 74.71%.

*4)* From the three classifier models, the unigram with stop words model produces the best accuracy values.

### O. Comparison of CNN and NBC Classifier Model Accuracy Values in Data Test

The final result of this research is to find out which classifier model produces better accuracy in sentiment analysis. Based on the classification report, the accuracy testing of the test data carried out in the previous stage shows that the CNN classifier model produces an accuracy value of 0.88 or 88% and the NBC classifier model produces the greatest accuracy value with the unigram with stop words model which produces an accuracy value of 0.78 or 78%. The following is a table of the results of the classification report comparison test data test from the two classifier models.

TABLE II.      COMPARISON OF CNN AND NBC CLASSIFICATION REPORTS AGAINST DATA TEST

| | Precision | | Recall | | f1-Score | | Support | |
|---|---|---|---|---|---|---|---|---|
| | CNN | NB | CNN | NB | CNN | NB | CNN | NB |
| Netral | 0.90 | 0.84 | 0.89 | 0.69 | 0.89 | 0.75 | 629 | 650 |
| Poitive | 0.89 | 0.77 | 0.87 | 0.83 | 0.88 | 0.80 | 659 | 627 |
| Negative | 0.85 | 0.75 | 0.89 | 0.82 | 0.87 | 0.78 | 669 | 680 |
| **Total** | **0.88** | **0.78** | **0.88** | **0.78** | **0.88** | **0.78** | **1957** | **1957** |

Based on the comparison results in Table II it can be seen that the results of precision and recall obtained from the CNN classifier model is 0.88 (88 %) while the precision results generated by the Naïve Bayes unigram with stop words classifier with 3000 features are 0.78 (78 %). These results indicate that the classifier model with Convolutional Neural Network algorithm can provide better accuracy results compared to the Naïve Bayes classifier model in sentiment analysis.

### V.   CONCLUSION AND FUTURE WORKS

The sentiment analysis conducted in this research uses English-language tweets obtained from Twitter using the Twitter API related to the topic "Turkey Crisis 2018". The whole sentiment analysis process starts from the data retrieval process, data classification with TextBlob which divides tweets into positive sentiments, negative sentiments, and neutral, and the use of Convolutional Neural Network and Naïve Bayes Classifier algorithms is done using the Python programming language.

The use of Deep Learning algorithm, Convolutional Neural Network in sentiment analysis has been successfully carried out in this research. The model architecture used in constructing this classifier model uses Keras Functional API model with the number of convolutional layers used is 3 layers with the addition of the kernel filter on each layer with the number 100 filters and kernel size will adjust the n-gram concept that will used in each convolutional layer, namely, bigram (2), trigram (3), and fourgram (4). The activation function used in the convolutional layer is ReLu. Three 1D max pooling layers are used in this model architecture to extract the maximum value from each filter. One fully connected layer with dropout is used to process the output from the max pooling layer with a total of 128 neurons. The output layer will consist of 3 neurons with the softmax activation function.

The CNN classifier model that has passed the training and evaluation process produces an accuracy value of 0.89 or 89% and in the test data testing process produces an accuracy of 0.88 or 88%. The accuracy results are then compared with the accuracy of the Naïve Bayes classifier model. This comparison of accuracy shows that the CNN classifier model has better accuracy values than the Naïve Bayes classifier model which produces an accuracy of 0.78 or 78%. From these results it can be concluded that the classifier model with Deep Learning algorithm produces better accuracy in sentiment analysis compared to the NBC classifier model.

Based on the results of the conclusions that have been described, it can be suggested that several things for further improvement and development include:

*1)* Retrieving tweet data from Twitter in greater numbers so that the classifier model can provide better accuracy in sentiment classification.

*2)* Comparing with other deep learning or machine learning algorithms.

*3)* The classifier model that has been built in this research is expected to be made into an application (front-end) either

desktop or website based to be utilized in analyzing sentiments on tweet data.

REFERENCES

[1] F.Ratnawati and E. Winarko, "Sentiment Analysis of Movie Opinion in Twitter Using Dynamic Convolutional Neural Network Algorithm, Indonesian Journal of Computing and Cybernetics System (IJCCS), 2018, 12 (1), https://doi.org/10.22146/ijccs.19237

[2] Y. Kim, "Convolutional neural networks for sentence classification," arXiv preprint arXiv:1408.5882, 2014

[3] B. Liu, "Sentiment Analysis and Subjectivity," in Handbook of natural language processing, 2010, 2, 627-666

[4] A.A Sattikar and R.V. Kulkarni, "Natural Language Processing For Content Analysis in Social Networking," International Journal of Engineering Inventions, 2012, 1(4), 6-9

[5] Madcoms, Facebook, Twitter, dan Plurk dalam Satu Genggaman. Yogyakarta: ANDI, 2010.

[6] Y. LeCun, Y. Bengio, & G. Hinton, "Deep Learning," Nature International Journal of Science: doi:10.1038/nature14539, 2015, May 27.

[7] V. Bui and L.C. Chang , "Deep learning architectures for hard character classification," In Proceedings on the International Conference on Artificial Intelligence (ICAI) , 2016, p. 108.

[8] P. Devikar, "Transfer Learning for Image Classification of various dog breeds." International Journal of Advanced Research in Computer Engineering and Technology (IJARCET), 2016, 5(12), 2707-2715

[9] R. Refianti, A.B. Mutiara, and R.P. Priyandini, "Classification of Melanoma Skin Cancer Using Convolutional Neural Network," IJACSA, 2019, 10 (3), 409-417

[10] A. Esteva, B. Kuprel, R.A. Novoa, J. Ko, S.M. Swetter, H.M. Blau and S. Thrun, "Dermatologist-level classification of skin cancer with deep neural networks", Nature, Vol. 542, pp 115-118, 2017

[11] T.J. Brinker, A. Hekler, J.S. Utikal, N.Grabe, D. Schadendorf, J. Klode, C. Berking, T. Steeb, A.H. Enk, and C.von Kalle, "Skin Cancer Classification Using Convolutional Neural Networks: Systematic Review", Journal of Medical Internet Research, Vol. 20, No.10, pp. 11936- 11946, 2018

# Digital Legacy: Posterity Rights Analysis and Proposed Model for Digital Memorabilia Adoption using Machine Learning

Amit Sudan[1]

ME Scholar

Department of Computer Science Engineering

Chandigarh University, Mohali, India

Dr. Munish Sabharwal[2]

Associate Dean & Professor

Department of Computer Science Engineering

Chandigarh University, Mohali, India

*Abstract*—**The paper informs about the digital legacy and its related concepts of posterity rights and digital memorabilia. It also deals with the right to exercise the digital posterity concerning the social networking profiles (SNP) on social networking sites (SNS). Digital Memorabilia is the compendium of people's social profiles and the digital artifacts accumulated in the name of people in online or virtual world, it can give people an online space to connect to and be remembered online even after their demise, showing the many dimensions of their real world personality. The paper proposes a model using multiple logistic regression technique of machine learning to predict the users that will opt for a digital memorial dependent upon different factors such as age, frequency of using SNPs, awareness about digital assets and digital legacy, awareness about privacy rights concerning digital assets and awareness about rights to posterity.**

*Keywords—Digital assets; digital legacy; digital posterity; digital executers; digital memorabilia; SNP (Social Networking Profiles); SNS (Social Networking Sites)*

## I. INTRODUCTION

From the past 10 years: demise, passing and online loss customs in the world are together forming an expanding field of attention in today's world [5]. Most of the research nowadays is nationally based and directed on discoursed examination of particular fields such as practices related to demise, interment traditions, and crypt traditions [9]. Although, in recent years, the area has expanded and has become more cross-punitive with the introduction of more networks across countries [1]. Work in online demise and memorial customs form a fast developing area of research, which directs on how death and misery are dealt with on several online platforms and social media such as Facebook, Twitter, etc. Moreover, this also questions how the online media [2] may be disguising our ways of mourning and harrowing. Every online platform has different features amongst which the online media share is one of the very important features for sharing and interacting with people whom we don't usually meet [8]. Now, we have entered into the social media phase, where people don't hesitate in uncovering truth and realities of their life to cite an example for the fame #MeToo is a living example. In fact, people share all their emotions, grieve and show support to the people they favor [2]. These memorials also allow people to participate in

their friends' and relatives' funeral process from any part of the world and at any time of the day or night [4]. In some sociologists' views, such people's exhibition of grief is significant for inner recuperation after deprivation [10]. Accessibility of low-cost or free space available online will allow pallbearer to include ample contents such as stories and discussions [15]. Facebook allows users with the chance to keep the deceased aside their lives by sharing posts on their walls during the birthdays and holidays in their lives or the grieving life [7]. These memorials also give the deprived the power to have the deceased's social media page if they want to be remembered of their good memories they once shared online with the deceased [12]. Continuous vows and conveying the feelings towards the person who are no more can be regarded as a remedy to the bereaved [14].

There is a need for a Digital Memorabilia of people's social presence in virtual space, which is a compendium of the digital artifacts of individual's online presence over a life span, showing the different facets of his personality and is live for an extended period, for the individual to be remembered for long on the online space by their friends.

In the following sections, we will focus on the previous works in the related area and the gap that is created by those research papers. The paper shows the technical aspect of the digital legacy: posterity rights and digital memorabilia by proposing a machine learning model using multiple logistic regression technique.

## II. LITERATURE REVIEW

Sudan A. et al., 2019 [1], in their review paper have explained different categories of digital assets, social media types as well as the concept of digital legacy. They have also explained different contexts of privacy rights which are concerned about people's digital legacy and what should be done to their assets after their demise. The digital posterity explains the passing of all the assets to their digital executers after the demise of the person.

According to the author, Cerrillo-i-Martínez, A., 2018 [2], digital footprint consists of three mechanisms: legal certainty, effectiveness and transparency. They must also respect the desires communicated by the user, their digital executers and provide enough certainty to allow a digital resources user a

never-ending rest in the online world. In the research paper of Peoples, C., & Hetherington, M., 2015 [3], they had created a survey to capture perceptions of users on digital cloud footprints. The results of which shows that users are generally not aware about their digital footprints and digital legacy. The survey includes people of every age group and of different places who came from a range of employment backgrounds.

This part of literature focuses on user interface frameworks and models designed for the digital legacy and its associated technologies. Byrd G., 2016 [4] in his study shows the high-level interaction between the digital legacy user interface, its users and the other online services such as cloud services. He made a functional design where users can add an account, amend an account, add a site, etc. Users can manage their account such as password and other information related to it. Users can also have an option to design their own digital memorial page where they can record their information. Whittaker, S., Bergman, O., & Clough, P., 2010 [5] in their paper have examined the effects of technologies related to digital photography which people had stored online for longer term. Due to poor organization of the digital contents, this study performed poorly. Another framework that ensures people to understand how to protect and pass on their digital legacy to their digital executers is given by Norris, J., & Taubert, M., 2016 [6]. The authors have made six steps framework that is associated with digital assets and digital end of life. It shows three categories which are digital assets, connected devices and digital legacy. In the field of digital legacy, another authors named Gulotta, R., Faste, H., & Forlizzi, J. (2012) [7] have created a tool called Revelado where users can store their information online so that their information can be accessed by their future generations and be remembered online forever. Kang, Y. S., & Lee, H., 2010 [8], brings out the author's attempt to propose a model to find out customer's satisfaction so as to design some investment strategy of retaining customers.

Some studies have highlighted the importance of public thoughts and reviews about digital legacy and posterity rights. Waagstein, A., 2014 [9] has collected data in the form of questionnaires mostly in semi-structured form. The questions were mostly related to digital legacy and digital artifacts. The authors concludes by discovering patterns and by making a summary of the interviews performed and in-depth readings were performed on some statements. The study by Gulotta, R., et al., 2013 [10] brings out the viewpoint of parents and focuses on finding the point of view of parents about the passing of digital materials in future. On the basis of their responses a system can be designed that can be used as provocative and speculative artifacts. The author had used diagrams and themes to interpret the findings. In the view of college students, Pempek, T. A. et al., 2009 [11], have highlighted experiences of college students of social networking on Facebook. They have proposed different factors such as frequency of Facebook by college students, gender, etc. They conducted surveys to find out the purpose of using Facebook by these college students. Another work which was done on this is by Massimi, M., & Baecker, R. M., 2010 [12] where the authors have presented the survey in the form of questionnaires to examine the use of technology and

other digital techniques to remember the deceased. Correa, T., et al., 2010 [13] shows the relationship between social media and personality predictors with respect to various factors such as gender, age, etc. The author has proposed various hypotheses in response to social media and personality predictors. These hypothesis are extraversion, emotional stability and openness. Both the personality predictors and social media showed how much these hypotheses have had an impact related to digital media. Petrelli, D., & Whittaker, S., 2010 [14] have conducted some fieldwork and compare the physical and digital work. They concluded the work with some digital limitations and design guidelines associated with it. Rubin, H. J., & Rubin, I. S., 2011 [15], in their book have conducted qualitative interviews to identify the gathering style of data. They describe detailed qualitative interviewing to underline philosophy related to project design and analysis.

This section of literature focused on the various factors related to social networking sites and social networking profiles. Lin, K. Y., & Lu, H. P., 2011 [16], the authors have focused on various factors that affects user's joining social networking sites by applying some network externalities and motivation theory. This is applied to find out why people are that desperate to join social networking sites. The factors involved here are age, gender, occupation, education and Facebook services. To find the reason behind the increase in usage of social media, the authors Lee, J., & Suh, E., 2013 [17] have used three theories to examine people's characteristics. These theories are Technology Acceptance Model (TAM), Innovation Diffusion Theory (IDT) and Network externality. Based on these theories, they find out some positive significant effects of SNS. In another research Sago, B., 2013 [18], has highlighted the various factors that influence adoption of social media and frequency. He examined the adoption factors for four platforms Facebook, Pinterest, Twitter, Google+ and the factors used are awareness, enjoyment, knowledge, reasons used, usefulness and ease of use. Kane GC et al., 2009 [19], in their study have stated the importance of social media in person's life. They highlighted how the social media platforms promote relationships. The authors have taken the example of health care industry to show the importance of social media platforms. Munish Sabharwal et al., 2012 [20] conducted a study with the objective to find out whether the selected Indian scheduled banks have presence on the Social Networking Media or not.

Few studies also focus on the life of famous media personalities after death, Sherlock, A., 2013 [21] has stated the reason and importance of conservation of famous personalities and the effect of social media on their careers. Even after their death, their followers will not go unwane due to digital technologies proposed by the author.

This section review studies highlighting the importance of digital artifacts in relation to digital legacy and the problems associated with it. Banks, R., 2011 [22] has highlighted the importance of managing the digital artifacts and also explans how to inherit those contents in the future. He wants to explore the technology that could help the people realize their potential. In another article, Banks, R., Kirk, D., &Sellen, A., 2012 [23] state the importance of artifacts in the life of people

as it can act as a trigger to remember someone after his demise. They highlighted on such artifact in this paper which is their heirloom. In this paper they suggested a design case study for the process of inheriting person's assets.

The study by the authors Romano, J. et al., 2011 [24] have focused on the life they have lived even after their death. They have thrown the light on the life of the person after their demise but online. They have pointed out different plans such as what could be done for the artifacts left behind by the person.

The next` two papers discuss the privacy rights of digital legacy. Edwards, L., & Harbina, E., 2013 [25], in their article have emphasized on the privacy rights of digital legacy of the deceased. They have given different defamation and moral rights for the regulation of post-mortem privacy. Gotved, S., 2014 [26] have offered a systematic way to keep track of people's timeline and their digital context related to physical death of the person. Bellamy, C. et al., 2013 [27] has pointed out the difficulties which are involved in conserving and leaving digital legacy online after the demise of a person. They pointed out several problems related to digital legacy, one of which is passing on digital music and books as it could lead to copyright issues. The next paper focuses on the sentiments and artifacts of the person, Kirk, D. S., & Sellen, A., 2010 [28] highlighted the sentiments and artifacts related to the person and the nature of thing. The authors explained the practices to keep sentimental artifacts of the person. Wiegand, D. L. M. et al., 2008 [29] in their article address issues related to dying research. They imposed some challenges related to informed consent, data collection, etc.

This part of literature describes the different memorialization practices and issues dealing with it. Walter, T. et al., 2012 [30] research is divided into two parts; the first part explains the practices related to dying and memorialization and the second part describes the concepts related to these practices. Odom, W. et al., 2010 [31], this paper describes the problems and issues about death and memorialization. The authors conducted in-depth interviews about the issues related to bereavement.

The result of the above literature is that most of the papers talked about the survey concerning the awareness of digital legacy of the people, whereas some of them talked about designing some digital memorial of people using their social networking profiles but none of them pointed out the technical aspect related to the digital legacy: digital posterity and digital memorabilia.

The studies by Munish, first [32], facilitated the researcher in overall preparation of literature review and planning for the overall research and the second [33], assisted in analysis.

### III. METHODOLOGY

#### A. Data Collection

To get people's opinion with regard to digital awareness, a questionnaire is made and has been distributed in the form of survey on the basis of different age groups, gender and different online platforms the respondents engage in.

The data used for this study were collected by forming the questionnaires related to different aspects of digital legacy such as the first section answers their personal questions, the second answers the matter related to digital legacy and the privacy rights related to digital legacy and the third answers the matter of the digital posterity and rights concerning their digital posterity. The questionnaires were distributed to people of different age groups and of different fields through Google forms. Based on the data, a model would be created for digital memorial of people based on their social networking profiles.

#### B. Analysis

The bar chart between posterity rights vs. age group is as follows in Fig. 1.

Below are the summarized responses in the form of pie charts and bar charts which we got from the questionnaires distributed through Google forms over the web which is shown in below figures from Fig. 2 to Fig. 13.



Fig. 1.    Bar chart between Age Group and Posterity Rights.



Fig. 2.    Response of People: Purpose of using IT Devices.



Fig. 3.    Response of People: Social Networking Applications.

Awareness of Digital Assets and Digital Legacy

Fig. 4.    Response of People: Awareness of Digital Assets and Digital Legacy.



Awareness of Privacy Rights of Digital Assets

Fig. 5.    Response of People: Awareness of Privacy Nights of Digital Assets.



Creation of First SNP

Fig. 6.    Response of People: Creation of First SNP.



Percentage of Respondents having SNP

Fig. 7.    Response of People: Percentage of Respondents having SNP.



Type of Device to Access SNP

Fig. 8.    Response of People: Type of Device to Access SNP.



Frequency of Accessing SNP

Fig. 9.    Response of People: Frequency of Accessing SNP.



Awareness of Right to Posterity

Fig. 10.  Response of People: Awareness of Right to Posterity.



Percentage of Respondents to Exercise Right to Posterity

Fig. 11.  Response of People: Percentage of Respondents to Exercise Right to Posterity.



Options to Exercise Right to Posterity

Fig. 12.  Response of People: Options to Exercise Right to Posterity.



Time Period of SNP

Fig. 13.  Response of People: Time Period of SNP.

## IV. RESULTS AND DISCUSSIONS

To assess the people's opinion on posterity rights, we had performed logistic regression analysis to predict a machine learning model from the data set we have got from Google forms.

First step was to identify dependent and independent variables. We will take dependent variable a categorical data in the form of YES or NO as if given an opportunity, would you like to exercise your rights to posterity (it is about announcing in advance what should be done with your SNPs or deciding the legal inheritor of your SNPs).

The independent variables were taken after analyzing the responses from respondents and factors such as age, frequency of using SNPs, awareness about digital assets and digital legacy, awareness about privacy rights concerning digital assets and awareness about rights to posterity from the given set of variables from the dataset were chosen as independent variables.

We applied multiple logistic regression to predict the relationship between posterity rights and various independent variables one of them being taken as age group.

Multiple Logistic Regression was applied on the collected dataset using Anaconda framework with Python with SciKit learn API.

The result of the logistic regression is given below in Fig. 14.

Below is the classification report of the model shown in Fig. 15.



Fig. 14. Result of Multiple Logistic Regression.



Fig. 15. Computation of other Parameters Such as Precision, Recall, F-Measure and Support based on Factors.

## V. CONCLUSION

The accuracy of the model is 81%. So, the model we have predicted has done fairly well. The study results indicated that our research model reveal good descriptive ability to predict user's persistent purpose whether to exercise their rights to posterity or not under various factors such as age, frequency of using SNPs, awareness about digital assets and digital legacy, awareness about privacy rights concerning digital assets and awareness about rights to posterity, giving a new way for researchers to inspect in future research work in related areas.

## VI. FUTURE SCOPE

The data set can be large so that more accurate model can be predicted in future and to get high accuracy, we can apply other technique other than regression.

Further research should endeavor to acquire more samples for more various SNS user type to validate our research model and to examine the differences among users. Moreover, we can add more factors or constructs such as self-efficacy, altruism etc. to give model a more precise view.

### REFERENCES

[1] Sudan A et.al (2019). Privacy rights for digital assets and digital legacy right for posterity: A Survey has presentated in 2nd International Conference on Data & Information Sciences on March 29-30, 2019 given at Raja Balwant Singh Engineering Technical Campus, Bichpuri, Agra, Uttar Pradesh, INDIA.

[2] Cerrillo-i-Martínez, A. (2018). How do we provide the digital footprint with eternal rest? Some criteria for legislation regulating digital wills. Computer Law & Security Review.

[3] Peoples, C., & Hetherington, M. (2015, November). The cloud afterlife: Managing your digital legacy. In Technology and Society (ISTAS), 2015 IEEE International Symposium on (pp. 1-7). IEEE.

[4] Byrd, G. (2016). Immortal Bits: Managing Our Digital Legacies. Computer, (3), 100-103.

[5] Whittaker, S., Bergman, O., & Clough, P. (2010). Easy on that trigger dad: a study of long term family photo retrieval. Personal and Ubiquitous Computing, 14(1), 31-43.

[6] Norris, J., & Taubert, M. (2016). P-221 Working with hospices to ensure patients' digital legacy wishes are adhered to.

[7] Gulotta, R., Faste, H., & Forlizzi, J. (2012). Revelado: Exploring the Preservation of our Digital Data.

[8] Kang, Y. S., & Lee, H. (2010). Understanding the role of an IT artefact in online service continuance: An extended perspective of user satisfaction. Computers in Human Behavior, 26(3), 353-364.

[9] Waagstein, A. (2014). An exploratory study of digital legacy among death aware people. Thanatos, 3(1), 46-67.

[10] Gulotta, R., Odom, W., Forlizzi, J., & Faste, H. (2013, April). Digital artifacts as legacy: exploring the lifespan and value of digital data. In Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (pp. 1813-1822). ACM.

[11] Pempek, T. A., Yermolayeva, Y. A., & Calvert, S. L. (2009). College students' social networking experiences on Facebook. Journal of applied developmental psychology, 30(3), 227-238.

[12] Massimi, M., &Baecker, R. M. (2010, April). A death in the family: opportunities for designing technologies for the bereaved. In Proceedings of the SIGCHI conference on Human Factors in computing systems (pp. 1821-1830). ACM.

[13] Correa, T., Hinsley, A. W., & De Zuniga, H. G. (2010). Who interacts on the Web? The intersection of users' personality and social media use. Computers in Human Behavior, 26(2), 247-253.

[14] Petrelli, D., & Whittaker, S. (2010). Family memories in the home: contrasting physical and digital mementos. Personal and Ubiquitous Computing, 14(2), 153- 169.

[15] Rubin, H. J., & Rubin, I. S. (2011). Qualitative interviewing: The art of hearing data. Sage.

[16] Lin, K. Y., & Lu, H. P. (2011). Why people use social networking sites: An empirical study integrating network externalities and motivation theory. Computers in human behavior, 27(3), 1152-1161.

[17] Lee, J., & Suh, E. (2013). An Empirical Study of the Factors Influencing Use of Social Network Service. In PACIS (p. 181).Munish Sabharwal et. al, "Indian Banks: Presence and Interactivity level on Social Networking Media", IFRSA Business Review, ISSN (Online): 2249-5444 ISSN (Print): 2249-8168 Impact factor (2012): 0.1351, Vol. 2 Issue 4, pp.360-365, Dec 2012.

[18] Sago, B. (2013). Factors influencing social media adoption and frequency of use: An examination of Facebook, Twitter, Pinterest and Google+. International Journal of Business and Commerce,3(1),1-14.

[19] Kane, G. C., Fichman, R. G., Gallaugher, J., & Glaser, J. (2009). Community relations 2.0. Harvard business review, 87(11), 45-50.

[20] Munish Sabharwal et. al, "Indian Banks: Presence and Interactivity level on Social Networking Media", IFRSA Business Review, ISSN (Online): 2249-5444 ISSN (Print): 2249-8168 Impact factor (2012): 0.1351, Vol. 2 Issue 4, pp.360- 365, Dec 2012.

[21] Sherlock, A. (2013). Larger than life: Digital resurrection and the re-enchantment of society. The Information Society, 29(3), 164-176.

[22] Banks, R. (2011). The future of looking back (Microsoft Research). Microsoft Press.

[23] Banks, R., Kirk, D., &Sellen, A. (2012). A design perspective on three technology heirlooms. Human–Computer Interaction, 27(1-2), 63-91.

[24] Romano, J. (2011). Your Digital Afterlife: When Facebook, Flickr and Twitter Are Your Estate, What's Your Legacy? New Riders.

[25] Edwards, L., &Harbina, E. (2013). Protecting post-mortem privacy: Reconsidering the privacy interests of the deceased in a digital world. Cardozo Arts & Ent. LJ, 32, 83.

[26] Gotved, S. (2014). Research Review: Death Online-Alive and Kicking!. Thanatos, 3(1/2014).

[27] Bellamy, C., Arnold, M., Gibbs, M., Nansen, B., & Kohn, T. (2013). Life beyond the timeline: creating and curating a digital legacy. In meeting of Prato Community Informatics Research Network (CIRN), Prato, Italy.

[28] Kirk, D. S., &Sellen, A. (2010). On human remains: Values and practice in the home archiving of cherished objects. ACM Transactions on Computer-Human Interaction (TOCHI), 17(3), 10.

[29] Wiegand, D. L. M., Norton, S. A., & Baggs, J. G. (2008). Challenges in conducting end- of-life research in critical care. AACN Advanced Critical Care, 19(2), 170-177.

[30] Walter, T., Hourizi, R., Moncur, W., &Pitsillides, S. (2012). Does the internet change how we die and mourn? Overview and analysis. OMEGA-Journal of Death and Dying, 64(4), 275-302.

[31] Odom, W., Harper, R., Sellen, A., Kirk, D., & Banks, R. (2010, April). Passing on & putting to rest: understanding bereavement in the context of interactive technologies. In Proceedings of the SIGCHI conference on Human Factors in computing systems (pp. 1831-1840). ACM.

[32] Sabharwal, M. (2016, March). Contemporary research: intricacies and aiding software tools based on expected characteristics. In 10th National Research Conference on Integrating Technology in Management Education, AIMA, New Delhi (pp. 28-29).

[33] Sabharwal, M. (2018). The Use of Soft Computing Technique of Decision Tree in Selection of Appropriate Statistical Test For Hypothesis Testing. In Soft Computing: Theories and Applications (pp. 161-169). Springer, Singapore.

# Smart Sustainable Agriculture (SSA) Solution Underpinned by Internet of Things (IoT) and Artificial Intelligence (AI)

Eissa Alreshidi[1]

Assistant Professor at College of Computer Science & Engineering, University of Hail, Hail, P.O. Box 2440, Saudi Arabia

*Abstract*—The Internet of Things (IoT) and Artificial Intelligence (AI) have been employed in agriculture over a long period of time, alongside other advanced computing technologies. However, increased attention is currently being paid to the use of such smart technologies. Agriculture has provided an important source of food for human beings over many thousands of years, including the development of appropriate farming methods for different types of crops. The emergence of new advanced IoT technologies has the potential to monitor the agricultural environment to ensure high-quality products. However, there remains a lack of research and development in relation to Smart Sustainable Agriculture (SSA), accompanied by complex obstacles arising from the fragmentation of agricultural processes, i.e. the control and operation of IoT/AI machines; data sharing and management; interoperability; and large amounts of data analysis and storage. This study firstly, explores existing IoT/AI technologies adopted for SSA and secondly, identifies IoT/AI technical architecture capable of underpinning the development of SSA platforms. As well as contributing to the current body of knowledge, this research reviews research and development within SSA and provides an IoT/AI architecture to establish a Smart, Sustainable Agriculture platform as a solution.

*Keywords—Smart Agriculture; Internet of Things; IoT; Artificial Intelligence; AI; Fragmentation; Smart Sustainable Agriculture solutions*

## I. INTRODUCTION

Agriculture forms a critical activity vital to the survival of humanity for approximately many thousands of years [1].This relationship has resulted in the advancement of agricultural activities, initially through the time-consuming methods of traditional agriculture [2]. The current recent rapid increase in in the global population (predicted to rise to 8.9 billion by 2050) has now led to an urgent need to balance demand and supply through the use of new technologies [3] to increase food production [4, 5]. This development places pressure on natural resources, with agriculture now consuming 70% of the world's fresh water supply for the purposes of irrigation. Limited resources and the impact of climate change will therefore lead to considerable challenges in producing sufficient high quality food to support the population [6]. Smart Agricultural is a global initiative to preserve resources and maintain sustainable agriculture [7]. Recently, researchers have adopted the Internet of Things (IoT) [8, 9], with a number of studies emphasizing the adoption and implementation of IoT in agriculture, farming, and irrigation [10]. Around the globe, many private companies and organizations are now focusing on investigating new technologies to create a smarter agriculture environment. These include mechanical and economic aspects, engineering, food retailers and computing. However, agricultural processes are fragmented, resulting in a number of issues, i.e. difficulties in operating and managing smart machines, data sharing and management, data analysis and storage [11, 12]. It is therefore important to facilitate cooperation when developing standards for smart agriculture, while also enhancing interoperability among different stakeholders, systems and technologies [13].

The use of IoT and AI technologies has the potential to result in a positive transformation of traditional agriculture [3], including: (a) improved use of data collected from smart agriculture sensors; (b) managing and governing the internal processes within the smart agriculture environment (including the management of the harvesting and storage of crops); (c) waste reduction and cost saving; (d) increasing business efficiency by means of automating traditional processes; and (e) improving the quality and volume of products [14]. A major challenge is to provide farmers with the required information in a rapid manner [15]. AI therefore has significant potential to address the urgent challenges faced by traditional agriculture. There has, over previous decades, been considerable research and application of AI, including in: (a) smart agriculture; (b) robotics; (c) agricultural optimization management; (d) automation; (e) agricultural expert systems; (f) agricultural knowledge-based systems; and (g) decision support systems [16].

There remains a lack of research and development in relation to Smart Sustainable Agriculture (SSA), accompanied by complex obstacles arising from the fragmentation of agricultural processes, i.e. the control and operation of IoT/AI machines; data sharing and management; interoperability; and large amounts of generated data analysis and storage. Therefore, this study firstly, explores existing IoT/AI technologies adopted for SSA and secondly, establishes an IoT/AI technical architecture to underpin SSA platforms, in order to tackle fragmentation in traditional agriculture processes and enrich the research and development of future smart agriculture worldwide via establishment of a Smart, Sustainable Agriculture platform as a solution.

There now follows an outline of the methodology underpinning this research, supported by related work highlighting the history of smart agriculture, smart and advanced computing technologies and examples of IoT/AI technologies in current agricultural practices. This is supported

by an in-depth discussion of Smart Agriculture and IoT/AI SSA technical architecture, along with the most significant outcomes from this study. The paper finishes with concluding remarks and plans for future work.

## II. BACKGROUND

Due to the lack of literature concerning the development of IoT frameworks for SSA, this study focuses on: (a) the history of smart agriculture, its potential and challenges; (b) smart and advanced computing technologies; and (c) existing smart, sustainable agricultural frameworks.

### A. Smart Sustainable Agriculture

There has recently been considerable research into SSA, employing various different terms, including Precision Farming, Smart Irrigation and Smart Greenhouse. This paper commences with an examination of these concepts to determine the definition of SSA used in this study.

Precision Farming refers to a method of managing farms and conserving resources through the use of IoT and Information and Communication Technologies (ICT). It obtains real-time data concerning the condition of farm elements, (i.e. crops, soil and air) to protect the environment while ensuring profits and sustainability [15]. Smart Irrigation is a method of improving the efficiency of irrigation processes and reducing water losses, while conserving existing water resources using IoT-based smart irrigation systems [16]. Drones are employed in many agricultural applications, including monitoring field crops and livestock, and scanning large areas, while sensors on the ground collect a huge range of information [13]. Smart greenhouses promote the cultivation of crops with the least degree of human intervention possible, through use of continuously monitored climatic conditions (i.e. humidity, temperature, luminosity and soil moisture), triggering automated actions based on the evaluated changes and implementing corrective action to maintain the most beneficial conditions for growth [17].

Farm Management Systems (FMS) can assist farmers with a variety of collected information, by managing and controlling various tracking devices and sensors. The collected information is analysed for the undertaking of complex decision-making tasks before being placed in a storage medium. This enables the use of the most effective agricultural data analysis practices [18]. Soil Monitoring Systems help to track and improve the quality of soil through the monitoring of its physical, chemical, and biological properties. Livestock monitoring systems provide real-time assessment of the productivity, health and welfare of livestock, to promote the health of animals [19]. The IoT/AI SSA platform Cloud offers real-time information to farmers to facilitate decision-making and reduce operational costs, while at the same time enhancing productivity. Following a review of a considerable amount of research, we define SSA as the utilization process of IoT/AI technologies to establish, monitor, manage, process and analyse data generated from various agricultural resources, such as field, crops, livestock and others to ensure the sustainability and quality of agricultural products and further enrich decision-making taken by stakeholders.

### B. Smart and Advanced Computing Technologies

This section provides an overview of appropriate technologies underpinning the development of smart, sustainable agriculture platforms, including: IoT; Big Data Analytics (BDA); Cloud Computing (CC); Mobile Computing (MC); and Artificial Intelligence (AI):

*1) Internet of Things (IoT):* IoT is a technology aimed at connecting all intelligent objects within a single network, i.e. the Internet. It involves all kinds of computer technologies, both (a) hardware (i.e. intelligent boards and sensors) and (b) software (i.e. advanced operating systems and AI algorithms). Its primary target is the establishment of applications for devices, in order to enable the monitoring and control of a specific domain. It has been widely adopted in many areas, i.e. industrial business processes; home machines; health applications; and smart homes and cities. IoT connectivity encompasses people, machines, tools and locations, aiming to achieve different intelligent functions from data sharing and information exchange [17]. However, it is primarily used in agriculture for management of agricultural products within gathered real-time data, alongside: (1) searching; (2) tracking; (3) monitoring; (4) control; (5) managing; (6) evaluating; and (7) operations within a supply chain [1, 9].

*2) Big Data Analytics (BDA):* BDA refers to the large volume of data gathered from different datasets sources over a long period of time, i.e. sensor, Internet and business data. The datasets used in this technology surpass the computational and analytical capabilities of typical software applications and standard database infrastructure. Its primary task is to capture, store, analyze and search for data, as well as seeking to identify concealed patterns in the gathered data. Thus, BDA involved the utilization of: (a) tools, (i.e. classification and clustering); (b) techniques, (i.e. data mining, machine learning and statistical analysis); and (c) technologies (i.e. Hadoop and spark). These go beyond traditional data analytical approaches, being employed to extract beneficial knowledge from a considerable amount of data, in order to facilitate timely and accurate decision–making [17]. However, the use of BDA in agriculture focusses on management of the supply chain of agricultural products, in order to enhance decision-making and minimize the cost of production cost. It is also employed for the analysis of the properties of different types of soil for classification and further enhancement. Furthermore, it is useful for the improved prediction and production of crops.

*3) Cloud Computing (CC):* CC has is a recent and rapidly growing phenomenon within IT [18]. The Cloud is not restricted to a particular business domain, but has been implemented to underpin and support various software applications and platforms [19]. It offers easy access to the Cloud provider's high-performance and storage infrastructure over the Internet, with one of its main benefits being to conceal from users the complexity of IT infrastructure management [20] [21]. NIST [22] defined CC as "a model for enabling convenient, on-demand network access to a shared pool of

configurable computing resources (e.g. networks, servers, storage, applications and services) that can be rapidly provisioned and released with minimal management effort or service provider interaction". The Cloud can be seen as high virtualization method for datacenter infrastructure distributed over a wide geographical area, linked by means of high bandwidth network cables providing a variety of virtualized services. These include entire infrastructures, as well as small software applications and different types of services, i.e. high-performance computing and large scalable storage services based on a pay-per-use model. CC can be divided into four main layers: (1) hardware; (2) infrastructure; (3) platform; and (4) application [23]. The delivery of Cloud services can generally be divided into three different models: (1) Infrastructure-as-a-Service (IaaS); (2) Platform-as-a-Service (PaaS); and (3) Software-as-a-Service (SaaS) [24, 25]. CC is considered the most effective method of storing agricultural data, along with IoT [1].

*4) Mobile Computing (MC):* MC refers to infrastructure in which data processing and data storage take place externally to the mobile device [26]. MC applications transfer computing power, processing and data storage from mobile devices in the Cloud [27, 28]. MC has had a considerable impact on modern daily life, due to the availability and low cost of purchasing and communication. It is now widely used in every field, including the agricultural sector [29], in which MC systems collect and send daily data to farmers, informing them of both the production status and weather conditions [29]. It is crucial to use automatic Radio Frequency Identification (RFID) efficient traceability systems to store and access data on electronic data chips in a more rapid and accurate manner. It has been primarily applied to the logistics of industrial products, for the purposes of identification and to check delivery processes [30].

*5) Artificial Intelligence (AI):* AI has been employed in smart systems over a long period of time[31], being the science of creating intelligent machines to facilitate everyday life [32]. AI covers many areas, including computer vision, data mining, deep learning, image processing and neural networks [16, 33]. AI technologies are now emerging to assist and improve efficiency and tackle many of the challenges facing the agricultural industry, including soil health, crop yield and herbicide-resistance. According to Sennaar [34], agricultural AI Cloud applications fall into three main categories, as discussed below.

*a) Robots:* these are developed and programmed to handle fundamental agricultural tasks (i.e. harvesting crops) more rapidly and with a higher capacity than human workers. Examples of robotic applications include: (a) See and Spray (i.e. a weed control robot) and (b) Harvest CROO (i.e. a crop harvesting robot). Agricultural robots have the potential to become valuable AI applications, i.e. milking robots.

*b) Monitoring Crop and Soil:* this employs computer vision and deep-learning algorithms for processing captured data by sensors monitoring crop and soil health, i.e. the PEAT machine for diagnosing pests and soil defects, based on deep

learning application known as Plantix that identifies potential defects and nutrient deficiencies in the soil. A further example is Trace Genomics, a machine learning based service for diagnosing soil defects and providing soil analysis services to farmers. This uses machine learning to provide farmers with a sense of both the strengths and weaknesses of their soil, with the emphasis being on the prevention of poor crops and optimizing the potential for healthy crop production. A SkySquirrel technology is an example of the use of drones and computer vision for crop analysis.

*c) Predictive Analytics:* This analysis captured data, based on machine learning models capable of tracking and predicting various environmental impacts on crop harvest, i.e. changes in weather. Examples of such AI technologies include (a) aWhere (i.e. prediction of weather and crop sustainability) and (b) Farmshots (i.e. monitoring of crop health and sustainability). Crop and soil monitoring technologies are important applications for addressing issues related to climate change. IoT/AI technologies (such as drone and satellite) that generate a large amount of data on a daily basis have the potential to enable agricultural production to forecast changes and detect opportunities. It is predicted that, over the coming years, IoT and AI applications will attract a considerable degree of interest from large industrial agricultural enterprises [34].

The benefits and advantages of the agricultural use of IoT are as follows: (a) efficiency of input; (b) cost reduction; (c) profitability; (d) sustainability; (e) food safety and environmental protection [35]. However, Ferrández-Pastor, et al. [36] considered SSA to contain a number of barriers potentially hindering its adoption: (a) initial expectations and advantages remaining unfulfilled; (b) complexity of technology and incompatibility of components; (c) a lack of products; and (d) the high cost involved in the establishment and maintenance of such facilities. To ensure the adoption and improvement of smart technologies in the agriculture sector, it is vital for farmers to be trained and given up-to-date knowledge of IoT/AI technologies. Furthermore, it is crucial to test and validate IoT/AI applications, due to the high risk involved in the adoption of these technologies in a critical sector, along with the influence of environmental factors.

*C. Examples of IoT/AI Technologies in Current Agriculture Practices*

There are many types of IoT and AI sensors and applications in current agricultural studies and development. Table 1 provides an overview of the most commonly employed IoT/AI platforms/technologies found in smart agriculture.

*D. Examples of an Existing AI/IoT Research in Smart Agriculture*

There are a number of specific challenges that need to be considered before investing in smart agriculture, primarily those falling into the following categories: (1) hardware; (2) data analysis; (3) maintenance; (4) mobility; and (5) infrastructure [56]. Nonetheless, there are many research efforts in the field of IoT/AI to support the creation and establishment of SSA, as shown in Table 2.

TABLE I.        EXAMPLE OF IOT/AI APPLICATIONS IN SMART AGRICULTURE

| Category | Tool/Company | Description |
|---|---|---|
| Climate conditions Monitoring | allMETEO [37] | A portal to manage IoT micro weather stations, to gather real-time data access and create a weather map. It also provides an API for easy real-time data transfer into developed or existing infrastructure. |
| | Smart Elements [38] | A collection of products that improve efficiency by eliminating manual checking. They work by deploying a wide range of sensors generating a report back to an online dashboard, allowing rapid and informed decisions based on real-time conditions. |
| | Pycno [39] | A software and sensor allowing continuous data collection and flow from the farm to smartphone. It also contains a dashboard to apply the latest phenological and disease models to monitor trends and assess risk to agricultural products. |
| Greenhouse automation | Farmapp [40] | A process of monitoring pests and diseases, generating reports for mobile applications. It records the data quickly and more efficiently than traditional methods (i.e. paper), allowing a smooth implementation. The stored data is synchronized with the server, enabling the following metrics to be immediately observed: (1) a satellite map with recorded points; (2) the current sanitary status of the farm; (3) comparative heatmaps to easily compare previous measures with the current situation; and (4) charts and reports concerning pests and diseases. |
| | Growlink [41] | A platform that tightly integrates hardware and software products, enabling smarter working, including providing wireless automation and control, data collection, optimization, and monitoring and visualization. |
| | GreenIQ [42] | A system to control irrigation and lighting from all locations and to connect IoT devices to automation platforms. |
| Crop management | Arable [43] | A device that combines weather and plant measurements, sending data to the Cloud for instant retrieval from all locations. It offers continuous indicators of stress, pests and disease. |
| | Semios [44] | A platform focused on yield improvement. It enables farmers to assess and respond to insects, disease and the health of crops using real-time data, forming on-site sensing, big data and predictive analytics solutions for sustained agricultural products. |
| livestock monitoring and management | SCR/Allflex [45] | An advanced animal monitoring system, aimed at the collection and analysis of critical data, including for individual animals. It delivers, when needed, the heat, health and nutrition insights required by farmers for effective decision making. |
| | Cowlar [46] | A smart neck collar for monitoring dairy animals to gather information on temperature, rumination, activity and other behavior. The intelligence algorithm in the system allows for the detection of health disorders before the appearance of visual symptoms. It can monitor body movement patterns and gait to provide accurate oestrus detection alerts. It uses a solar power base unit, along with a waterproof and non-invasive monitoring system, both comfortable for the animal and requiring minimum maintenance. |
| End-to-end farm management systems | FarmLogs [47] | This system monitors field conditions, facilitating the planning and managing of crop production. It also markets agricultural products. |
| | Cropio [48] | A decision-making tool used to optimize fertilization and irrigation to control the amount of fertilizer and reduce the use of water. It combines weather information and satellite data to monitor crops and field forecasts. |
| Predictive Analytics | Farmshots [49] | A system analyzing satellite and drone images of farms fields to map potential sign of diseases, pests and poor nutrition. It turns images into a prescription map to optimize farm production and view analytics on farm performance. Generated data in the Cloud can be exportable into nearly all agricultural software for prescription creation. |
| | aWhere [50] | A platform employed for weather prediction and information on crop sustainability. Its goal is to deliver complete information and insight for real-time agricultural decisions on a daily basis and at a global level. |
| Crop and Soil Health Monitoring | Plantix [51] | A machine learning based tool to control and manage the agriculture process, disease control, and the cultivation of high-quality crops. |
| | Trace Genomics [52] | A soil monitoring system performing complex tests (i.e. DNA) on soil samples. It uses a machine learning process known as 'genome sequencing' that generates a health report for a soil sample by reading its DNA and comparing it to a large soil DNA database. |
| Agriculture machines /drones | SkySquirrel [53] | A drone system aimed at helping users to improve their crop yield and reduce costs. Users pre-program a drone's route, and, once deployed, the device will leverage computer vision to record images to be used for analysis. Once the drone completes its route, users can transfer the data to a computer and upload it to a Cloud drive. It uses algorithms to integrate and analyze the captured images and data to provide a detailed report on the health and condition of crops. |
| | See & Spray [54] | A robot designed to control weeds and protect crops. It leverages computer vision to monitor and precisely spray weeds and infected plants. |
| | CROO [55] | A robot that assists in the picking and packing of crops. The manufacturer claimed that this robot can harvest eight acres in a single day and replace the work of thirty human laborers. |

TABLE II.    IoT/AI Research and Development in Smart Sustainable Agriculture

| Researcher/s | Year | Summary |
|---|---|---|
| Ray [57] | 2017 | The researchers undertook a review of various potential IoT applications, including the specific issues and challenges associated with IoT deployment to improve farming. They comprehensively analyzed the specific requirements the devices and wireless communication technologies associated with agricultural IoT applications. They presented different case studies to explore existing IoT based solutions operated by various organizations and individuals, followed by categorizing them based on their deployment parameters. Furthermore, they identified a number of factors for the improvement and future road map of work using IoT. |
| Mekala and Viswanathan [58] | 2017 | The researchers surveyed a number of conventional applications of Agricultural IoT Sensor Monitoring Network technologies utilizing CC. Their study aimed at understanding the diverse technologies to build smart, sustainable agriculture. They addressed a simple IoT agriculture model with a wireless network. |
| Kamilaris, et al. [59] | 2017 | The researchers reviewed work in agriculture employing the practice of big data analysis to solve various different problems. Their review emphasized the opportunities provided by big data analysis for the development of smarter agriculture, the availability of hardware and software, as well as the techniques and methods for big data analysis. |
| Rajeswari, et al. [29] | 2017 | The researchers investigated a number of different features, i.e. humidity, temperature sensing, server-based remote monitoring system detection and soil moisture sensing. They used sensor networks to measure temperature, moisture and humidity in place of manual checking. They deployed several sensors in different locations within farms, using a single controller. Their major objective was to collect real-time data of the agriculture production environment to establish an easy access agricultural advice, in order to identify weather or crops patterns. |
| Antonacci, et al. [5] | 2018 | The researchers attempted to provide nanotechnology-based (bio) sensors to support farmers in delivering an analysis that is accurate, fast, cost-effective, and useful in the field to identify water and soil nutrients/pesticides, soil humidity, and plant pathogens. |
| Cadavid, et al. [60] | 2018 | The researchers proposed an extension to a popular open-source IoT platform, known as 'Thingsboard'. This formed the core of a Cloud-based Smart Farming platform and deliberate sensors, a decision support system, and a configuration of remotely autonomous and controlled machines (e.g. water dispensers, rovers or drones). |
| Soto-Romero, et al. [61] | 2019 | The researchers designed an easily insertable cylindrical sensor with internal electronics to offer a low power electronic architecture to measure and communicate wirelessly with a LoRa, Sigfox network or mobile phones. |
| Nóbrega, et al. [62] | 2019 | The researchers reviewed the proposed stack and details of the recent developments within smart agriculture, focusing on IoT/Machine-2-Machine interaction. They described the design and deployment of a gateway addressing the requirements of the SheepIT service, evaluating this gateway using real scenarios in terms of performance, thus demonstrating its feasibility and scalability. |

## III. Methodology

As well as undertaking the literature review, the current researcher enhanced this study by informally interviewing experienced farmers. The study aims to establish an IoT/AI SSA architecture, as well as exploring the potential of the use of IoT and AI as a backbone to establish an SSA platform. A review was employed to identify, analyze and study key books, journals, reports, and white papers, in order to achieve the above-noted aim. The lack of existing studies in this area ensures that this current research also contributes to the body of knowledge by establishing an IoT/AI framework for the adoption of smart technologies, in order to establish smart sustainable agricultural practices. Fig. 1 shows the methodology adopted for this study.

## IV. Results and Discussion

Based on the research aim outlined in the Methodology, this section is divided into (A) Domains of Smart Sustainable agricultural model; (B) B. Proposed IOT/AI SSA platform as a solution; and (C) Proposed IoT/AI technical architecture for SSA platform.

### A. Domains of Smart Sustainable Agricultural Model

The results from the literature review revealed that several domains need to be considered when adopting the smart agricultural model. Fig. 2 demonstrates the interrelation and complexity of data flow between different Smart, Sustainable Agriculture domains.



Fig. 1.   IoT/SSA Research Methodology.

Fig. 2.    Chaotic Data Flow and Interrelation of SSA Domains.

These domains are discussed individually, as follows:

- **Human resources:** This refers to people, policies and practices within the agricultural environment, which are as important as in any other domain, as are weather and technology. Human resources receive careful attention, due to their significant impact on production, as well as financial and marketing decisions. Whatever its size, an agricultural concern requires effective human resources management and planning, including hiring and keeping employees who are engaged, high-performing and effective communicators. Providing up-to-date knowledge potentially opens the means to adopt smart technologies in an agricultural environment.

- **Crops**: This refers to a plant that can be grown and harvested extensively for subsistence or profit: (1) food crops (i.e. for human consumption); (2) feed crops (i.e. for livestock consumption); (3) fibre crops (i.e. for cordage and textiles); and (4) oil crops (i.e. for consumption or industrial uses).

- **Weather:** This plays a major role in determining the success of agricultural processes. Most field crops and livestock are solely dependent on climactic conditions to provide life-sustaining water and energy. Adverse weather can cause losses in agricultural products, particularly during critical stages of growth. The elements of weather (solar radiation, temperature, precipitation, humidity and wind) influence the physiology and production of agricultural plants and animals. Severe weather (i.e. tornadoes, drought, flooding, hail and wind storms) can cause considerable damage and destruction to fields and livestock.

- **Soil:** This forms a critical aspect of successful agriculture, being the source of nutrients used to grow crops, which are subsequently passed into plants and then to humans and animals. Healthy soils produce healthy and rich food supplies; however the health of soil tends to decline over time, forcing farmers to move to new fields. Soil health depends on regional conditions and climates, with soil nutrients more likely to deteriorate in dry climates, particularly if irrigated, which, if not managed carefully, can result in salinization, i.e. a build-up in the level of salts and chemicals contained in water. Healthy and rich soil can be achieved through the use of IoT sensors to monitor its chemical status, using specific sensors (e.g. moisture sensors), whose data readings are transferred to the data management and analysis layer for analysis, assisting decisions concerning the need for fertilizer.

- **Pests:** These consist of any living creature that is invasive, or damaging to crops, livestock or human structures. Pests often occur in high quantities, to the detriment of agricultural products. It is vital to control and monitor these creatures by means of IoT/AI technologies, to avoid serious diseases, including plague and malaria, as well as plant and livestock diseases.

- **Fertilization:** Soils naturally contain many nutrients, i.e. nitrogen, phosphorous, calcium and potassium. Crops are unable to function effectively and produce high quality food when their nutrient level is low. The natural levels of nutrients in the soil need to be enhanced by the addition of nutrients once crops have been harvested. Fertilizers have been used since the beginning of agriculture, but it is now recognized that their extensive use can, if not correctly controlled, harm the environment. Therefore, farmers use IoT sensors to read and test soil samples for baseline testing to enable them to add fertilizers using correct and appropriate measurements. Fertilization is an important method of maintaining sustainable agricultural production systems.

- **Agricultural Products:** These are derived from cultivating crops or livestock to sustain or enhance human life. Human beings also use a wide collection of agricultural products on a daily basis, i.e. food and clothing. Agricultural products fall into the following groups: (a) grains; (b) foods; (c) fuels; (d) fibres; (e) livestock; and (f) raw materials. Food is the most extensively produced agricultural product.

- **Irrigation/Water:** Water demand in agriculture is now rising globally and particularly in Mediterranean countries, increasing the pressure to preserve available freshwater resources. Smart, sustainable agriculture processes should therefore focus on new and efficient techniques to improve agricultural productivity, which promote considerable savings in terms of food consumption and wasted water.

- **Livestock:** These are animals raised in a domesticated agricultural environment, for the purposes of labour and to produce commodities such as eggs, meat, milk, fur, wool and leather. Animal husbandry is a component of current agriculture and refers to the breeding, maintenance and slaughter of livestock.

- **Machines:** Agricultural equipment is any kind of machinery used to assist with farming. Such machines

can be light or heavy, i.e. tractors. Modern farm machinery is seen as the important driver for increased agricultural sustainability, efficiency and competitiveness. Smart technologies can reduce the impact of farming practices within global agriculture. The current development of agricultural machinery addresses environmental challenges, while increasing productivity and bringing economic benefits. These smart agricultural machines should be: (a) fast, accurate, versatile and intelligent; (b) produce less $CO^2$ emissions and (c) make use of bioenergy.

- *Fields:* This refers to an area of land used for agricultural purposes, i.e. crops, cultivation or for livestock. Many fields have borders composed of a strip of bushes used to provide both food and cover, in order to ensure the survival of wildlife. Monitoring field activities using IoT devices can have a significantly positive impact on controlling objects within the field.

### B. Proposed IOT/AI SSA Platform as a Solution

This platform would prove a valuable medium to facilitate data flow and sharing among SSA domains. Many researchers have developed IoT architectures, but their efforts have tended to target specific areas of IoT/AI, i.e. sensors or weather monitoring systems [58, 60, 61]. This current paper is proposing a holistic IoT/AI platform to cover all areas within a SSA environment, performing the following tasks: (a) manage and govern data flow between SSA domains; (b) facilitate the integration of the different components of SSA architecture; (c) tackle interoperability issues caused by the utilization of different tools and software; (d) provide easy-to-use interfaces for interaction; (e) provide an ability to generate reports based on real-time data and keep it updated; (f) store generated data in sustainable storage place (i.e. the Cloud), to enable it to permanently recorded for future reuse; (g) isolate different layers to improve the development process in the future; and (h) the platform should consist of several nodes, so that, in the case of any failure, other nodes can keep the system up and running. Fig. 3 shows how the SSA-IoT/AI platform would be used at the center of the SSA domain to facilitate business process and data flow and to share within a smart, sustainable agricultural environment.



Fig. 3. Proposed SSA-IoT/AI Platform.

### C. Proposed IoT/AI Technical Architecture for SSA Platform

Fig. 4 shows the overall AI/IoT technical architecture for SSA. It consists of two main components: (a) the first component: SSA Layers and AI/IoT technologies; and (b) the second component: data lifecycle within SSA architecture and data process location. In order to provide additional detail concerning the framework, the following description gives further details about each component:

*1)* The First component: Smart Sustainable Agriculture (SSA) layers and (SSA) AI/IoT technologies

The first component of smart, sustainable agriculture (SSA), AI/IOT framework, consists of the following layers: (a) Physical Hardware and Storage layer; (b) AI and Data Management layer, and Governance layer; (c) Network layer; (d) Security layer; (e) Application layer; (f) IoT and sensing layer; and (g) SSA domain layer. Fig. 4 demonstrates the following description of each layer, highlighting its primary role:

- *Physical hardware and Storage layer*: This layer consists of powerful hardware to host virtualized machines, as well as dedicated traditional storage medium, cloud storage solutions or hybrid storage solutions. It contains the supporting hardware for IoT devices that exist in Internet of Thins and Sensing layer.

- *Artificial Intelligence and Data Management Layer*: This layer is responsible for managing processes and controlling the business logic, focussing on three main activities: (i) data analysis and processing, using data mining and intelligence statistical analysis on generated data; (ii) data classification and transformation, using ontologies, and machine learning to classify and transform analysed data; and (iii) data interpretation, representing the transformed data into knowledge to make machines smarter.

- *Network Layer:* This layer contains all the different types of network connections exist, with this cloud consisting of wired or wireless network connections/devices. It is important that this layer uses the latest networking technologies to keep up with the most recent developments in other sectors. Its main job is to facilitate the transaction of all data to and from different layers within the architecture. The technologies of IoT in this layer include Internet, WIFI and GSM/CMDA. It is responsible for data accessibility and availability throughout other layers. Further, it manages data and its flow within between all layers.

- *Security Layer:* This layer is responsible for data security transferred among different layers and should be the means of addressing any security concerns and vulnerabilities within all other IoT/SSA layers, i.e. malware, spyware and viruses. It should employ up-to-date security solutions and AI algorithms to block and quarantine any threat to the platform.

- *Applications Layers:* This layer gathers different applications related to smart sustainable agriculture. It is built based on AI and data management layers. Many

smart, sustainable agricultural applications could be developed and integrated into this layer, i.e. crop monitoring applications and drone-controlled applications. This layer focuses on the supervision aspect of the data flow and migration between all layers and can provide an authorized institution with full or partial governance on data migration, transactions and access.

- *Internet of Things (IoT) and Sensing Layer:* This forms the first interaction layer with SSA domains. It uses and hosts various types of IoT devices (e.g. sensors), capable of collecting data from real-world objects, sharing it to provide real-time data. Many sensors in the Cloud are hosted and integrated within this layer, i.e. humidity sensors, moisture sensors and weather monitoring systems. Furthermore, this layer is responsible for operating robotic and drone actuators to assist in the mobility of intelligent machines within the agricultural area. It thus allows intelligent machines to move between locations, in order to cover a wide area.

- *SSA Domain Layer:* This layer hosts various different Smart Sustainable Architectural domains and forms the main source for data generated from various agricultural domains, including: fields; machines; human resources; and crops. It forms the basis for IOT/AI SSA platforms, as it contains various raw data formats without interference.

*2)* The Second component: data lifecycle within SSA Architecture and its process location:

Fig. 5 demonstrates that data lifecycle remains in line with SSA Architecture layers. This commences with the original data source, i.e. SSA domain layers. The acquisition and capturing of data is undertaken at the layer containing sensors and actuators. These captured data are then passed to the application layer for business logic and control, following which, the data must be checked for security issues before moving to data analysis and processing. This is followed by data classification and transformation of the analyzed and processed data before it moves to data interpretation and the resulting building decisions. The final stage is to store the data for future retrieval.

There are two locations for sharing and processing of generated data. Firstly, on-site. Here, the generated data is more likely to be shared and processed within the location of the agricultural field, covering the four architecture layers: (1) SSA; (2) the domain layer; (3) the IoT layer; and (4) the application layer. Secondly, off-site: here the sharing and processing of the generated data must be outsourced to the physical location 'data centre', in which data is processed and analyzed away from the field. It can cover the four architecture layers: (1) network; (2) AI and data management; (3) physical; and (4) storage. It is also important to highlight that the Security layer forms a common layer between on-site and off-site processing locations.



Fig. 4.   Overall AI/IoT Platform Technical Architecture for SSA.



Fig. 5.   IoT/AI Data Lifecycle.

## V. CONCLUSION

This paper has established the importance of employing recent and advanced computing technologies in the agricultural sector, in particularly AI and IoT. Agriculture is considered central to the survival of human beings. Supporting the current practices of traditional agriculture with recent IoT/AI technologies can improve the performance, quality and volume of production. This study has reviewed the existing IoT/AI technologies discussed within the main research journals in the area of agricultural. Furthermore, it categorized the main domains of smart, sustainable agriculture, i.e. human resources; crops; weather; soil; pests; fertilization; farming products; irrigation/water; livestock; machines; and fields. The major contribution of this paper concerns the AI/IoT technical architecture for SSA, leading to an emphasis on the research and development of a unified AI/IoT platform for SSA, to positively resolve issues resulting from the fragmentary nature of the agricultural process. Future work will include investigation of the process of implementing AI/IoT technologies for SSA by applying the proposed AI/IoT technical architecture in the form of the prototype of a unified platform on real test cases. This will identify the relevant strengthens and weaknesses for further improvement and enhancement.

### REFERENCES

[1] H. Channe, S. Kothari, and D. Kadam, "Multidisciplinary model for smart agriculture using internet-of-things (IoT), sensors, cloud-computing, mobile-computing & big-data analysis," Int. J. Computer Technology & Applications, vol. 6, no. 3, pp. 374-382, 2015.

[2] L. Taiz, "Agriculture, plant physiology, and human population growth: past, present, and future," Theoretical and Experimental Plant Physiology, vol. 25, pp. 167-181, 2013.

[3] C. Smita and G. Shivani, "Smart Irrigation Techniques for Water Resource Management," in Smart Farming Technologies for Sustainable Agricultural Development, C. P. Ramesh, G. Xiao-Zhi, R. Linesh, S. Sugam, and V. Sonali, Eds. Hershey, PA, USA: IGI Global, 2019, pp. 196-219.

[4] V. Scognamiglio, Antonacci, A. , Lambreva, M. D., Arduini, F. , Palleschi, G. , Litescu, S. C., Johanningmeier, U. and Rea, G. , "Application of Biosensors for Food Analysis," in Food Safety, 2016.

[5] A. Antonacci, F. Arduini, D. Moscone, G. Palleschi, and V. Scognamiglio, "Nanostructured (Bio) sensors for smart agriculture," TrAC Trends in Analytical Chemistry, vol. 98, pp. 95-103, 2018.

[6] P. S. Kumar and G. J. Joshiba, "Water Footprint of Agricultural Products," in Environmental Water Footprints: Springer, 2019, pp. 1-19.

[7] E. Bennett et al., "Toward a more resilient agriculture," Solutions, vol. 5, no. 5, pp. 65-75, 2014.

[8] O. Elijah, T. A. Rahman, I. Orikumhi, C. Y. Leow, and M. N. Hindia, "An Overview of Internet of Things (IoT) and Data Analytics in Agriculture: Benefits and Challenges," IEEE Internet of Things Journal, vol. 5, no. 5, pp. 3758-3773, 2018.

[9] S. de Wilde, "The future of technology in agriculture," 2016.

[10] P. Jayashankar, S. Nilakanta, W. J. Johnston, P. Gill, and R. Burres, "IoT adoption in agriculture: the role of trust, perceived value and risk," Journal of Business & Industrial Marketing, vol. 33, no. 6, pp. 804-821, 2018.

[11] J. Leventon et al., "Collaboration or fragmentation? Biodiversity management through the common agricultural policy," Land Use Policy, vol. 64, pp. 1-12, 2017.

[12] E. Timofti, D. Popa, and B. Kielbasa, "Comparative analysis of the land fragmentation and its impact on the farm management in some EU countries and Moldova," Scientific Papers: Management, Economic Engineering in Agriculture & Rural Development, vol. 15, no. 4, 2015.

[13] P. R. Crosson, "Sustainable agriculture," in Global Development and the Environment: Routledge, 2016, pp. 61-68.

[14] M. Aleksandrova. (2019) Technologies and IoT have the potential to transform agriculture in many aspects. Namely, there are 5 ways IoT can improve agriculture. Eastern Peak.

[15] I. Mohanraj, K. Ashokumar, and J. Naren, "Field monitoring and automation using IOT in agriculture domain," Procedia Computer Science, vol. 93, pp. 931-939, 2016.

[16] D. I. Patrício and R. Rieder, "Computer vision and artificial intelligence in precision agriculture for grain crops: A systematic review," Computers and Electronics in Agriculture, vol. 153, pp. 69-81, 2018/10/01/ 2018.

[17] S. E. Bibri, "The IoT for smart sustainable cities of the future: An analytical framework for sensor-based big data applications for environmental sustainability," Sustainable Cities and Society, vol. 38, pp. 230-253, 2018.

[18] A. Kumar, L. Byung Gook, L. HoonJae, and A. Kumari, "Secure storage and access of data in cloud computing," in ICT Convergence (ICTC), 2012 International Conference on, 2012, pp. 336-339, New Jersey, USA: Institute of Electrical and Electronics Engineers (IEEE).

[19] M. Armbrust et al., "A view of cloud computing," Communications of the ACM, vol. 53, no. 4, pp. 50-58, 2010.

[20] W. Jiyi, P. Lingdi, G. Xiaoping, W. Ya, and F. Jianqing, "Cloud Storage as the Infrastructure of Cloud Computing," in Intelligent Computing and Cognitive Informatics (ICICCI), Kuala Lumpur, Malaysia, 2010, pp. 380-383, New Jersey, USA: Institute of Electrical and Electronics Engineers (IEEE), 2010.

[21] J. Repschlaeger, S. Wind, R. Zarnekow, and K. Turowski, "A Reference Guide to Cloud Computing Dimensions: Infrastructure as a Service Classification Framework," in System Science (HICSS), 2012 45th Hawaii International Conference on, 2012, pp. 2178-2188, Washington, USA: IEEE Computer Society, 2012.

[22] P. Mell and T. Grance, "The NIST definition of cloud computing," National Institute of Standards and Technology, vol. 53, no. 6, pp. 1-3, 2009.

[23] Q. Zhang, L. Cheng, and R. Boutaba, "Cloud computing: state-of-the-art and research challenges," Journal of Internet Services and Applications, vol. 1, no. 1, pp. 7-18, 2010.

[24] D. Chornyi, J. Riediger, and T. Wolfenstetter, "Into The Cloud," 2010.

[25] A. Marinos and G. Briscoe, "Community cloud computing," in Cloud Computing: Springer, 2009, pp. 472-484.

[26] K. Gai, M. Qiu, H. Zhao, L. Tao, and Z. Zong, "Dynamic energy-aware cloudlet-based mobile cloud computing model for green computing," Journal of Network and Computer Applications, vol. 59, pp. 46-54, 2016.

[27] C. V. Raja, K. Chitra, and M. Jonafark, "A Survey on Mobile Cloud Computing," International Journal of Scientific Research in Computer Science, Engineering and Information Technology, vol. 3, no. 3, 2018.

[28] A. u. R. Khan, M. Othman, F. Xia, and A. N. Khan, "Context-Aware Mobile Cloud Computing and Its Challenges," IEEE Cloud Computing, vol. 2, no. 3, pp. 42-49, 2015.

[29] S. Rajeswari, K. Suthendran, and K. Rajakumar, "A smart agricultural model by integrating IoT, mobile and cloud-based big data analytics," in 2017 International Conference on Intelligent Computing and Control (I2C2), 2017, pp. 1-5: IEEE.

[30] K. Sugahara, "Traceability system for agricultural productsbased on RFID and mobile technology," in International conference on computer and computing technologies in agriculture, 2008, pp. 2293-2301: Springer.

[31] J. V. Abellan-Nebot and F. R. Subirón, "A review of machining monitoring systems based on artificial intelligence process models," The International Journal of Advanced Manufacturing Technology, vol. 47, no. 1-4, pp. 237-257, 2010.

[32] S. J. Russell and P. Norvig, Artificial intelligence: a modern approach. Malaysia; Pearson Education Limited, 2016.

[33] S. S. Kale and P. S. Patil, "Data Mining Technology with Fuzzy Logic, Neural Networks and Machine Learning for Agriculture," in Data Management, Analytics and Innovation: Springer, 2019, pp. 79-87.

[34] K. Sennaar. (2019, April). AI in Agriculture – Present Applications and Impact. Available: https://emerj.com/ai-sector-overviews/ai-agriculture-present-applications-impact/

[35] K. Lakhwani, H. Gianey, N. Agarwal, and S. Gupta, "Development of IoT for Smart Agriculture a Review: Proceedings of ICETEAS 2018," 2019, pp. 425-432.

[36] F. Ferrández-Pastor, J. García-Chamizo, M. Nieto-Hidalgo, J. Mora-Pascual, and J. Mora-Martínez, "Developing ubiquitous sensor network platform using internet of things: Application in precision agriculture," Sensors, vol. 16, no. 7, p. 1141, 2016.

[37] allMETEO. (2019, April). allMETEO. Available: https://www.allmeteo.com/

[38] S. Elements. (2019, April). Smart Elements. Available: https://smartelements.io/

[39] Pycno. (2019, April). Pycno. Available: https://www.pycno.co/

[40] Farmapp. (2019, April). Farmapp. Available: https://farmappweb.com/

[41] Growlink. (2019, April). Growlink. Available: http://growlink.com/

[42] GreenIQ. (2019, April). GreenIQ. Available: https://easternpeak.com/works/iot/

[43] Arable. (2019, April). Arable. Available: https://arable.com/

[44] Semios. (2019, April). Semios. Available: http://semios.com/

[45] SCR/Allflex. (2019, April). SCR/Allflex. Available: http://www.scrdairy.com/

[46] Cowlar. (2019, April). Cowlar. Available: https://cowlar.com/

[47] FarmLogs. (2019, April). FarmLogs. Available: https://farmlogs.com/

[48] Cropio. (2019, April). Cropio. Available: https://about.cropio.com/#agro

[49] Farmshots. (2019, April). Farmshots. Available: http://farmshots.com

[50] aWhere. (2019, April). aWhere. Available: https://www.awhere.com

[51] Plantix. (2019, April). Plantix. Available: https://plantix.net/en

[52] T. Genomics. (2019, April). Trace Genomics. Available: https://www.tracegenomics.com/#/

[53] SkySquirrel. (2019, April). SkySquirrel. Available: https://www.skysquirrel.ca/#productnav

[54] S. Spray. (2019, April). See & Spray. Available: http://smartmachines.bluerivertechnology.com

[55] CROO. (2019, April). CROO. Available: https://harvestcroo.com

[56] S. Chen, H. Xu, D. Liu, B. Hu, and H. Wang, "A vision of IoT: Applications, challenges, and opportunities with china perspective," IEEE Internet of Things journal, vol. 1, no. 4, pp. 349-359, 2014.

[57] P. P. Ray, "Internet of things for smart agriculture: Technologies, practices and future direction," Journal of Ambient Intelligence and Smart Environments, vol. 9, no. 4, pp. 395-420, 2017.

[58] M. S. Mekala and P. Viswanathan, "A Survey: Smart agriculture IoT with cloud computing," in 2017 International conference on Microelectronic Devices, Circuits and Systems (ICMDCS), 2017, pp. 1-7: IEEE.

[59] A. Kamilaris, A. Kartakoullis, and F. X. Prenafeta-Boldú, "A review on the practice of big data analysis in agriculture," Computers and Electronics in Agriculture, vol. 143, pp. 23-37, 2017.

[60] H. Cadavid, W. Garzon, A. Pérez Ruiz, G. López, C. Mendivelso, and C. Ramírez, "Towards a Smart Farming Platform: From IoT-Based Crop Sensing to Data Analytics: 13th Colombian Conference, CCC 2018, Cartagena, Colombia, September 26–28, 2018, Proceedings," 2018, pp. 237-251.

[61] G. Soto-Romero, J. Roux, C. Escriba, J.-Y. Fourniols, and G. Soto-Romero, "A new bi-frequency soil smart sensing moisture and salinity for connected sustainable agriculture," Journal of Sensor Technology, 2019.

[62] L. Nóbrega, P. Gonçalves, P. Pedreiras, and J. Pereira, "An IoT-Based Solution for Intelligent Farming," Sensors, vol. 19, no. 3, p. 603, 2019.

# Reconstruction of Fingerprint Shape using Fractal Interpolation

Abdullah Bajahzar[1]

Department of Computer Science and Information
College of Science, Majmaah University
Zulfi 11932, Saudi Arabia

Hichem Guedri[2]

Electronics and Microelectronics Laboratory
Physics Department, Faculty of Sciences, Monastir
University, 5019 Monastir, Tunisia

*Abstract*—One of the severe problems in a fingerprint-based system is retaining the fingerprint images. In this paper, we propose a method to minimize the fingerprint images size and retain the reference points. The method is divided into three parts, the first part is about digital image preprocessing that allows us to eliminate the noise, improve the image, convert it into a binary image, detect the skeleton and locate the reference point. The second part concerns the detection of critical points by the Douglas-Peucker method. The final part presents the methodology for the fingerprint curves reconstruction using the fractal interpolation curves. The experimental result shows the accuracy of this reconstruction method. The relative error (ER) is between 2.007% and 5.627% and the mean squared error (MSE) is between 0.126 and 0.009 at a small iterations number. On the other hand, for a greater number of iterations, the ER is between 0.415% and 1.64% and MSE is between 0.000124 and 0.0167. This clearly indicates that the interpolated curves and the original curves are virtually identical and exceedingly close.

*Keywords*—*Fingerprint images; enhancement; thresholding process; thinning algorithms; minutiae extraction; Douglas-Peucker algorithm; fractal interpolation; iterated function system (IFS)*

## I. INTRODUCTION

In nowadays, the world has become electronically connected and more dynamic because the technological revolution. The traditional identity representations such as passwords and cards cannot be trusted to identify a person. The cards may be lost or stolen, and passwords or PIN numbers may be compromised. In addition, passwords and cards can be easily shared, thus our personal information cannot be secured [1-3]. There are innovative techniques for identifying individuals called biometric identifiers [1]. The biometric identification uses distinct anatomical properties such as fingerprints, facial, iris and behavior (such as speech). The biometrics offers an effective solution to identify people because biometric identifiers cannot be lost or assumed, and represent essentially the individual's body identity.

The fingerprint is considered one of the most practical to identification persons. The fingerprint recognition requires minimal effort from the user, since it allows us to capture unique information necessary for the recognition process, and provides us with relatively excellent performance. Another reason of the fingerprints popularity is the relative minimum price of fingerprint sensors, which for peaceful integration into many companies. But, the downside of this technique is that

they require too much memory or storage space to save all the fingerprint images.

There are specific techniques to describing, analyzing, understanding, modeling, and controlling complex processes. Fractal modeling makes it possible to manage complex elements by starting with a reduced number of classical geometric shapes. Fractal interpolation can be found in many applications at both the 1D, 2D and 3D levels. Moreover, the use of fractal models, both in terms of the signal and the image, has become commonplace and constitutes an active search, motivated by the plethora of possible applications based on this concept. In this paper, a study is developed to reconstruct fingerprint curves using the fractal interpolation method.

Some pretreatment and enhancement steps are frequently performed to simplify the minutiae extraction task [4-8]. The first step of the algorithm concerns the fingerprint image segmentation; this phase requires the conversion of the gray-scale fingerprints image into a binary image [9-11]. The binary images obtained by the binarization process are generally subjected to a thinning step [12-13] which makes it possible to reduce the thickness of the peak line to one pixel [14-15]. Once a binary skeleton has been obtained, a simple image scan can detect the minutiae [16-18]. The further step is to determine the control points implementing the Douglas-Peucker algorithm [19-24]. In the final step, the reconstruction approach was presented by using the fractal interpolation algorithm and the iterated function system (IFS) [25-32].

The rest of the paper is arranged as follows. In Section 2, brief description of the related work. The proposed method is illustrated in Section 3; Section 4 describes a general overview in the fingerprint image segmentation and detection of the minutiae, the Douglas-Peucker algorithm is presented in Section 5, Section 6 describes the fractal interpolation, the illustrative simulation and experimental results are presented in Sections 7, finally followed by the discussion and the conclusion in Sections 8 and 9, respectively.

## II. RELATED WORK

A brief discussion is presented in this section on related work; fractal reconstruction techniques have been widely studied. Lai et al. [25] have used fractal interpolation for the compression and reconstruction of two-dimensional digital elevation model (2D DEM). In their proposed method, they have used the improved Douglas - Peucker (IDP) technique to extract feature points. They used fractal interpolation for

reconstruction, and they proposed a probability-based method to accelerate the fractal interpolation execution. Li et al. [26] have used the fractal interpolation technique for the seismic data reconstruction. They have examined the fractal interpolation method based on previous work and theoretically analyzed a special type of fractal interpolation function. They have ascertained that the numerical results of the fractal interpolation method have a high accuracy and efficiency, the largest error between the theoretical seismograms and the reconstructed seismograms within limits 3%. Cader and Krupski [27] have introduced a new interpolation method for fractal curves. They have studied curves that have a very irregular character, this type of curve has coarse characteristics and has a complex structure at Different metrics, and they numbered among fractals or stochastic fractures - multifractals. They proposed another alternative method of using fractal curves for the complex curves approximation, his method is better than (FIF) for the multifractal structures interpolation. It bases this on the classical notion of interpolation node and introduces non-local values for its description (the fractal dimension). Guedri et al. [28] have introduced the fractal interpolation technique for the human retina blood vessels reconstruction. They have studied the segmentation technique of the retinal image (such as binarization and skeletonization)

and the line simplification by using the Douglas-Peucker method. And finally, they have used fractal interpolation and IFS for the blood vessel curves reconstruction.

### III. PROPOSED APPROACH

The proposed method considers three necessary phases; preprocessing phase of fingerprint images, phase for the characteristic points determination using the Douglas-Peucker algorithm and the fractal interpolation phase. The structure of these phases discussed in the flowchart is illustrated in Fig. 1.

In the proposed method, the first objective is the enhancement and Binarization the fingerprint image which consists of transforming a multilevel image into a black and white image (two levels only). Subsequently the phase of thinning; at this point, the fingerprint is presented as a set of connected curves, while keeping its original topology. After obtaining the image skeleton, the most interesting minutiae points are detected. Subsequently, a lines simplification step applied to reduce the lines complexity (the control point's detection) by implementing the Douglas-Peucker algorithm. In the final step, the fractal interpolation and iterated function systems (IFS) has been used to visualize the result of the fingerprint image reconstruction.



Fig. 1. Schematic Diagram for a Fractal Reconstruction the Fingerprints Image.

## IV. PRETREATMENT OF FINGERPRINT IMAGES

In this approach, the grayscale image is converted to a binary image before the minutia detection. The proposed method is composed of diverse parts; it includes the following steps:

- Enhancement: In the First step, to eliminate unnecessary noise, the Gabor filter algorithm is used, which is based on the technique proposed by Lin Hong [4] to enhance fingerprint images [4-6]. This image is convolved with a Gabor filter with uniform symmetry to improve fingerprints [7-8].

- Binarization: the binarization phase is a very important step in image segmentation; the thresholding process is used to create a binary image from grayscale image [9]. In this approach, each pixel is examined and assigned his luminance. If a pixel is in a region where the low luminance (lower luminance to a set ready level value), then it is replaced by the value 0 (black color). However, if a pixel is in a region where the luminance is clearly defined, it is replaced by the value 1 (white color) [10-11].

- Skeletonisation: The purpose of skeletonization is to extract a shape feature from an object and present this object through a connected curves set or digital arcs with a reduced data amount or a simplified form, and remain the original object topology [12-15]. There are different categories of skeletonisation methods (Distance Skeleton Algorithms, Critical Points Algorithms, and Thinning Algorithms); in this work, the thinning technique is used to create the skeleton image. This technique based on the thinning process, it peels the object contour until the medial line with a one-pixel width, and the extracted skeleton retains the topology of the original form.

- Minutiae extraction: After obtaining the skeleton image, the objective principle is to detect the most interesting minutiae from this image. Rutovitz's approach [16] is implemented to detect minutiae and their end point or bifurcation types [17-18], as shown in Fig. 2(a) and Fig. 2(b), by calculating the Crossing-Number indicator (CN) according to the each pixel neighborhood, according to the CN value the point type is determined.

- The crossing-number of a pixel P is calculated by the following equation (1) [16]:

$$CN = \frac{1}{2} \sum_{i=1}^{8} |P_i - P_{i+1}|$$

(1)

- $P_i$ is the value of the pixels in the 3 * 3 neighborhood of P.

- $P_1$, $P_2$,...., $P_8$ are the 8 pixels in the neighborhood of P

- Indeed, the value CN allows us to identify the nature of a minutiae according to the result obtained during the computation of CN, as illustrated in the Table 1.



Fig. 2. Example of the Two Most used Minutiae, (a) end Point, (b) Bifurcation.

TABLE I. IDENTIFICATION OF THE MINUTIAE NATURE FROM THE CN VALUE

| CN value | Nature of a minutia |
|---|---|
| 0 | Isolated point |
| 1 | Endpoint |
| 2 | Connective point |
| 3 | Bifurcation point |
| >=4 | Crossing point |

## V. SIMPLIFICATION ALGORITHM

Initially, a line simplification algorithm is used to reduce the curve point number by a point's series, while preserving its shape. The points types selected should provide valuable indications for utilizing it later in the reconstruction of this curve. There are many line simplification algorithms such as the Area-Based Method, the Douglas-Peucker algorithm, Visvalingam Algorithm,... Depending on White study [22] of three simplification algorithms based on Marino work [23], he showed that line simplification produced by Douglas-Peucker they are the best lines examples with 86% of all tests. Because of these results, the Douglas-Peucker method represents the best of the three suggested methods. Also the McMaster [24] study which showed that the Douglas-Peucker algorithm was mathematically and perceptually superior to other algorithms, this algorithm detected more critical points and kept the original curve topology because it produces the least displacement compared to the original curve. The simplification of a 2D line using the Douglas - Peucker algorithm defined by this algorithm (1):

| **Algorithm 1. Douglas-Peucker Algorithm** |
| --- |
| **INPUT**⟶Curve C (an ordered set of points (P1, ..., Pn)) and the threshold value Ɛ (Ɛ> 0). <br> **OUTPUT**⟶ ResultPoint (control points) <br> **function DouglasPeucker(C, Ɛ)** <br> k⟶ length (P1, ..., Pn) <br> Dmax⟶ NULL  // initialization of the maximum orthogonal distance <br> index⟶ NULL <br> **For** i=2 **to** k-1 **do** <br> Dnew = orthogonal distance(P(1),P(i),P(n)); <br>    **if** Dnew > Dmax <br>     index = i <br>     Dmax = Dnew <br>    **End If** <br>   **End For** <br> **if Dmax > Ɛ** <br>     C1=P(1) to P(index) <br>     C2= P(index) to Pn <br>     ResultPoint1=function DouglasPeucker(C1, Ɛ) <br>     ResultPoint2=function DouglasPeucker(C2, Ɛ) <br>     ResultPoint={ ResultPoint1, ResultPoint2} <br> **else** <br>     ResultPoint={P1,Pn} <br> **End If** |

## VI. FRACTAL INTERPOLATION

In life, they are many complex phenomena and natural forms that cannot be characterized in a classical way. All these complex phenomena may have specific characteristics. They can be complex in the sense of multiscale, scale invariance or self-similarity [30-32].

Mandelbrot gives mathematical models as well as the tools of analysis of these phenomena in his book "fractal geometry of nature" [32]. In addition, there are several studies that are interested in fractal study, such as Barnsley studies [30-31], which found that fractal phenomena can be obtained with deterministic models. And that they are based on sets of affine transformations, are called IFS (Iterated Function Systems) [30-31]. The iterative application of these transformations produces a form that has the property of being similar to it at several scales. Further details on this aspect will be given in the following section of the document.

### A. Presentation

Consider the interpolation points:

$$P_i = \{(x_i, y_i): i = 0, 1, \dots, N\} \in \Re^2 \text{ and } N \in \Box$$

A geometric transformation $W$ can be given by equation (2) the coordinates $x_i^{'}$ and $y_i^{'}$ of a point $P_i^{'}$, image of $P_i$ by W, according to the coordinates $x_i$ and $y_i$ of the point $P_i$ [25-27].

$$W = \begin{cases} x_i^{'} = f(x_i, y_i) \\ y_i^{'} = g(x_i, y_i) \end{cases} \tag{2}$$

Where $x_i^{'} = f(x_i, y_i)$ and $y_i^{'} = g(x_i, y_i)$ are 2-variable functions $x_i$ and $y_i$

### B. Affine Transformations

An affine transformation can scale, distort, rotate, and convert data differentially. a transformation can be described as a function that sends a point $P_i(x_i, y_i)$ from the plane to another point $P_i^{'}(x_i^{'}, y_i^{'})$, let $x_i^{'}$ and $y_i^{'}$ the coordinates of a point $P_i^{'}$, according to the coordinates $x_i$ and $y_i$ of $P_i$ [27-29].

The transformation is called affine if it can be written in the following equation (3):

$$W_i = AP_i + T = \begin{bmatrix} a_i & b_i \\ c_i & d_i \end{bmatrix} * \begin{pmatrix} x_i \\ y_i \end{pmatrix} + \begin{pmatrix} e_i \\ f_i \end{pmatrix} \tag{3}$$

It is represented from the matrix $A$ where the sub-matrix 2x2 denotes the combined effects of the transformations, rotation and scale change, while the vector $T$ denotes the combined effects of the translation transformations.

It can be written as an equation system following:

$$W_i = \begin{cases} x_i^{'} = a_i x_i + b_i y_i + e_i \\ y_i^{'} = c_i x_i + d_i y_i + f_i \end{cases}$$

With a, b, c, d, e and f are real numbers.

Where $x_i$ and $y_i$ are the source control points coordinates and $x_i^{'}$ and $y_i^{'}$ the transformed coordinates [29].

- The IFS fundamental idea: The IFS fundamental idea: The IFS model is exclusively based on the self-similarity notion. The idea is to encode a fractal object (in a whole or in part with a similar structure) by a transformation set translating this property; it is this similarity, which allows utilizing the iterative method to generate a curve or an image. The transformations set will be called IFS and the corresponding fractal will be the attractor of the IFS.

The self-similarity notion is expressed by the equation (4):

$$K = \bigcup_{i=1}^{N} w_i(K) \tag{4}$$

In general, the fractal object K is equal to the transformation union of itself.

Let f be a continuous fractal function passing through a certain given points number of the form $f(P_i, F_i) \in \Re^2, i = 0, \ldots\ldots\ldots, N$ and with $P_0 < P_1 < \ldots\ldots\ldots < P_N$

In this study, we rely on the fractal interpolation principle studied by Barnsley [31-32]. An interpolation function corresponding to this data set declares above is a continuous function $f[P_0, P_N] \rightarrow \Re$ passing through the interpolation points $(P_i, F_i) \in \Re^2$ and checking $f(P_i) = F_i$ with $i = 0, 1, \ldots\ldots\ldots N$. The function graph f (he is the attractors set) having N transformations of the form:

$$\begin{pmatrix} x_i^{'} \\ y_i^{'} \end{pmatrix} = W_i \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} a_i & 0 \\ c_i & d_i \end{pmatrix} * \begin{pmatrix} x_i \\ y_i \end{pmatrix} + \begin{pmatrix} e_i \\ f_i \end{pmatrix}$$

The transformations $W_i$ are defined with the five real numbers $a_i$, $c_i$, $d_i$, $e_i$ and $f_i$.

The resolution of the preceding constraints makes it possible to define the transformations parameters $W_i$ of the IFS generating an attractor passing by the points set $P_i$ with i = 0,1,………,N. The transformations parameters are thus defined by the following equations (5) [29]:

$$a_i = \frac{x_i - x_{i-1}}{x_N - x_1}$$

$$c_i = \frac{y_i - y_{i-1}}{x_N - x_1} - d_i * \frac{y_N - y_1}{x_N - x_1} \tag{5}$$

$$e_i = \frac{x_N x_{i-1} - x_1 x_i}{x_N - x_1}$$

$$f_i = \frac{x_n y_{i-1} - x_1 y_i}{x_N - x_1} - d_i * \frac{x_N y_1 - x_1 y_N}{x_N - x_1}$$

There remains therefore a freedom degree symbolized by the parameter $d_i$ and representing a "vertical scale factor" or "contraction factor"[29]. Its value is independent of the interpolation points and controls the shape and roughness between the interpolation points of the graph. This parameter is calculated using the fractal dimension (DF) [33-34] and is expressed by the following equation (6):

$$d_i = (N-1)^{DF-2} \tag{6}$$

With the fractal dimension (DF) is calculated by the box-counting method [33-34].

*C. Implantation of Fractal Interpolation and IFS*

Let us consider a points set $P_i$ {$(x_i, y_i)$} i = 0,1, ..., N that we are trying to interpolate (rhe control points detected by the Douglas-Peucker algorithm). The IFS theory can be used to interpolate [29-30]. The *N* affine transformations are used for this purpose. This *N* transformations partition the interval [P0, PN] into [P0, P1] ∪ [P1, P2] ∪. . . ∪ [PN-1, PN]. We present in the remainder of the present section the fractal interpolation algorithm [28-29].

The implementation of this algorithm is a function called IFS and uses affine transformations in $R^2$, as you can see in algorithm (2).

| Algorithm 2. IFS Algorithm |
|---|
| **INPUT** ⟶ $P_i$, $P_0$, N, K |
| **OUTPUT** ⟶ $P_{new}$ |
| **function IFS($P_i$ , P0, N, K)** |
| h ⟶ length ($P_i$) |
| **For** i=1 **to** h **do** |
| $W_i$ ⟵ The parameters ai, ci, di, ei, fi are calculated using relations (5) - (6) |
| **End For** |
| **For** i=1 **to** h **do** |
| $P_{new}$ (i) ⟵ $W_i(P_0)$ |
| **End For** |
| N ⟵ length ($P_{new}$) |
| **For** j=1 **to** K **do** |
| **For** i=1 **to** h **do** |
| N ⟵ N+i |
| $P_{new}$ (N) ⟵ $W_i(P_{new})$ |
| **End For** |
| **End For** |

The input parameters are:

- P$_i$: The control points set.

- P$_0$: The initial points set (a matrix with the coordinates of the initial M points).

- K: integer type determines the iterations number.

The output is a point set P$_{new}$, in matrix form.

The algorithms begin by calculating the IFS map parameters Wi with using the data points set or control points Pi. Subsequently, we initialize the algorithm with a randomly selected points set. In the first step, apply an IFS map to these selected points; a new point set will emerge. At each step, repeat the same technique to the new points set obtained during the previous iteration until we obtain the original curve or image.

### D. Evaluation of Performances

An interpolation function represent a approximating function, there is no interpolation type can be assumed as fully accurate. It is therefore meaningful to recognize to what extent the measured value may deviate from the initial value, the quantity variation, it is called error analysis [35-39]. To evaluate the interpolation accuracy as well as the efficiency of the fractal interpolation algorithms, There are two kinds of error computation, the first represent the relative error (RE) [35-37] and the second is the mean square error (MSE) [37-39].

Relative error (RE) used as an accuracy measure—represent the ratio of a measurement absolute error to the current measurement. Put differently, the RE (also called fractional error) is obtained by dividing the absolute error in the quantity by the quantity itself. RE is expressed as a percentage and he includes no units [35-37]. The relative error (RE) of the potentials is defined by the following equation (7):

$$RE = \frac{\square x}{x} \qquad (7)$$

Where x is the true value of a quantity and $\square x$ is the absolute error.

The mean squared error (MSE) or mean squared deviation (MSD) [37-38] is an estimator measures the errors mean squares, that is, the mean squared difference between the origin and what is interpolated [38-39]. It is described by as follows equation (8):

$$MSE = \frac{1}{n} \sum_{i=1}^{n} \left( P_i - \hat{P_i} \right)^2 \qquad (8)$$

### VII. RESULT

To evaluate the proposed approach performance, this algorithm uses fingerprint images from the Fingerprint Verification Competition database (FVC2006)[1]. This database consists of four different sub-databases (DB1, DB2, DB3, and DB4), the first three acquired with different sensors and the last created with a synthetic generator; each fingerprint image was

captured at a 500 dpi resolution with various sizes. The size of each database to be used in the test was set to 110 fingers wide (w) and eight impressions per finger depth (d) (880 fingerprints in total).And they are devised in two sets:

Set A: consisting of fingers numbered from one to 100

Set B: made up of fingers numbered 101 to 110 and made available to users.

The results obtained were implemented using the software platform MATLAB (R2015a). The program was assessed on a personal computer running at 2.20 GHz for CPU and 4 GB of RAM.

### A. Segmentation Steps

Initially, the fingerprint image is extracted from the database, the raw image or original image. The initially, the fingerprint image is extracted from the database, the raw image or initial image. Thereafter, the Gabor filter is used to enhance fingerprint images through eliminating its unnecessary noises. Then the binarization phase to transform the greyscale image into a binary image by the thresholding method. Subsequently, the binarized images are thinned implementing a thinning algorithm to reduce the lines width to one-pixel width of the skeleton. This thinning process does not change the original fingerprint topology that ensures efficient of the minutiae points extraction. Ultimately, scanning the fingerprint skeleton image allows detecting the pixels that correspond to minutiae by calculating the Crossing Number Indicator (CN).

Taking the example of the fingerprint image 101_1.tif [1], as you can see in Fig. 3(a) (raw image) taken from DB1 in FVC2006 to verify and applying the same segmentation steps, the results are presented in (Fig. 3), where Fig. 3(b) shows the enhanced version of the fingerprint image . Fig. 3(c) represents its binary image, Fig. 3(d) shows its thinned version, Fig. 3(e) and Fig. 3(f) represent the Endpoints and the bifurcation points, respectively.

The results of six fingerprint images are presented in the table below (Table 2). The first column is for illustration the minutiae number of type 1 (Endpoints) and column 2 for the minutiae number of type 2 (bifurcation). Column 3 and 4 shows the run time in second for each segmentation step.

Table 2 shows the results obtained during the segmentation steps. It should be noted that the execution time doesn't depend on the image size, but depends on another criterion, such as the image quality in the enhancement task, the execution time increases for low quality images. Second criterion, the end and bifurcations points number contained in the image. The average execution time for the enhancement task of the fingerprint images is between 4 and 6 second. On the other hand, the execution time for the minutia extraction is between 0.5 and 0.7 second. As you can see, the enhancement task is extremely costly in time terms. The run time obtained and that found by S. Sojan and R. K. Kulkarni [8] are almost identical (average execution times between 0.4 s and 0.6 s). The minutiae points extracted are visible in column 2 and 3. About 104 up to 173 endpoints and 20 up to 59 bifurcations were detected, as can be seen from Table 2.

---

[1] Fingerprint Verification Competition 2006 (FVC2006) web site (http://bias.csr.unibo.it/fvc2006/)

Fig. 3. (a) Original Fingerprint Image (Raw image) (b) The Enhanced Image (c) Binarized Image (d) Thinned Image (e) Endpoints (Red Points) (f) Bifurcations (Bleu Points).

TABLE II. MINUTIAE EXTRACTION AND THE RUN TIME

| images | number of endpoint | number of bifurcation | execution time of enhancement task (Second) | execution time(Second) |
|---|---|---|---|---|
| **101_1.tif** | 126 | 29 | 4.94 | 0.548 |
| **101_2.tif** | 104 | 32 | 5.15 | 0.577 |
| **110_4.tif** | 138 | 41 | 6.319 | 0.606 |
| **102_5.tif** | 174 | 28 | 4.12 | 0.702 |
| **108_3.tif** | 134 | 59 | 4.71 | 0.638 |
| **108_5.tif** | 173 | 20 | 5.79 | 0.6164 |

### B. Douglas-Peucker Algorithm

An approach to simplifying data in curves is to use a line simplification algorithm. One of the most used algorithms is the Douglas-Peucker algorithm. This algorithm allows data compression and reduces the data points number by eliminating redundant points, while retaining its shape, which will save storage space and transmission cost. Douglas-Peucker is a recursive algorithm based on perpendicular distance and a given tolerance value ε.

The image pixel unit is used when it comes to the tolerance value. The figure below (Fig. 4) shows the simplified fingerprint curves with the Douglas-Peucker algorithm at a few different tolerance values between 0.5 to 2 Pixels.

As Fig. 4 shows, even with a low-precision simplification that results in a much smaller set of results, the overall curve shape remains the same. The first figure (Fig. 4(a)) shows a reduction for ε = 0.5 pixels of 75.37%, from 5113 to 1259 points. While, the twelfth figure (Fig. 4(c)) shows a reduction of 91.82% for ε = 1, from 5113 to 418 points. Whereas, for ε = 2 the Fig. 4(e) shows a reduction of 94.15 %, from 5113 to 299 points.

Table 3 illustrates the Douglas-Peucker algorithm results on the test images set. This table describes for four different tolerance values (between 0.5 and 2 pixels) the reduced points number in the third column, the simplification rate in the next column and the last column, the run time.

The values in the Table 3 clearly show that the control points number decreases if the tolerance value ε increases, and vice versa. As an example, for the first image test the control points number increases from 299 points (for ε = 2) to 1259 points (for ε = 0.5), and similarly for the simplification rate, he passes from 75% to 94 % for the same values of ε. From the

obtained results, the simplification rate is between 75 and 95% according to the value of ε. In addition, the execution time of the Douglas-Peucker algorithm is between 1s and 2.3s.

*C. Interpolation Fractal Tests*

After implementing the Douglas-Peucker algorithm and the characteristic point's detection, the fractal interpolation algorithm already stated in section VI is used to reconstruct the fingerprint curve. The IFS technique result for different iterations number (between 50 and 500 iteration) is shown in Fig. 5. It can be noted that the reconstruction quality is affected according to the iterations number, for the narrow iterations number (between 50 and 300) the reconstruction quality is poor, as you can shows in Fig. 5(a) and Fig. 5(b). On the other hand, for an important iteration number (between 400 and 500), as shown in Fig. 5(c) and Fig. 5(d), the interpolated curve is identical to the original curve.

Fig. 5 shows examples obtained results at various iterations. At 50 iterations, the deformation is weak but visible. At 300 iterations, the final result is almost obtained. On the other hand, at the 500 iterations, the final result is perfectly obtained. The tests models details as well the numerical results are presented in the figure below.

Fig. 6 shows the relative error (RE) as a function of the iterations number. Indeed, there is an average relative error REM = 5.48% at the iteration 50, REM = 3.165% and REM = 2.13% at the iteration 100 and 200, respectively. In addition, the average relative error is between 0.5 and 1.5% at iterations between 500 and 300, which indicates clearly that the interpolated curves and the original curves are practically identical and exceedingly close.



**(a)**                **(b)**                **(c)**



**(d)**                **(e)**

Fig. 4.   Simplification Rate for different Tolerance Values ε. (a) ε=0.5 Pixel (b) ε=0.8 Pixel, (c) ε=1 Pixel, (d) ε=1.5 Pixel, (e) ε=2 Pixel.

TABLE III.  SIMPLIFICATION RATE FOR DIFFERENT TOLERANCE VALUES E AND RUN TIME

| image | Data points number | $\varepsilon$ | Reduced points number | Simplification rate in% | Run time |
|---|---|---|---|---|---|
| 101_1.tif | | 0.5 | 1259 | 75.3765 | 1.345205 |
| | | 0.8 | 511 | 90.0059 | 1.231993 |
| | | 1 | 418 | 91.8248 | 1.196033 |
| | | 1.5 | 342 | 93.3112 | 1.176604 |
| | | 2 | 299 | 94.1522 | 1.161165 |
| 101_2.tif | | 0.5 | 2251 | 76.4564 | 2.331704 |
| | | 0.8 | 817 | 91.4549 | 2.040004 |
| | | 1 | 634 | 93.3689 | 2.017986 |
| | | 1.5 | 494 | 94.8332 | 1.933262 |
| | | 2 | 415 | 95.6594 | 1.932721 |
| 110_4.tif | | 0.5 | 2278 | 75.5579 | 1.763994 |
| | | 0.8 | 820 | 91.2017 | 1.439569 |
| | | 1 | 654 | 92.9828 | 1.392574 |
| | | 1.5 | 512 | 94.5064 | 1.361353 |
| | | 2 | 441 | 95.2682 | 1.339857 |
| 102_5.tif | | 0.5 | 2383 | 74.1400 | 1.812687 |
| | | 0.8 | 701 | 92.3928 | 1.542794 |
| | | 1 | 552 | 94.0098 | 1.467193 |
| | | 1.5 | 449 | 95.1275 | 1.536463 |
| | | 2 | 386 | 95.8112 | 1.488915 |
| 108_3.tif | | 0.5 | 3014 | 77.0903 | 2.091803 |
| | | 0.8 | 1111 | 91.5552 | 1.854681 |
| | | 1 | 851 | 93.5315 | 1.698356 |
| | | 1.5 | 697 | 94.7020 | 1.565364 |
| | | 2 | 578 | 95.6066 | 1.506735 |
| 108_5.tif | | 0.5 | 2735 | 74.3457 | 1.804292 |
| | | 0.8 | 804 | 92.4585 | 1.712217 |
| | | 1 | 606 | 94.3157 | 1.651870 |
| | | 1.5 | 491 | 95.3944 | 1.542717 |
| | | 2 | 410 | 96.1542 | 1.444087 |

Fig. 5. Fractal Reconstruction for different Iterations Values (White Points: Bullet Points, Red Points: Erroneous Points), (a) Iterations Number= 50, (b) Iterations Number= 100, (c) Iterations Number= 200, (d) Iterations Number= 300, (e) Iterations Number= 400, (f) Iterations Number= 500.



Fig. 6. The Relative Error Curve.

## Mean Squared Error (MSE)



Fig. 7. The Mean Squared Error Curve.

Fig. 7 represents the Mean Squared Error MSE measure that reflects the distance between the original curves points and corresponding interpolated curves points. According to Fig. 7. above, it can be noted that when the iterations numbers are great (between 300 and 500), the mean squared error is small (between 0,000124 and 0.012); on the other hand, when the iterations numbers are small (between 50 and inferior to 300), the mean squared error is relatively large (between 0.041 and 0.126), which corresponds exactly to the same remark in the previous section as the interpolated curves and the original curves are virtually identical and exceedingly close for the high iterations numbers.

## VIII. DISCUSSION

The results concerning the curves approximation with fractal interpolation is presented in this section. Validation tests were performed on real fingerprint image. With respect to the curvature error evolution using fractal interpolation, the obtained results show that the relative error (ER) and the mean squared error (MSE) are large when the iterations number is small (between 2.007% and 5.627% for RE and 0.126 and 0.009 for MSE). On the other hand, these errors (between 0.415% and 1.64% for RE and 0.000124 and 0.0167 for MSE) are weak when the iterations number is considerable. The method proposed by Jian-Kai et al. [40] for the average relative error between predicted and genuine values was 1.32%. For the method proposed by MAČĖNAITĖ et al [41], he calculated two errors types (Mean squared errors (MSE) and relative errors (RE)) between the autocorrelation functions of the real profilogram and its model. They obtained a relative error between 8.068% and 20.382%, mean squared errors (MSE) between 0.134 and 0.339.The method proposed by Guedri et al. [28] the relative errors between the ideal centre line and the simulated blood vessels centre line ranged from 1.341 to 12.608 %.

The results obtained and the comparison shows that the proposed method has succeeded in reconstructing very precisely the fingerprint curves; this study shows that the fingerprint curves can appropriately preserve its particular structure when the iterations numbers exceed 300 iterations.

The advantage of such an approach lies in the high simplification rate (can reach up to 96%), which makes it possible to reduce the memory size to storing the images of the fingerprint, which allows reduce the costs and the transmission time. In addition, the fractal interpolation can reconstruct the data points with higher resolution than initially model and the reconstructed model has more natural and real details. On the other hand, the proposed algorithm vulnerable point lies in the calculation complexity. This requires a specialized computer or calculator to minimize the run time.

## IX. CONCLUSION

This paper explains how to used fractal theory to data compression and interpolation of fingerprint curves. This method is divided into three essential parts; the first part is the image segmentation of the fingerprint (enhancement, image binarization, skeletonization and minutiae extraction). Subsequently, the second part integrates the Douglas-Peucker linear simplification method to reduce the images size. In closing, the third part is interested in the fingerprint curves reconstruction by fractal interpolation. This Research has shown that the data reduction rate of the fingerprint curves can reach 96%, it can be noted that the Douglas-Peucker algorithm has a high reduction rate and also a rapid reduction to 2 seconds, moreover , they have retains detailed characteristics for the fingerprint curves. Further, the general structure of reconstructed curves by fractal interpolation can obtain an excellent quality, validation tests were performed. The research shows that the error obtained between 0.415% and 1.64% for RE and between 0.000124 and 0.0167 for MSE with the iterations number larger than 300. This clearly indicates that the interpolated curves and the original curves are virtually identical and extremely close.

In the future work we will implement this method proposed within hardware platform for example on embedded architectures and reconfigurable such as FPGA programmable elements of the type "soft-core" (heart CPU generalist or core DSP).

REFERENCES

[1] N.Soundharadevi and M.Pushparani, "Analysing on multimodal biometric frame work with face, iris and fingerprint images," Shanlax International Journal of Arts, Science & Humanities, Vol. 4 (1), pp.165-172, September 2016.

[2] V. Conti, C. Militello and S. Vitabile, "Biometric authentication overview: a fingerprint recognition sensor description," Int J Biosen

Bioelectron, vol. 2(1), pp. 26–31, 2017. DOI: 10.15406/ijbsbe.2017.02.00011.

[3] A.mira Saleh, A. Bahaa and A. Wahdan, "Fingerprint Recognition," Advanced Biometric Technologies, InTech, 2011 pp.201- 224, DOI: 10.5772/23476

[4] L. Hong, Y. Wan, and A. K. Jain, "Fingerprint image enhancement: Algorithm and performance evaluation," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol 20(8), pp 777-789, 1998.

[5] A. Nagar, S. Rane and A. Vetro, "Alignment and Bit Extraction for Secure Fingerprint Biometrics," Conference: Media Forensics and Security II, part of the IS&T-SPIE Electronic Imaging Symposium, San Jose, CA, USA, January 18-20, 2010. DOI: 10.1117/12.839130

[6] T. Vidhya and T. K. Thivakaran, "Fingerprint Image Enhancement Using Wavelet Over Gabor Filters," Int. J. Computer Technology & Applications(IJCTA), Vol. 3 (3),pp 1049-1054, 2012.

[7] S. Mohammedsayeemuddin, S. K. Gonsai, and D. Vandra, "Efficient Fingerprint Image Enhancement Algorithm Based On Gabor Filter," International Journal of Research in Engineering and Technology, Vol. 3( 4),pp. 809-813, Apr 2014.

[8] S Sojan and R. K. Kulkarni, "Fingerprint Image Enhancement and Extraction of Minutiae and Orientation," International Journal of Computer Applications, vol 145(4),pp 14-19, 2016.

[9] W. Hussain, T. Munawar, M. Shahzaib, M. Masood, "Automated Enhancement of Compromised Fingerprint Images," Journal of Biochemistry, Biotechnology and Biomaterrials ( JBCBB), Vol. 1(2), pp. 27-33, 2016.

[10] A. Vij and A. Namboodiri, " Learning Minutiae Neighborhoods: A New Binary Representation for Matching Fingerprints," IEEE Conference on Computer Vision and Pattern Recognition Workshops 23-28 June 2014, Columbus, OH, USA DOI: 10.1109/CVPRW.2014.15

[11] G. A. Bahgat, A. A. Hefnawy, A. H. Khalil, N. S. A. Kader and S. Mashali, "Developed Fingerprint Segmentation Technique based on Mean and Variance Intensity Thresholding," Proceedings of the 4th IIAE International Conference on Industrial Application Engineering 2016, pp. 250- 257. DOI:10.12792/iciae2016.047.

[12] W. Abu-Ain, S. N. H. S. Abdullah, B. Bataineh, T. Abu-Ain and K. Omar, "Skeletonization Algorithm for Binary Images," Procedia Technology, vol. 11, pp. 704–709, 2013. doi: 10.1016/j.protcy.2013.12.248.

[13] M. S. Al-Ani, "A Novel Thinning Algorithm for Fingerprint Recognition," International Journal of Engineering Sciences, Vol. 2(2), pp. 43-48, 2013.

[14] R. F. L. Carneiro and J. A. Bessa, "Techniques of Binarization, Thinning and Feature Extraction Applied to a Fingerprint System", International Journal of Computer Applications, Vol. 103(10), pp. 1-8, 2014.

[15] V. Jain, A. K. Singh, "An approach for Minutia Extraction in Latent Fingerprint Matching," International Journal of Innovations in Engineering and Technology (IJIET), Vol. 6(1), pp. 51-58, 2015.

[16] D. Rutovitz, "Pattern recognition," Journal of the Royal Statistical Society, Vol. 129 (4), pp. 504–530, 1966.

[17] Atul S. Chaudhari, G. K. Patnaik, S. S. Patil, "Implementation of Minutiae Based Fingerprint Identification System using Crossing Number Concept," International Journal of Computer Trends and Technology (IJCTT), Vol. 8(4), pp. 178-183, 2014.

[18] I. K. Virdaus, A. Mallak, S.-W. Lee, G. Ha, and M. Kang, "Fingerprint Verification with Crossing Number Extraction and Orientation-Based Matching," Proceedings of The International Conference on Next Generation Computing,Bangkok, Thailand, February 2016, pp. 113-115.

[19] D. H. Douglas and T. K. Peucker "Algorithms for the reduction of the number of points required to represent a digitized line or its caricature," The Canadian Cartographer, vol. 10 (2), pp. 112-122, 1973.

[20] Z. Shukai, L. Zhengjiang, Z. Xianku, S. Guoyou and C. Yao, "A method for AIS track data compression based on Douglas-Peucker algorithm," Journal of Harbin Engineering University , 5, pp. 595-599, 2015.

[21] T. Tienaah, E. Stefanakis, D. Coleman, "Contextual Douglas-Peucker Simplification," Geomatica, Vol. 69(3), pp. 327-338, 2015. DOI: 10.5623/cig2015-306

[22] E.R. White, "Assessment of line-generalization algorithms using characteristic points," Am. Cartogr, vol. 12, pp. 17–28, 1985.

[23] J.S. Marino, "Identification of characteristic points along naturally occurring lines," Cartogr. Int. J. Geogr. Inf. Geovisualization, vol. 16, pp. 70–80, 1979.

[24] R. B. McMaster, "Automated Line Generalisation", Cartographica, vol. 24 (2) , pp.74 -111, 1987.

[25] S. Ri, "A new idea to construct the fractal interpolation function," Indagationes Mathematicae, Vol. 29(3), pp. 962-971, 2018. DOI : https://doi.org/10.1016/j.indag.2018.03.001

[26] J. Tang, Z. Feng, Y. Guo, "The Structure of Fractal Interpolation Curve in Plane," International Journal of Nonlinear Science, Vol. 26(1),pp.34-40, 2018.

[27] S. Ri, "A new nonlinear fractal interpolation function," Fractals, Vol. 25(6), pp. 1750063-1- 1750063-12, 2017 DOI: 10.1142/S0218348X17500633

[28] H. GUEDRI, J. MALEK and H. Belmabrouk, "Reconstruction of the human retinal blood vessels by fractal interpolation," Journal of Theoretical and Applied Information Technology (JATIT), Vol. 83(2), pp. 227-233, 2016.

[29] C. Chen, T. Lee, Y. M. Huang and F. Lai, "Extraction of Characteristic Points and Its Fractal Reconstruction for Terrain Profile Data," Chaos, Solitons Fractals Extraction of Characteristic Points and Its Fractal Reconstruction for Terrain Profile Data, vol. 39(4), pp. 1732-1743, 2009.

[30] M.F. Barnsley, "Fractal Functions and Interpolation," Constructive Approximation , vol. 2(1), pp. 303–332, 1986.

[31] M.F. Barnsley, Fractals Everywhere, 2nd edn. Academic Press, Inc., Boston , 1993.

[32] B. Mandelbrot, The Fractal Geometry of Nature, Freeman, San Francisco, 1982.

[33] S. B. Sujata, V. Khushbu, M. P. Jaymin and J. B. Subhash, "Finger Print Fractal Dimension as a Supplementary Quantitative Measure Distinguishing Fingerprints and Gender," Biostat Biometrics Open Acc J. vol.2(3), pp. 555-587, 2017. DOI: 10.19080/BBOAJ.2017.02.555587.

[34] C. Sahu, V. Jain, "A Novel Approach to Fractal Dimension based Fingerprint Recognition System," International Research Journal of Engineering and Technology (IRJET), vol. 3(4), pp.67-71 ,2016.

[35] K. Chen, S. Guo, Y. Lin and Z. Ying, "Least Absolute Relative Error Estimation," Journal of the American Statistical Association, vol. 105:491, pp. 1104-1112, 2010.

[36] W. Ford, "Floating Point Arithmetic," Numerical Linear Algebra with Applications Using MATLAB, pp. 145-162, Elsevier 2015.

[37] C. Chen, J. Twycross and J. M. Garibaldi, "A new accuracy measure based on bounded relative error for time series forecasting," PLoS ONE, vol. 12(3): e0174202, 2017. Doi:https://doi.org/10.1371 /journal.pone. 0174202

[38] J . Rougier, "Ensemble averaging and mean squared error," Journal of Climate, vol. 29(4), pp. 8865-8870, 2016. https://doi.org/10.1175/JCLI-D-16-0012.1

[39] E. Holst and P. Thyregod, "A statistical test for the mean squared error," Journal of Statistical Computation and Simulation, vol. 63(4), pp. 321-347, 1999. DOI: 10.1080/00949659908811960

[40] J-K. Liang, C. Cattani and W-Q. Song, "Power Load Prediction Based on Fractal Theory," Advances in Mathematical Physics, vol. 2015, Article ID 827238, 6 pages. http://dx.doi.org/10.1155/2015/827238.

[41] L. MAČĖNAITĖ , M. LANDAUSKAS and V. P. PEKARSKAS, "Surface Roughness Simulation Using Fractal Interpolation to the Profilogram," MATERIALS SCIENCE (MEDŽIAGOTYRA). Vol. 18(2), pp. 138-144, 2012. http://dx.doi.org/10.5755/j01.ms.18.2.1916.

# Analytical and Comparative Study of Different Types of Two-Leg Chopping up Regulator

Walid Emar[1]

Electrical Engineering Department
Isra University, Faculty of Engineering
Amman, 11622 Jordan

Omar A. Saraereh[2]

Electrical Engineering Department
The Hashemite University, Faculty of Engineering
Zerqa, 13133 Jordan

*Abstract*—The main focus of this article is to analyze and simulate the two-leg parallel connection of a chopping up regulator with flattering inductive smoothers or with an interphasing centre-tap transformer supplied by a three-phase diode rectifier and a DC link in between. The article deals with the problem of reducing total harmonic distortion, minimizing THD and EMI with low switching frequency. The Simulated three phase a.c. load model is added at the end to investigate the current and voltage harmonics. The main objective of this paper is the investigation of the problem and active impact of replacing flattering inductive smoothers used to reduce voltcurrent fluctuating waveforms of the chopping up regulator by new topology known as interphasing centre-tap transformer with magnetic coupling. The comparison of these two variations of the study is then done based on their technical parameters and economical investment viewpoint. The considered technical parameters are to be current distribution into individual legs, amount of voltcurrent ripple and area of discontinuous currents. The investment costs governed by the material requirements are essential for transformer and smoother production design. The outcome of using the interphasing centre-tap transformer is successive cancelation of voltcurrent fluctuating waveforms generated at the output of the chopping up regulator. This is proved by an experiment with 35 in input and power chopping up400-V/90-A.Software simulations in Simplorer and Matlab/Simulink or software program and regimen prototypes have been arranged to confirm the results.

*Keywords—Chopping up regulator with a flattering inductive smoother; magnetic coupling; connection with interphasing centre-tap transformer*

## I. INTRODUCTION

In many manufacturing applications it is necessary to change the dc direct voltcurrent source into a dc variable one. A chopping up regulator directly changes from dc to dc and like an ac transformer, it is used to step up or step down a dc voltage source value. Due to its capability to deliver smoothly adjustable dc voltage, dc chopper ups have revolutionized the regime industrial control devices and frequency inverters with unidirectional and with power levels extending from fractional horsepower up to several megawatts [1,2].

The chopping up regulator is required for the use in traction motor control applications in electric cars, marine vehicles, marine lifts and forklift. The chopping up regulator moreover can be used in braking regime of motors to return back to the source energy which results in energy savings for transport systems with regular stops. In addition, the Chopper

ups are used together with a capacitive filter respectively with an inductive smoother to generate dc voltage respectively direct current sources.

Regulated inverters are widely used these days to control the voltage at the output of solar systems. Another form for stepping up the voltage would be a PWM inverter followed by a chopping up regulator [3].

Due to the periodic chopper up function, its current waveforms and its output voltage contain certain harmonics [4]. This type of harmonics has an impact on load performance, chopper up losses and rises the range of discontinuous (interrupted) currents in which the chopper up can operate, and as a result, its operation is more complex and the control characteristics will in many cases be changed which requires additional control system unitsand special controllers.

In practice, inductive smoothers are used to lower the harmonic level, but the magnitude of these inductive smoothers increases with increasing input and output curve fluctuations [5, 6].

To improve the performance of chopping up regulators in reducing variable waves generated at its input or output, the chopping up regulator structure is usually improved by adding two or more legs to act in parallel. The outputs of these individual legs are combined by inductive smoothers without magnetic coupling. Generally, the switching technique of these legs introduces a phase displacement of the control signals between the individual legs and at the same time keeps them working with the same switching frequency. This type of connection is referred to as a several-leg parallel connection of a chopping up regulator with flattering inductive smoothers [7-9].

A better solution is attained if a magnetic coupling of each two flattering inductive smoothers is introduced. This new topology of flattering inductive smoother is recognized as an interphasing centre-tap transformer and the connection of chopping up regulators in this case is known as a several-leg parallel connection of chopping up regulators with an interphasing centre-tap transformer. In this case the number of legs must be even.

The result is a successive cancellation of the voltcurrent curves fluctuations with little design demands on the smoothers. It also makes the control of the regulator easier

since there will be no need for additional controllers to equally divide the currents into the individual legs. The magnetic coupling between each two smoothers takes care of that [10-12].

In this paper, an analysis of an interleaved chopping up regulator supplied from an ac network via a single phase bridge rectifier circuit is carried out. An interleaved converter with two-leg converters connected in parallel was considered for this work. The performance of the conventional chopping up regulator, two-leg interleaved chopping up regulator with inductive smoothers or with an interphasing centre-taped transformer were compared by means of simulation and design. Simulation of the system was carried out using MATLAB/Simulink. A design example is presented to illustrate the novel design concept. The formatter will need to create these components, incorporating the applicable criteria that follow.

## II. BASIC CONFIGURATION OF CHOPPING UP REGULATOR

Fig. 1 shows the topology of a fundamental connection of a chopping up regulator. The main parts are an a.c. power supply via a single phase bridge rectifier, input capacitive filter (capacitance $C_f$ with inductance $L_f$) flattering inductive smoother, L, power electronic switch (BJT, IGBT, MOSFET or THYRISTOR), diode, and output capacitive filter. The load is considered to be a resistor of resistance $R_a$ with an smoother of inductance $L_a$.

The principle of operation of an elementary chopping up regulator shown in Fig. 1(a) is currently known. A chopping up regulator varies its average output voltage $V_o$ that appears across the load which is relative to its input $V_s$, by varying the proportion of its operating time during which the output is connected to the input. In other words, the unspecified switching device S operates with a regular periodic time $T$, and is closed for a time $t_{on} = T - t_{off}$.

### A. Continuoius Operating Regime

The regulator simulation parameters for the current continuous regime are listed in Table I as:

When the switch is closed, the diode is reverse biased. Kirchhoff's voltage law around the path containing the source, smoother L, and closed switch S is:

$$V_s = L \frac{di_L}{dt} = v_L \qquad (1)$$

The peak to peak current ripple of the smoother current is then computed from:

$$\Delta i_L = \frac{V_s}{L} t_{on} = \frac{V_s}{fL} k, \quad k = \frac{t_{on}}{T} \qquad (2)$$

k is a duty time ratio of the chopping up regulator. When the switch is opened, due to the energy stored in the smoother L, the source current, cannot change instantaneously, so the diode becomes forward-biased and conducts the source current. Assuming that the load voltage $V_o$ is

$$V_s - V_o = L \frac{di_L}{dt} = v_L \implies \Delta i_L = \frac{V_o - V_s}{fL} t_{off} = \frac{V_o - V_s}{fL} (1 - k) \qquad (3)$$

TABLE I.    TECHNICAL PARAMETERS OF THE SIMULATED REGULATOR

| Parameter | Symbol | Real value |
|---|---|---|
| Smoothing filter: | $L$ | 30mH |
| Regulator operating period: | $T$ | 20ms |
| Supply voltage | $V_s$ | 35V |
| Load resistance: | $R_a$ | 5Ω |
| Load inductance: | $L_a$ | 2mH |
| Output capacitive filter | $C$ | 5mF |

Where $\Delta i_L$ is the peak to peak current ripple in the smoother current, when the switch is open. Under steady-state operation conditions, the net change in smoother voltage must be zero. Using Eqs. (2) and (3), it yields [2]:

$$V_o = \frac{V_s}{1-k} \implies k = 1 - \frac{V_s}{V_o} \qquad (4)$$

Substituting from Eq. 4 into Eq. 2 or Eq. 3, results in the following value for the smoother current ripple:

$$\Delta i_L = \frac{V_o - V_s}{fL} (1 - k) \qquad (5)$$

By expressing the operating period of the converter as $= t_{on} + t_{off}$, then from Eqs. 2 and 3, we may describe the current ripple of the smoother $i_L$ or the source current $i_s$ as a function of the output load voltage as follows:

$$\Delta i_L = \frac{V_s (V_o - V_s)}{f L V_o} \qquad (6)$$

The maximum value of $\Delta i_L$ as a function of the output load voltage $V_o$ occurs when $V_o$

$$\Delta i_{Lmax} = \frac{V_o}{fL} \qquad (7)$$

### B. Discontinuous Operating Regime

The chopping up regulator may also operate in discontinuous operating regime due to small values of the smoother inductance ($L=1mH$) or the use of low frequency switches [16].

In some cases, the discontinuous smoother current, is desirable for control reasons when the output is regulated. The smoother current and load voltage ripples are determined from the fact that the average smoother voltage is zero.

The source and diode currents for discontinuous current regime have the basic waveforms as shown in Fig. 1(c). When the main switch, S, is on, the voltage across smoother, L, is $V_s$. But when the switch is off and the smoother current, is decreasing, the smoother voltage is $V_s$- $V_o$. The smoother current, falls down until it touches the zero-axis before the end of the switch operating period, T. With the switch and diode are off, the smoother current, is zero. The average voltage across the smoother is:

$$V_s k + (V_s - V_o)k_1 = 0 \implies V_o = V_s(\frac{k}{k_1} + 1) \qquad (8)$$

Where $k_1 = t_k/T$, $t_k$ is the instant of interruption at which the smoother current reaches its zero value when the diode becomes off. The maximum value of the discontinuous smoother current is the same as the peak to peak current ripple in this current, when the switch is closed:

$$I_{Lmax} = \frac{V_s}{fL}\left(1 - \frac{V_s}{V_o}\right) \tag{9}$$



(a) Basic Chopping up Regulator Supplied from an ac Network via a Rectifier Circuit.



(b) Continuous Voltage and Current Waveforms in the basic Chopping up Regulator.



Fig. 1. (c) Discontinuous Voltage and Current Waveforms in the basic Chopping up Regulator.

The average value of the regulator output current, $i_o$, is the same as the average current in the load, thus:

$$I_{oavg} = \frac{V_o}{R} = \frac{1}{2}I_{Lmax}\, k_1 = \frac{1}{2}\frac{V_s\, k}{fL}k_1 \tag{10}$$

After solving, it results in:

$$k_1 = \frac{V_o}{V_s}\frac{2\,f\,L}{R\,k} \tag{11}$$

Substituting for $k_1$ into Eq. 8, results in the following value for the average load voltage, $V_o$ within the discontinuous current regime:

$$V_o = \left[\frac{1}{2}\left(1 + \sqrt{1 + \frac{2k^2R}{f\,L}}\right)\right]V_s \tag{12}$$

The boundary operating conditions between continuous and discontinuous current regimes occurs when the interruption and zero value of the smoother current occurs exactly at $t_k = T$. The average value of the smoother current on the limits of discontinuous current regime may be determined as follows:

$$\boldsymbol{I_{flim}} = \frac{V_s}{2fL}\left(1 - \frac{V_s}{V_o}\right) \tag{13}$$

The maximum values of limits (13) occur when$V_o$

$$I_{flimmax} = \frac{V_s}{2fL} \tag{14}$$

The average value of discontinuous smoother current within the region of discontinuity is obtained in [1] as:

$$I_{Ldisav} = \frac{k^2 V_o \left(\frac{V_o}{V_s}-1\right)}{2fL} = \frac{\left(1-\frac{V_s}{V_o}\right)^2 \left(\frac{V_o^2}{V_s}-1\right)}{2fL} \tag{15}$$

From the above mentioned expression 13, it may be determined the smoother current curve at the boundary of discontinuous currents shown in Fig. 5.

- Advantages

The fundamental chopping up regulator is a low cost regulator with a simple topology that can be easily regulated and is able to provide high gains. The gating signals of the switch can be generated with well-developed integrated circuits or microprocessors. In order for the fundamental chopping up regulator to work properly it must get a smooth input current. High efficiency can be achieved with medium and low switching times.

- Disadvantages

The circuit components are not ideal (voltage drops on power electronic devices, capacitor series resistances, smoother resistances, switching losses) and the output is limited. The off time of the switches should be limited to below 0.90 or 0.95 in order to avoid short circuiting the switch and the dc source. Also as the switching time gets closer to unity, the output voltage becomes more sensitive to any changes in this time. This can make it to be more difficult to adjust the output voltage at higher gains. The switch must also have a high voltage rating due to having the output voltage across it when it is off. This sometimes requires a switch that will have a slower switching time, or have a higher forward resistance.

## III. Conventional Two-Leg Chopping up Regulator with Magnetically Uncoupled Inductive Smoother

Fig. 2 depicts a circuit diagram for a conventional two-leg connection of chopping up regulator with magnetically uncoupled inductive smoothers $L_1$, $L_2$. Such a chopper up has two legs connected in parallel, $S_1$, $D_1$ and $S_2$, $D_2$. They are supplied via a single phase bridge diode rectifier and are switched on non-simultaneously with the displacement time T / n. The load comprises a smoother $L_a$ and resistor $R_a$.

Inductive smoothers $L_1$ and $L_2$are passive electrical elements (known as magnetically uncoupled filters) wounded from a wire coil around a magnetic core of steel to confine its magnetic flux within the total number of coil turns. They are designed to create magnetic fields in the core as a result of flowing current and induced voltages across their terminals. This leads to a considerable stronger magnetic induction than would be produced by a simple wire coil without a magnetic core.



Fig. 2. Two-Leg Chopping up Regulator with Magnetically Uncoupledinductive Smoothers.

### A. Peak to Peak Current Ripple

Traditionally, the several-leg chopping up regulator with magnetically uncoupled inductive smoother, as presented in Fig. 2, highly reduces the total current fluctuating waveforms flowing into the output capacitors and significantly the power increase as compared to the basic structure of the chopping up regulator. These results are listed for identical and linear design conditions of the inductive smoothers while $L_1 = L_2 = L$. Therefore, the averaged currents into both legs are divided into their inductive smoothers, $L_1$, $L_2$ equally by soft control of every leg main switch with a time displacement of T/2.

### B. Continuous Operating Regime

Under the condition of a continuous conduction regime of operation, the two-leg regulator is simulated for the same parameters mentioned in Table I, whereas both smoothers have same inductance values, $L_1 = L_2 = L = 30mH$. The voltage and current waveforms of the system for $0 \leq k \leq 1/2$are shown in Fig. 3.

The second area of economic interest where consumption increases are the loss of electricity due to transmission and distribution of electricity to end users. Technical losses that are not caused by human causes can be divided into losses in the lines and losses in voltage transformation. It can be seen from formula (1) that line losses are directly proportional to the quadrant of the maximum current, so the effort to control consumption is therefore to limit the peak of the load during the day, especially morning and evening [17-18].

Considering the switching duty ratio as $= t_{on}/T$ , then the ripple $\Delta i_L$ of the phase current flowing into individual smoothers should be same whether the switch is on or off, and it is obtained as:

$$T = t_{on} + t_{off} = \frac{L\,\Delta i_L}{V_s} + \frac{L\,\Delta i_L}{V_o - V_s} \Rightarrow \Delta i_L = \frac{V_s\,(V_o - V_s)}{f\,L\,V_o} =$$
$$\frac{k\,(1-k)}{f\,L}V_o = \frac{k\,V_s}{f\,L} \tag{16a}$$

Fig. 3. Current and Voltage Steady State Waveforms of Two-Leg Connection of Chopping up Regulator with Flattering Inductive Smoothers for $0 \leq k \leq 1/2$.

The ripple of the source current $i_s$, may be determined for $0 \leq k \leq 1/2$ from the sharp growth of this current during on-regime of one switch and off-regime of other switches. Thus, when switch $S_1$ is on:

$$\Delta i_s = \frac{V_s}{L} t_{on} - \frac{V_o - V_s}{L} t_{on} = \frac{2 V_s - V_o}{f L} \left(1 - \frac{V_s}{V_o}\right) = \frac{V_o}{f L} (1 - 2k)k = \frac{V_s}{f L} \frac{(1 - 2k)k}{1 - k} \qquad (17)$$

The ripple of the source current in region $\frac{1}{2} \leq k < 1$ is obtained using similar idea as in [9]:

$$\Delta i_s = 2 \frac{V_o - V_s}{L} \left(\frac{T}{2} - t_{on}\right) = 2 \frac{V_o - V_s}{f L} \left(\frac{1}{2} - \frac{V_s}{V_o}\right) = \frac{V_o}{f L} (3k - 2k^2 - 1) = \frac{V_s}{f L} \frac{(3k - 2k^2 - 1)}{1 - k} \qquad (18)$$

The maximum value of the ripple of the source current occurs in region $0 \leq k \leq 1/2$ when $k = 1/4$ or when $V_o = \frac{4V_s}{3}$ and in region $\frac{1}{2} \leq k < 1$ when $k = 3/4$ or when $V_o = 4V_s$, thus:

$$\Delta i_{smax} = \frac{V_o}{8 f L} = \frac{V_s}{6 f L} \qquad (19)$$

The boundary between continuous and discontinuous current regime occurs as soon as the minimum value of the smoother current touches the zero axis exactly at $k_1 = 1 - k$. Therefore, the average value of the phase smoother current flowing into each leg of the chopper up at the boundary of discontinuity is given exactly as in Eq. 13.

### C. Discontinuous Operating Regime

In the previous analysis the switching frequency and the flattering smoothers of the regulator were considered to be sufficiently large and therefore, all currents in the circuit have normal continuous waveform. Otherwise, the source, phase and load currents are discontinuous as shown in Fig. 4.



Fig. 4. Discontinuousphase Smoother Current Waveforms for $0 \leq k \leq 1/2$.

Fig. 5. Source Characteristics of the Regulator with the Border Curves in the Interrupted Region.

The average value of each smoother current within the region of discontinuity of phase currents has the same value as that in Eq.15 which is obtained for the smoother current of fundamental connection. Therefore, the average value of the source current of the two-leg regulator within the region of discontinuity of phase currents is given as:

$$I_{sav} = 2I_{Ldisav} = \frac{k^2 V_o (\frac{V_o}{V_s}-1)}{fL} = \frac{\left(1-\frac{V_s}{V_o}\right)^2 (\frac{V_o^2}{V_s}-1)}{fL} \qquad (20a)$$

Eq. 20a determines the input current characteristics in the area of discontinuous phase and source currents shown in Fig. 5. If the source current is operating at the limits of discontinuity which may occur if and only if $k_1 = 1/2$, then the average value of such current at the boundary of the discontinuous source current region is determined as follows:

$$I_{slim} = \frac{i_L(t_{on})}{2} = \frac{V_s (V_o - V_s)}{4 f L V_o} \qquad (21)$$

Where $i_L(t_{on})$ is the value of smoother phase current, $i_L$ at time $t = t_{on}$. The boundary between the continuous and discontinuous source current regions determined by relationship (21) is indicated in Fig. 5.

$$\Rightarrow I_{slimmax} = \frac{V_s}{8 f L} = \frac{V_o}{16 f L} \qquad (22)$$

Eqs. 21 and 22 are used to determine the input characteristic curves that describe the behavior of the currents within the region of discontinuity as shown in Fig. 5.

The peak to peak current ripples of the basic and two-leg chopping up regulators are shown as a function of the duty ratio, k in Fig. 6. Eqs. 6 and 19 have a good practical meaning because they are used to make an approximate calculation of the source current ripple slightly, which helps in determining the desired inductive smoother design values by selecting the appropriate ripple in it. From Fig. 6 It can be seen that this fluctuating ripple decreases with the number of legs of the regulator.

It should be noted that the number of inductive smoothers for n-leg connection is n smoothers, but compared to the simple regulator each smoother is adjusted to a value of 1/n of

the source current. This advantage results in the use of smaller inductive smoothers on the output of the regulator in order to meet the same requirements for the total current ripple. Therefore, the distortion in the output voltage can also be improved due to smaller output coils, resulting in using a smaller capacitance at the output [1-2].



(a) Plot of Current Ripple of Source and Phase Current Versus Duty Time Ratio, *k* as it Varies from Zero to Unity.



(b) Plot of Current Ripple of Source and Phase Current Versus Duty Time Ratio, *k* with Additional Inductance at the Input Source $V_s$.



Fig. 6. (c) Plot of Source Current Curves Versus Source Voltage $V_s$ within the Region of Interrupted Currents.

Inserting an additional inductance, $L_s$ in series with the source $V_s$, will result into the following expressions and graphs for the peak-to-peak ripple of the source and phase currents for $0 \le k \le 1/2$ as follows:

$$\Delta i_L = \frac{V_s\,(L+L_s) - L\,V_o}{L\,(2\,L_s + L)}\,t_{on} = \frac{V_s\,(L+L_s) - L\,V_o}{f\,L\,(2\,L_s + L)}\,k \qquad (16b)$$

Similar expressions may be obtained for $1/2 \le k \le 1$. Fig. 6(b) illustrates that the peak-to-peak smoother current ripple reaches its maximum value beyond its operating region at $= \frac{L_s + L}{2L} > 1/2$, therefore we consider its maximum value obtained at $k = 1/2$ as shown in the figure. Fig. 6(c) presents the waveforms of the source and phase interrupted currents after adding an inductance at the input:

$$I_{sav} = 2I_{Ldisav} = \frac{V_s(\frac{L_s+L}{L} - k)k}{f(2L_s + L)} \qquad (20b)$$

## IV. TWO-LEG REGULATORS WITH INTERPHASING CENTRE-TAP TRANSFORMER

Finally, complete content and organizational editing before formatting. Please take note of the following items when proofreading spelling and grammar:

In recent years, most power and regulation researches have been focused on the use of several-leg chopping up regulators to improve power electronic processors performance and power factor for personal computers, office equipment, space systems, laptops and telecommunication equipment as well as motor drive and power systems.

In addition, demand for electronic power processors for renewable energy systems has increased, since the improvement and growing production of industrial and commercial energy products and technologies and the growing dilemma of fossil fuels as a source of primary energy sources have expanded [10-14].

In general, the basic structure of the chopping up regulator topology with a single inductive smoother gives an acceptable performance in some of the aforementioned applications, but in some others such as distributed power conversion systems and power factor correction circuits, the performance of this configuration is enhanced by adding one or more legs with magnetically coupled smoothers that are operated in parallel or in series.

A connected coil is a device that is primarily used to store energy during the regulator switching duty cycle. The power entering the coupled coil is not the same as the power leaving it in a given instant. Coupled flattering inductive smoother are used to reduce regulator volume by using one core instead of two or more, to improve regulation of power regulators [2-5].

The aim of the work is to improve the performance of chopping up regulator by means of several-leg connection with directly coupled flattering inductive smoothers known as interphasing centre-tap transformer [2, 3, 12].

### A. Peak to Peak Ripple

Fig. 7 shows the schematic diagram of the two-leg chopping up regulator with directly coupled flattering inductive smoothers. The coupled flattering inductive smoothers $L_1$ and $L_2$ share the same winding directions.

This type of regulators analysis can be best explained under a condition of a very small value of leakage inductances. Then, the windings of the two inductive smoothers have identical numbers of turns $N$. The inductance of each coupled flattering smoother is divided into two main parts [13-16]:

$$\begin{aligned} L_1 &= L_{r1} + L_m \\ L_2 &= L_{r2} + L_m \\ L_m &= \mu\sqrt{L_1\,L_2} \end{aligned} \qquad (23)$$

Where $L_1$, $L_2$ are the self inductances of the flattering smoothers, $\mu$ is the magnetic coupling coefficient, $L_{r1}$, $L_{r2}$ are leakage inductances of the two flattering inductive smoother, $L_m$ is mutual inductance. In order to simplify the analysis, let's consider the inductive smoothers to have the same inductance values, $L_1 = L_2 = L$ and $L_{r1} = L_{r2} = L_r$.

The aim of the work is to improve the performance of several-leg regulators by means of a several-leg parallel connection with an interphasing centre-tap transformer [2, 3, 6].

Concerning the switching technique of such regulator, when switch $S_1$ is on and switch $S_2$ is off, $i_1$ flows into $S_1$ and a similar large current $i_2$ flows into $D_2$ because of the magnetic coupling, and each leg carries approximately half the source current continuously and substantially equal voltages across the two halves of the winding. On the other side, when both switches are off, both currents will decrease exponentially into the two individual legs of the regulator. Therefore, there will be no need for controllers to distribute the currents equally into both individual legs as in the case of connection with smoothing inductive smoothers without magnetic coupling [1, 14].



Fig. 7. Two-Leg Chopping up Regulator with Directly Coupled Flattering Inductive Smoothers.

## B. Continuous Operating Regime

The regulator simulation parameters for the continuous current regime are listed in Table II.

TABLE II.    TECHNICAL PARAMETERS OF THE SIMULATED REGULATOR.

| Parameter | Symbol | Real value |
|---|---|---|
| Transformer winding inductances: | $L_1=L_2$ | 30mH |
| Smoothing filter: | $L$ | 30mH |
| Regulator operating period: | $T$ | 20ms |
| Supply voltage | $V_s$ | 35V |
| Load resistance: | $R_a$ | 5Ω |
| Load inductance: | $L_a$ | 2mH |
| Output capacitive filter | $C$ | 5mF |
| Magnetic coupling | $k$ | 0.5 |

Thus, assuming ideal devices with zero voltage drop and that the voltages across the individual legs of the regulatorare $V_1$ and $V_2$, respectively, then it can be concluded that during the on-regime of switch $S_1$ and off-regime of switch $S_2$:

$$V_1 = 0 = V_s - L\frac{di_{s/2}}{dt} = V_s - L_r\frac{di_1}{dt} - L_m\frac{di_m}{dt}$$

$$V_2 = V_s + L\frac{di_{s/2}}{dt} = V_s + L_r\frac{di_2}{dt} - L_m\frac{di_m}{dt} \qquad (24)$$

$$\frac{i_s}{2} = i_1 + i_m = i_2 - i_m$$

Identical expressions may be derived for the on-regime of $S_2$ and off-regime of $S_1$.

When both switches, $S_1$ and $S_2$ are off simultaneously, the regulator diodes $D_1$ and $D_2$ conduct the phase currents to the load, and therefore:

$$V_1 = V_2 = V_s - L\frac{di_{s/2}}{dt} = V_s - L_r\frac{di_1}{dt} - L_m\frac{di_m}{dt} = V_s - L_r\frac{di_2}{dt} + L_m\frac{di_m}{dt} \qquad (25)$$

The voltage and current waveforms are illustrated for the regime of non-simultaneous conduction of switches $S_1$ and $S_2$ in Fig. 8 for $0 \le k \le 1/2$. During on-state of $S_1$, it can be concluded that:

$$\Delta i_s = \frac{2V_s}{L}t_{on} \qquad (26)$$

And under the condition of large value of output capacitive smoother, then during the off-state of $S_1$:

$$\Delta i_s = \frac{2(V_o - V_s)}{L}t_{off} \qquad (27)$$

To avoid any damage that may occur due to the short circuiting of the supply to the chopper up when the main switch is on and to fully damp the fluctuations that appear in the source current ripple, it is suggested to add a flattering smoother $L$ into the circuit in series with the main input supply, $V_s$ [2, 15]. Therefore, under steady state conditions, the source current ripple $\Delta i_s$ will be determined from the steep rise respective from the steep fall of the source current during the on-time respective the off-time of switch $S_1$. This

fluctuating waveform for the regime of non-simultaneous conduction of switches in region $(0 \le k \le 1/2)$ is obtained as follows:

$$T = t_{on} + t_{off} \Rightarrow \Delta i_s = \frac{V_s}{2fL}\left(\frac{k-2k^2}{1-k}\right) = \frac{V_o}{2fL}(k - 2k^2) \qquad (28)$$

For the regime of simultaneous conduction of switches in region $(1/2 \le k \le 1)$, the source current ripple $\Delta i_s$ is given as follows:

$$\Delta i_s = \frac{V_s}{2fL}\left(\frac{3k-2k^2-1}{1-k}\right) = \frac{V_o}{2fL}(3k - 2k^2 - 1) \qquad (29)$$



Fig. 8.   Current and Voltage Steady State Waveforms of Two-Leg Chopping up Regulator with an Interphasing Centre-Tap Transformer for $0 \le k \le 1/2$.

The maximum value of $\Delta i_s$ occurs in eq. 27at $k=1/4$and in eq.28at$k=3/4$. Thus:

$$\Delta i_{smax} = \frac{V_o}{16fL} \tag{30}$$

Since the current rippling each phase is during the on-state of its active switches equal to the ripple in the off-state of this switch, then the maximum value of the magnetization current of the transformer occurs when $k = \frac{1}{2}$ as follows:

$$I_{\mu max} = \frac{V_o}{4 f L_\mu} \tag{31}$$

Where $L_\mu$ is the magnetization mutual inductance between both smoothers of the transformer.

### C. Discontinuous Operating Regime

In the previous analysis the smoother inductance was considered to be very large and due to the use of high frequency devices in the regulator, all currents in the circuit have normal continuous waveform and they do not even touch the zero-axis during the regulator operating period. But if the smoother has a small inductance and the switching frequency of the regulator is low, the source, the phase and the load currents are not sustained throughout the regulator operating period and they become zero before the end of this period as shown in Fig. 9.

The value of the inductance of smoothing smoothers is 20-times smaller than that in the regime of continuous currents. Therefore, during the on-time of switch $S_1$, the discontinuous phase current is expressed as follows:

$$i_1(t) = \frac{V_s}{L}t \tag{32}$$

At $t = t_{on}$ the phase current is transferred from switch $S_1$ to its freewheeling diode $D_1$ which means that:

$$i_1(t) = i_1(t_{on}) - \frac{V_o - V_s}{L}(t - t_{on}) = \frac{V_s}{L}t_{on} - \frac{V_o - V_s}{L}(t - t_{on}) \tag{33}$$

This conduction interval of $S_1$ ends at the moment of interruption, $t = t_k$, when the phase current is distinguished.

This moment would be obtained from Eq. 33 by replacing $t$ with $t_k$ and putting$i_L(t_k) = 0$:

$$t_k = \frac{V_o}{V_o - V_s}t_{on} \tag{34}$$

Similar expressions could be obtained for all other phase currents of the regulator. The average value of the phase currents within the region of discontinuity is then obtained using Fig. 10(a) as:

$$I_{av1} = \frac{i_1(t_{on})t_x}{2T} = \frac{V_oV_s}{(V_o - V_s)fL}k^2 \tag{35}$$

The average value of the source current is then obtained as the sum of the average values of n phase currents, thus:

$$I_{avs} = nI_{av1} = n\frac{V_oV_s}{(V_o - V_s)fL}k^2 \tag{36}$$

Eq. 35 and 36 could be used to determine the input characteristics of the regulator that describes the behavior of the phase and source currents within the region of discontinuous operating regime. The variable parameter of such characteristics is the duty ratio, k, whereas in this operating regime$V_s \neq (1 - k)V_o$.

The boundary conditions for the phase currents at the limits of discontinuous operating regimes may be derived using the notations of Fig. 11 as:

$$I_{flim} = \frac{V_o}{2fL}(1 - k)k \tag{37}$$



Fig. 9.   Discontinuous Current Waveforms with an Interphasing Transformer for $0 \leq k \leq 1/2$.

Hence, the regulator operates at the limits of discontinuous and continuous region of currents when the lower value of the smoother currents hits zero-axis exactly at the end of its conduction period. Using Eq. 37, the average value of the source current at the limits of discontinuity of phase currents is:

$$I_{slim} = n\frac{V_o}{2fL}(1-k)k \tag{38}$$

However, if the discontinuous waveforms of phase currents overlap each other, as shown in Fig. 11, then the source current, $i_s$ would also have a discontinuous waveform. Therefore, the source current would touch the limits of interruption exactly at $t_k = \frac{T}{2}$. Substitution this value for $t_k$ into Eq. 34, we get for the duty ratio the following value:

$$k = \frac{V_o - V_s}{nV_o} \tag{39}$$

The boundary condition for the source current at the limits of discontinuity is then given as:

$$I^*_{slim} = \frac{\Delta i_1}{2} = \frac{kV_s}{2fL} \tag{40}$$

Substituting from Eq. 39 into Eq. 40 yields:

$$I^*_{slim} = \frac{V_s(V_o - V_s)}{2nfLV_o} \tag{41}$$

Similar expressions may be obtained for the source current in the conduction region of $0 \le k \le 1/2$. The maximum unrestricted source current limit value is $I^*_{slimmax} = \frac{V_o}{32fL}$. The source characteristic curves of the regulator with an interphasing transformer are illustrated in Fig. 10.

Adding an inductance, $L_s$, in series with the input source, $V_s$, results into the following expressions and waveforms for the source and smoother currents ripple at $0 \le k \le 1/2$:

$$\Delta i_s = \frac{2V_s}{fL_m}k + \frac{V_s}{2fL_s}k(1-2k) \tag{42}$$



(a) Source characteristic curves with their borders in the region of interrupted currents.



(b) Peak-to-peak ripple of source current with additional inductance at the input and different ratios: a) $L_s < 4L_m$, b) $L_s > 4L_m$, c) $L_s = 4L_m$



Fig. 10. Source Characteristic Curves with their Borders in the Region of Interrupted Currents for Real Values of $L_s$ and $L_m$.

Maximum value for this ripple is obtained at $= \frac{4L_s + L_m}{4L_m}$. Since this switching ratio of maximum value is out of the first operating region of the converter, $0 \leq k \leq 1/2$, the real maximum ripple is considered to be reached at $k= ½$. Similar Expressions may be obtained for $1/2 \leq k \leq 1$:

$$\Delta i_s = \frac{2V_s}{f\,L_m}(1-k) + \frac{V_s}{2fL_s}(3k - 2k^2 - 1) \qquad (43)$$

Illustration of $\Delta i_s$ versus switching ratio $k$ for different values of the magnetization inductance $L_m$ and source inductance $L_s$ is shown below in Fig. 10(b). The average value of the source current at the limits of discontinuity within the operating region $0 \leq k \leq 1/2$ is then given as $\Delta i_s/2$. Fig. 10(c) presents the waveforms of the source and phase interrupted currents after adding an inductance at the input.

## V. DEMANDS TO MATERIAL OF FLATTERING INDUCTIVE SMOOTHER AND TRANSFORMER

Next issue will be to find a simple method for determining demands to materials of magnetically coupled inductive smoother to compare them with other connections of a chopping up regulator with a smoother. Taking into consideration a coil with a number of turns *N,* wrapped around a rectangular ferromagnetic core with a cross-sectional area *A* and current *i* flowing into it as shown in Fig. 11. Then, under the condition of linear core saturation curve and a neglected leakage flux, we can write for the inductance *L* of the coil, the possible maximum value of the non-saturation magnetic flux ($\phi$) density *B* within the ferromagnetic core and the corresponding maximum value of the peak magnetization current $I_{\mu max}$ that may flow within the magnetization inductance the following expression:

$$\emptyset = N S_{Fe} B = L I_{\mu max} \qquad (44)$$

Concerning the coil, it is necessary to design its winding with respect to the maximum rms value of the smoother current, $I_{rms}$, that can flow through it. After rearrangement we can get [1]:

$$LI_{\mu max} I_{rms} = N B S_{Fe} I_{rms} = N B S_{Fe} \rho S_{wire}$$
$$= \rho B S_{Fe} S_{cu} = \Sigma_{material} \qquad (45)$$

Wherein $\rho$ is the admissible density of a current flowing through a cylindrical coil wire of cross-sectional area $S_{wire}$, and $S_{cu}$ is the cross-sectional area of the whole winding of the coil with N turns and it may represent the material used in making the winding and $S_{Fe}$ represents the material required for the design of the core [2].

The expression for demands to the materials of a simple inductive smoother for the fundamental connection (Fig. 1) may be simplified, if we put $I_L = I_{rms} = I_{\mu max}$ where $I_L$ is the average value of the source current at the boundary of discontinuity. And thus:

$$\Sigma_{material} = L I_L^2 \qquad (46)$$

For the two leg connection without magnetic coupling we may get that:

$$\Sigma_{material} = L\frac{I_f^2}{4} \qquad (47)$$



Fig. 11. Magnetic Core with a Coil.

Where $I_f = I_{\mu max} = I_{rms}$ is the average current passing through each leg of the two-leg chopping up regulator.

Concerning the connection with an interphasing transformer, it may be concluded from the waveforms of the magnetization current, $i_m$ and magnetization voltage, $v_m$ in Fig. 9, that:

$$I_{\mu max} - I_{\mu min} = 2I_{\mu max} = \frac{v_m}{L_m} t_{on}$$
$$I_{\mu max} = \frac{v_m}{fL_m}\frac{k}{2} = \frac{V_o}{2fL_m}k = \frac{V_s}{2(1-D)fL_m}k \qquad (48)$$

According to the results of the simulation shown, for example, in Fig. 9 and 10, the maximum value of the magnetization current, $i_m$ occurs when $k=1/2$, therefore:

$$I_{\mu max} = \frac{V_s}{2fL_m} \qquad (49)$$

Thus

$$L_m = \frac{V_s}{2f I_{\mu max}} \qquad (50)$$

Substituting from Eq. 47 into Eq. 50 yields:

$$\frac{L_m}{L} = \frac{V_s I_{rms}}{2f\Sigma_{material}} \qquad (51)$$

Eq. 51 represents the percentage of the magnetization mutual inductance within the overall inductance of each smoother which is used as the appropriate design of the inductive smoother and their common core.

For the design of the inductance of the uncoupled inductive smoother of a fundamental connection of boost regulator shown in Fig. 1 we may write:

$$L = \frac{V_o}{4f\Delta i_{Lmax}} \qquad (52)$$

The value of the inductance of one smoother in the case of two-leg chopping up regulator without magnetic coupling is given as:

$$L = \frac{V_o}{8f\Delta i_{smax}} \qquad (53)$$

Concerning the chopping up topology with an interphasing transformer, the inductance of the smoother connected directly to the output load is given as:

$$L = \frac{V_o}{16\,f\Delta i_{smax}} \qquad (54)$$

The selection of the appropriate value for the inductance of a certain smoother should satisfy the critical inductance condition which happens when the current through the smoother decays to zero just prior to the next on time of the chopping up regulator switch.

## VI. Control Techniques of Boost Converter without Transformer

As it has been said, the distribution of currents into phases of two-leg chopper up regulator without magnetic coupling is achieved using suitable control techniques to proportionally adjust the switching ratio of each leg. Concerning the two-leg connection with magnetic coupling the interphase transformer takes care of distributing the currents equally into phases, and therefore there will be no need for controllers for this purpose [9].

The output voltage of the regulator is compared within a comparator with a reference voltage to produce an error signal. This error signal is then used to adjust the switching ratio of the regulator via a pulse modulator. An amplifier is then used to amplify the modulator output just to bring it to the required voltage and power level suitable for driving the chopper regulator [2].

The output voltage of the regulator can be changed by controlling the switching ratio $k$, and this may be using achieved Constant frequency method: The regulator, or switching, frequency $f$ (or chopping period $T$) is kept constant and the on-time $T_1$ is varied. The width of the pulse is varied and this type of control is known as pulse width modulation (PWM) control [1]. The PWM control circuit may be composed of a square wave oscillator, a flip-flop, an integrator and two gates.

The output voltage of the chopper regulator can be changed by controlling the switching ratio $k$, and this can be achieved by using the constant frequency method, where the switching frequency f of the regulator is kept constant and its on-time $t_{on}$ is varied. The pulse width of the triggering signal for each switch is therefore distinct based on the level of the required output voltage. This type of control technique is known as pulse width modulation control technique (PWM). The PWM control circuit usually comprises a square wave oscillator, flip-flop logic circuit, integrator, and two gates.

The frequency of the master square wave oscillator is chosen to be twice that of the regulator operating frequency. The triangular wave is generated by the help of an integrator which is the compared with an error signal to produce the pulse width modulation [13-14].

The current distribution equally into both legs of the regulator would be in accord with their respective forward voltage drops. Therefore, using one controller for both legs does not guarantee a same voltage drop across each leg and consequently a same average value of current through each phase. Therefore, it is advised to implement two PI controllers, one for each leg as it is shown in Fig. 12. Even though, these two controllers can help the equal distribution of currents average value but it does not help current sharing of their dynamic values [2].



Fig. 12. Control Circuit for Two-Leg Chopper Regulator without a Transformer.

## VII. Simulation Results

A two leg chopping up regulator with directly coupled flattering inductive smoother or with an interphasing centre-tap transformer is analyzed and simulated. All the components of the regulator including BJTs, diodes, coupled flattering inductive smoother are assumed to be ideal. They are replaced by prototype regime is of almost zero resistance in the on-state, approximately zero fall and rise switching time and infinite resistance in the off-state. Dynamic properties of the diodes and switches were ignored. It does not have a weighty effect on the main chopping up operation, particularly when the chopper up operating frequency is not closed to the cutoff frequency of the transistors.

Equations 36, 38 and 41 are the desired results of the analysis of the chopping up regulator since they represent the behavior of its input and phase currents in the discontinuous and continuous regions shown in Fig. 13.

This figure illustrates the source current peak-to-peak ripple of the chopping up regulator with fundamental connection, flattering inductive smoother, and with an interphasing centre-tap transformer. With magnetically uncoupled flattering inductive smoother, the maximum input current ripple is 0.125A. With direct coupling, the input current ripple is reduced to 0.0625A.

Wherein $\rho$ is the admissible density of the current flowing through a cylindrical coil wire of cross-sectional area $S_{wire}$, and $S_{cu}$ is the cross-sectional area of the whole winding of the coil with N turns.

Fig. 14 illustrates the decrease in the amount of material required for the transformer and filters design process and in the magnetization current as a result of the increase in switching frequency or strength of the magnetic coupling between its filters.

Fig. 13. Source Current Ripple Versus Duty Ratio.



Fig. 14. Plot of Material Design Versus (a) Magnetization Inductance. (b) Switching Frequency.

## VIII. CONCLUSION

This paper emphasis on the idea of design process, overall current and voltage reduction of two leg dc-dc chopping up regulators by using magnetically coupled flattering inductive smoother wrapped around a common core identified as an interphasing centre-tap transformer. Detailed analysis has been done, while simulation and experimental results have been done to validate the concept. In addition, it has been found that phase current fluctuating waveform will also be decreased

as compared to the fundamental regulator connection and to the connection with magnetically uncoupled inductive smoother. This may improve the dynamic performance of the system and increase its efficiency. Furthermore, to obtain better results, magnetic coupling coefficient should be carefully chosen.

Using the interphase transformer with the chopper up regulator can help not only to accomplish both conditions of equal sharing of currents into individual legs (dynamic and steady-state) but also, as it has been said at the beginning of this research, it abolishes the need for using any controllers.

An interleaved two-leg chopping up regulating topology, however, improves regulator performance at the cost of additional smoothers, power electronic devices, and input rectifiers. Since the smoother is the largest and heaviest component in a power boost converter, the use of a coupled smoother, where a core is shared by multiple regulators instead of using multiple discrete smoothers, offers a potential approach to reducing the filters size, volume and weight. Coupled smoother topologies can also provide additional advantages such as reduced core and winding losses as well as improved input and output current and voltage ripple characteristics. The designer may choose to reduce either the boost smoother volume or increase the switching frequency to reduce the size of the EMI filter. In some cases just adding an additional phase will reduce the size of the EMI filter. Multiphase connection of such regulators also reduces the RMS current in the chopping up regulator capacitor greatly reducing electrical over stress on the capacitor. However, further study should be done to show if the complexity and cost of the design will increase with each additional phase.

REFERENCES

[1] Walid Emar, Khader Mohammad, and Mahdi Washhais: Multileg Step Up Power Converter with Magnetically Uncoupled Filters, (International Journal of Power and Energy Systems), ACTAPRESS, Vol. 36, No. 1, 2016.

[2] Issam Trrad,Walid Emar, Ziad Sobih: Novel Hybrid Two-Phase Interphase-Reactor Boost Type Converter with Controlled Output Voltage and Sinusoidal Input Utility Voltage, IJREEE, Vol. 3, issue 2, 2015.

[3] Muhammad H. Rashid. Power Electronics, Circuits, devices, and Applications, Electrical and Computer Engineering, University of West Florida, Pearson Education International, third edition, 2004.

[4] Khadmun W , Subsingha W. High Voltage Gain Interleaved DC Chopper up regulator Application for Photovoltaic GenerationSystem, Energy Procedia, ELSEVIER, 10th Eco-Energy and Materials Science and EngineeringSymposium, Volume 34, 2013; Pages 390–398.

[5] N. Mohan, T. M. Undeland, and W. P. Robbins, "Power Electronics, Converters, Applications, and Design," John Wiley & Sons, 2003.

[6] Ying qiu, Helen Liu, and Xiyouchen: Digital Average Current regime control of PWM DC-DC Converters Without Current Sensor. In: IEEE Transactions on Industrial Electronics, May 2010, Vol. 57, pp. 1670-1677.

[7] Chander, S. Agarwal, P. and Gupta, I.: Auto-tuned, discrete PID controller or DC-DC converter for fast transient response. In: International Conference on Power Electronics, Dec 2011, pp. 1-7.

[8] Chuanlin Zhang, Junxiao Wang, Shihua Li, BinWu,and Chunjiang Qian.: Robust Control for PWM-Based DC–DC Boost Power Converters With Uncertainty Via Sampled-Data Output Feedback.In: IEEE Transactions on Power Electrionics, Jan 2015, Vol. 30, No. 1, pp. 504-515.

[9]   Ramesh kumar, K, Jeevananthan, S.: A Sliding Regime Control for Positive Output Elementary Luo Converter. In: Journal of Electrical Engineering,Nov 2010, Vol. 10, pp. 115-127.

[10]  Subramanian Vijayalakshmi, and Thangasamy SreeRenga Raja.: Time domain based digital PWM controller for DC-DC converter. In: Automatika, Dec 2014, Vol. 55, No. 4, pp. 434 – 445.

[11]  Rengamani Shenbagalakshmi, and Thangaswamy SreeRenga Raja.: Discrete prediction controller for DC-Converter. In: Acta Scientiarum. Technology Maringá, Jan 2014, Vol. 36, No. 1, pp. 41-48.

[12]  Lung sheng yang, tsorng-juu liang ,hau-cheng lee, j. chen, " Novel high step up DC-DC converter with coupled smoother and voltage doubler circuit" IEEE 2011 trans. Industrial electronicsvol 58.

[13]  Yi-ping hsieh, jiann-fuh chen, tsorng-juu liang lung-sheng yang "Novel high step up DC-DC converter with coupled smoother and switch-capacitor techniques" IEEE 2012 trans. Industrial electronics vol 59.

[14]  M. Jang, M. Ciobotaru, and V.G. Agelidis, "A Single-Phase Grid-Connected Fuel Cell System Based on a Boost-Inverter," IEEE Trans. Power Electronics. Appl., vol. 28, no. 1, pp. 279–289, Jan. 2013.

[15]  F. Lin Luo and H. Ye "Power Electronics: advanced conversion Technologies" CRC Press,  2010.

[16]  M. Amundarain, M.Alberdi, A.J. Garrido, and I. Garrido, "Regimeling and Simulation of Wave Energy Generation Plants: Output Power Control" IEEE Transactions on Industrial Electronics, Vol. 58, No. 1, Pp. 105-117, 2011.

[17]  K.G. Remya, Chikku Abraham, and Babita R. Jose, " A Slope Compensated Current Regime Controlled Boost Converter",

Communications in Computer and Information Science, 2012, Volume 305, Part 3, 69-76.

[18]  Mohan, Undeland and Robbins "Power Electronics: Converter Applications and Design" John Wiley & Sons, ISBN, 978-0471226932, 2002.

## AUTHOR'S PROFILE

Walid Emar received his B.Sc. and his M.Sc. degree in power electronics in 1996/1997 and his Ph.D. degree in power electronics and control in 2002 from the University of West Bohemia, Czech Republic. Currently, He is at Isra University, Jordan as a full-time associate professor for teaching Energy Management for Master degree students and control systems, electrical machines, in addition to power electronics and other subjects for undergraduate students. He is also engaged in research in control of power electronics and machinery control.

Omar A. Saraereh received his B.Sc. in telecommunication engineering from Mu'tah University, Jordan, in 1999, the M.Sc. degree in digital communication systems in U.K., and the Ph.D . degree in electrical and electronic engineering/mobile communications from Loughborough University, Loughborough, U.K., in 2005. He is currently an Associate Professor with the Department of Electrical Engineering, The Hashemite University, Jordan. He has published many papers in various international journals and conferences.

# CWNN-Net: A New Convolution Wavelet Neural Network for Gender Classification using Palm Print

Elaraby A. Elgallad[1]
Deanship of Information Technology
Tabuk University, KSA

Wael Ouarda[2], Adel M. Alimi[3]
Research Groups in Intelligent Machines
ENIS, BP 1173, Sfax, 3038, Tunisia

*Abstract*—**The human hand is one of the body parts with special characteristics that are unique to every individual. The distinctive features can give some information about an individual, thus, making it a suitable body part that can be relied upon for biometric identification and, specifically, gender recognition. Several studies have suggested that the hand has unique traits that help in gender classification. Human hands form part of soft biometrics as they have distinctive features that can give information about a person. Nevertheless, the information retrieved from the soft biometrics can be used to identify an individual's gender. Furthermore, the soft biometrics can be combined with the main biometrics characteristics that can improve the quality of biometric detection. Gender classification using hand features, such as palm contributes significantly to the biometric identification domain and, hence, presents itself as a valuable research topic. This study explores the use of Discrete Wavelet Transform (DWT) in gender identification, with SqueezeNet acting as a tool for unsheathing features, and Support Vector Machine (SVM) operating as discriminative classifier. Inference is made using mode voting approach. Notably, the two datasets that were crucial for the fulfillment of the study were the 11k database and CASIA. The outcome of the tests substantiated the use of voting technique for gender recognition.**

*Keywords*—*Deep learning; feature extraction; gender; voting*

## I. INTRODUCTION

Biometrics is a verification technique that incorporates science (biology) as well as technology. It is highly applicable in information assurance. The method uses human biological data, including DNA and fingerprints to ensure safe entry, data access, and protection. The systems normally have high-tech elements that are connected to produce powerful performance. Biometric systems are popular among institutions that have security systems as well as replacement systems including PIN, and ID replacement. The main distinction between biometric system and other conventional systems is the fact that the former requires the physical presence of the individual who is using the system. This physical presence creates an extra level of security and makes it difficult for identity thieves to use a false ID card or any other stolen mode of identification.

Biometric systems are principally divided into two categories on the basis of physiological traits and behavioral traits. Biometric systems with the physiological identifiers focus on authenticating the physical features of the individual being verified such as the face, fingerprints, structure of the fingers, iris, DNA, retina, and other human physical

characteristics that are uniquely different. On the other hand, behavioral identifiers refer to the inimitable traits of individuals such as their typing patterns, manner of walking, use of gestures, and so on. The system can rely on the behavioral traits to offer uninterrupted verification rather than a solitary one-time check.

Accordingly, this paper focuses on gender recognition in biometric identification and explains how it can be used to enhance the accuracy of biometric identification processes. The paper is organized into five sections. The first section provides an introduction to this topic. The second section comprehensively explores the existing literature on gender identification. The third section details the manner in which the suggested technique operates. The fourth section outlines the trials and outcomes of the proposed method. Lastly, the fifth and sixth sections discuss the findings of these results and eventually provide the conclusion of this paper.

## II. RELATED WORK

There literature on gender recognition remains scarce when compared to person identification that has attracted several academicians and experts in the field. The scarcity of related literature can be explained by the fact that person identification, as opposed to gender identification, has a practical security interest for the government and law enforcement agencies. Nevertheless, gender recognition has the potential to become the next hot topic when it comes to human-computer communication. Furthermore, gender identification can also be a valuable tool for various organizations that need intelligent advertising based on accurate gender identification. One of the pioneers is [1] who has explored the concept of using computers to classify gender based on the distinctive features of the hand. He came up with a system that splits the hand delineation into various fragments; it aligns with the fingers as well as the palms, and defines the features of every section using different techniques comprising boundary descriptors, file descriptors, region descriptors and ZMs.

The respective gender was allocated based on the consolidated evidence from the hand profile with the match score being determined using the score-level fusion and file descriptors. The best performance was a match score of 98%. Notably, the evidence provided by ZMs after utilizing the score-level fusion was almost similar to that of file descriptors. Afifi [2] shared a logical record of hand images that can be used in biometric detection and gender identification. Authors presented a collection of valuable metadata drawn from

Mahmoud's logical record. The researchers also proceeded to offer top-notch approaches for using the data record to in gender detection.

There is also a convolutional neural network (CNN) based on a two-steam data that indicates the gender identification challenge. CNN helps to disclose features that are then conveyed to a support-vector-machine that analyzes data and classifies the biometric recognition issue. There is a strong belief that the suggested data record will significantly contribute towards the creation of better gender detection and biometric systems that rely on hand images. Another researcher, Ming and Yubo [3] developed a palm geometric-oriented method for identifying whether a person is a male or a female with the support-vector machine. This method is less complicated and has been shown to be effective with regards to gender classification. They note that this method guarantees a high-level of accuracy without using any intricate computations, processes, and procedures that have been used less successfully in other biometric techniques.

Instead, it is characterized by few features, including the ability to record hand features without the users having to peg their hands on the device. Its efficiency, simplicity, and convenience make it appropriate for real-world execution. Meanwhile, Font-Aragonesin and Faundez-Zanuy [4] proposed a technique that uses anthropometric hand data to identify the respective genders of the sample being tested. Using their method, researchers collected data on unvaried number of several men and women in their visual hand database. Majority of the information was collected on men.

A simple approach was then devised to get the measurements of the users' hands. The information was transmitted through a Biometric Dispersion Matcher (BDM) to retrieve the suitable data. BDM functions as a quadratic discriminant classifier. This discriminant classifier begins with filtering out data that will not assist in disclosing the user's gender. It then proceeds to display a vector of the key computations. The technique had a performance rate of 95% where the ratio of men to women was 2:1, with a projection that the accuracy could become higher when data records increase.

Gender identification through the use of facial features remains a problematic issue in the field of computer vision. Previous efforts in the face image gender recognition concentrated on enhancing understanding the neural network. Key efforts included the double-layered neural network known as SEXNET that was the brainchild of Gollomb and Lawrence [5]. Another development in face image gender recognition was that of Yen [6], which sought to eliminate the inaccuracies that are caused by less important data such as changes in face lighting as well as adjustments in facial expressions. The authors developed a novel algorithm that could help to reduce the impact of the latter disturbances. The 2-D Gabor transform was effective in spotting the facial points, with the SVM discriminative classifier picking out distinctive features to determine the respective gender of the user.

Meanwhile, it is apparent that many researchers have focused on Local Binary Patterns (LBP) and the alternatives in computer vision that can be used in gender classification. One

of the suggestions shared by researchers, in this case [7], is Binary Robust Independent Elementary Features (BRIEF), Oriented FAST & Rotated BRIEF (ORB) as well as Binary Robust Invariant Scalable Keypoints (BRISK), which have been reported to be reliable and fast in gender classification. When compared to LBP, the latter variants provide speedy detection while still maintaining high performance rates as LBP. This makes them better than LBP when used in systems that require constant gender recognition. Meanwhile, there is no sufficient literature to validate that these options are more suitable than LBP.

Notably, in another study [8], they sought to determine if face images with deep learning could help in gender classification. Specifically, they used Local Receptive Field-Extreme Learning Machine (LRF-ELM) as well as Convolutional Neural Networks (CNN). The technique was tested using data records that were meant to reveal the age and gender of the sample. The outcome of the experiment revealed an accuracy score of 80% for LRF-ELM, and 87.13% for CNN.

Another approach that captured the interest of researchers is gender recognition through boxing action. Author in [9] tested this technique whereby a period detection procedure is applied and thereafter, an averaged profile is used to denote a boxing sequence of a period. The human classification is done using Nearest Neighbor Classifier (NNC), which basically compares the recorded data with the data set that is being tested. The NNC was based on Euclidian metric. The tests were done on the KTH-Dataset, which has a good performance rate of at least 80%. PCA was used to extract features, and finally SVM was used to allocate the respective genders. In [10], they have also recommended Global Local Feature Fusion (GLFF) as a reliable method for gender recognition. The study relied on findings from psychophysics as well as neurophysiology observations, which asserted that universal and local data is important in image perception.

The initial step in the GLFF technique is mining universal and local traits through Active Appearance Model (AAM) and LBP tool respectively. One then proceeds to merge the universal and local features through sequent selection. Eventually, the chosen traits are used to determine the receptive genders in the dataset with the help of SVM. The test was made on a sample of 20 men and 20 women with the results having an average accuracy rate of 80% and an indication that the accuracy levels could go up. Based on these test results, it was apparent that blending local and global features could significantly improve the precision of gender recognition and the performance of gender classification tools.

Furthermore, the study presented the significance of different body features on gender classification. The numerical analysis revealed that the head makes a substantial contribution in gender recognition with the buttocks and leg making little contribution towards accurate gender identification.

Conversely, in [11], they suggested the use of facial features to classify gender by combining the findings of different SVM classifiers. Typically, three SVM classifiers are used to achieve this goal. The first descriptor, Histogram of Oriented Grardients (HOG), analyzes the features that pertain

to shape with the help of histogram intersection kernels. The second visual descriptor, LBP, focuses on the texture of the facial features, also with the help of histogram intersection kernels. The third descriptor analyses the raw pixel values using a linear kernel. The researchers used this approach to examine data that had been mined from FERET. The good performance rate of this test stood at 92.6%, which surpasses the accuracy of some commercial tools such as Face++.

Ardakany and Louis [12] dealt with the challenges of gender recognition through genetic algorithms. The researchers mined relevant data on facial feature records available in FERET database using LBP and PCA. The algorithm allows extracting the features that enhance the capability of the SVM classifiers to accurately recognize whether one is male or female. The genetic algorithm condenses the features being analyzed from 142 to about 71, and thereby increasing level of accuracy to about 98.5, making the results completely reliable.

Meanwhile, in the study [13], a discriminatively-trained CNN was found to be effective in classifying the gender of pedestrians. The CNN has a complex tier of neural networks that fuse feature mining and grouping into a solitary framework. The researchers recorded an accuracy of 80.4% when they used simple architecture and nominal prepossessing of dataset that provided complete body features of the pedestrians that were sampled. The accuracy of the results is equivalent to other high-tech approaches that do not use hand-manufactured trait miners.

In [14], they sought to identify gender through analyzing the different walking style of individuals. The researchers were inspired by previous studies, which have shown that gait can be used in behavior recognition tests. SVM was used to deduce the respective genders of the sample data in the analysis of the sparse spatio temporal features. The method attained a decent performance rate of 87%.

Another popular proposal in the current literature is the one made by Haitpoglu and Kose [15], which suggests combining Speed-Up Robust Features (SURF) of bags-of-visual-words (BOW) and SVM. This proposal was tested on different parts of 3560 face samples retrieved from the FERET database to determine the degree of its efficiency and reliability. The results showed that the approach is efficient in classifying gender for records retrieved from FERET database.

However, it provided more accurate recognition on the frontal face images than those on the left and right side. Meanwhile, in [16], author proposed an automated gender classification technique that relies on CNN. This method begins with network training, which is accomplished through combining numerous face datasets retrieved from different databases including ~70000 facial images from the World Wide Web. Once the data for the networks had been recorded, they were assessed, and compared with various network architectures that displayed better performance.

From their experiment, inception-v4 network had the highest accuracy level at 98.2% with Audience dataset coming second at 84%. Crime and Pedrini [17] used a pre-defined face silhouette prototype to create geometric descriptor for classifying gender. The researchers used this method to analyze four sets of face data records and reported better performance than other techniques that use geometric descriptors.

Another study [18], focusing on how the unique manner of walking could help in gender recognition shared a more advanced gait energy image (GEI) titled D-GEI. This was a step in the right direction in terms of building the literature on gait analysis and its place in gender recognition. Most of the existing studies have not extensively dealt with gait analysis, which has resulted in poor performance of gait evaluation in human identification.

The procedure in the D-GEI method begins with creating a dynamic region, followed by establishing the dynamic region of frame, and then computing the weighted average of the dynamic region. These steps help in determining the D-GEI. Once the latter has been established, HOG is used in gradient computation, while SVM classifier helps to determine the respective genders of the sample. Author in [19] presented another method that could be used in gender recognition based on the facial features of the user. In this proposal, the domain-specific as well as trainable traits are merged to help in gender classification. Fifty-two facial features linked to the eyes, nose, as well as the mouth were mined to act as SURF descriptors.

Likewise, COSFIRE acted as the trainable traits. This method responds strongly to some of the notorious challenges associated with face profiles such as the expression variations, light adjustments, and different poses. It attains a high accuracy rate on some of the top datasets in gender classification: GENDER-FERET and LFW. The technique was also highly reliable when real-world data records were used. The datasets comprised 206 training (144 men and 62 women) images, as well as 200 test (139 men and 61 women) images that were taken when the sample population on the normal walking motion. The purpose of assessing the performance of this algorithm on real scenarios was to predict its accuracy on analyzing images retrieved from videos, that are likely to be more problematic that standard data records.

Despite the problematic issues in real scenarios, the algorithm was 91.5% accurate in classifying the genders. The researchers further observe that the COSFIRE is more reliable than the SURF descriptors and therefore, recommends that it can be used solve several types of visual pattern classification problem.

## III. METHODOLOGY

Fig. 1 indicates the CWNN-Net system which this paper seeks to propose. The system works together with Squeezenet to extract features while SVM remains as the discriminative classifier. The final interpretation is made using the mode voting method. There are two databases that are relied upon in this system are 11k database [2], and CASIA database [20] each containing sufficient hand images and palmprint images respectively. The images were changed to RGB and minimized to 227×227 pixels to enhance CNN feature extraction.

### A. Features Extraction

Feature extraction discloses the nature of shape in a given pattern and thereby, simplifying the process of sorting the pattern using a formal method. It normally entails minimizing

the number of random variables being analyzed until one is left with the main variables. In pattern recognition and in image processing, feature extraction is a particular type of reducing dimensionality. Its key objective is to mine the pertinent material from the original sample and present in a way that makes it easy for image processing and pattern classification, Kumar and Bhatia [21].

*1) Discrete Wavelet Transform (DWT):* The Continuous Wavelet Transform (CWT) changes an uninterrupted signal into one that has two uninterrupted variables: translation and scale. The subsequent signal after the modifications made by CWT has a less complicated interpretation and beneficial when it comes to time-frequency analysis. The continuous wavelet transform of continuous function, x(t) relative to real-valued wavelet, ψ(t) is described by:

$$W(a,b) = \frac{1}{\sqrt{a}} \int \psi\left(\frac{t-b}{a}\right) s(t) dt \qquad (1)$$

where $\psi$ is the analyzing wavelet, *a* represents a time dilation, *b* a time translation, and the bar stands for complex conjugate.

DWT is a robust signal processing tool. Although DWT has some similarities with CWT, the difference is more apparent when it comes to scale and position values: the former scales and position values apply powers of two.

The values of s and t are: s=$2^j$, τ= k*$2^j$ and (j, k) ∈ $Z^2$ as shown in (2):

$$\psi_{s,\tau}(t) = \frac{1}{\sqrt{2^j}} \psi\left(\frac{t-k*2^j}{2^j}\right) \qquad (2)$$

The focus in DWT as well as inverse DWT is disintegration and restoration. The disintegration and restoration are achieved through LPF and HPF. The effect of wavelet disintegration is a structured decomposition that occurs in tiers. The level of disintegration is selected as per the preferred cutoff frequency. Fig. 2 illustrates DWT undergoing a process of LPF and HPF [22].



Fig 1. CWNN-Net: A New Convolution Wavelet Neural Network for Gender Classification using Palm Print.



Fig 2. A Three-Level forward DWT via a Two-Channel Iterative Filter Bank [22].

*2) Discrete haar wavelet transform*

The original Haar definition is as follows:

$haar(0,t) = 1, for\ t \in [0,1);$

$$haar(1,t) = \begin{cases} 1, for\ t \in \left[0, \frac{1}{2}\right), \\ -1, for\ t \in \left[\frac{1}{2}, 1\right) \end{cases} \qquad (3)$$

One of the key characteristics of the Haar functions other than the Haar (0, t), the i–th Haar function can occur through the restriction of the (j − 1)–th function to be 50 percent of the interval where it is different from zero, through multiplying it with p2 and scaling over the interval [0, 1]. These traits make Haar function makes share a degree of link with the wavelet concept and can be used to define wavelet. In this context, two Haar functions represent universal functions, while the remaining functions represent the local functions. As such, the odd rectangular pulse pair, a Haar function in this setting, denotes the oldest and the most basic wavelet. The purpose of using DWT is to get data that is highly differentiated by offering distinct solution at separate segments of the time-frequency domain. Through wavelet transforms, one can create unequal segments of time-frequency plane in line with the time-spectral aspects of the signal. The wavelet technique has robust links to the traditional standard basis of the Haar functions. Scaling as well as shrinking a standard wavelet produces classical Haar functions.

Let ψ: R → R, the Harr wavelet function is defined by the formula:

$$\psi(t) = \begin{cases} 1, for\ t \in \left[0, \frac{1}{2}\right), \\ -1, for\ t \in \left[\frac{1}{2}, 1\right), \\ 0, otherwise. \end{cases} \qquad (3*)$$

for any Haar function (except function haar (0, t))

From basis (3*) may be generated by means of the formulas:

$$\psi_i^j(t) = \sqrt{2^j}\, \psi(2^j t - i)$$

$i = 0,1, \dots, 2^j - 1.\ and\ j = 0,1, \dots, log_2 N - 1 \qquad (4)$

The constant $\sqrt{2^j}$ is chosen so that:

the scalar produc $< \psi_i^j, \psi_i^j \geq 1, \psi_i^j(t) \in L^2(R)$

Let Φ: R → R, the Harr scaling function is defined by the formula:

$$\Phi(t) = \begin{cases} 1, & for\ t\ \in [0,1), \\ 0, & for\ t\ \notin [0,1). \end{cases} \qquad (5)$$

Similarly, to the properties of the wavelet function, for scaling function one can define the family of functions:

$$\Phi_i^j(t) = \sqrt{2^j}\ \Phi(2^j t - i)$$

$$i = 0,1, \ldots, 2^j - 1.\ and\ j = 0,1, \ldots, log_2\ N - 1 \qquad (6)$$

The constant $\sqrt{2^j}$ is chosen so that:

the scalar product $< \Phi_i^j, \Phi_i^j \geq 1, \Phi_i^j(t) \in L^2(R)$

In the two-dimensional case, three sets of detail coefficients residing in the horizontal, vertical, and diagonal directions. The subbands $LH_j$, $HL_j$, and $HH_j$, j = 1, 2... J are the detail coefficients, as noted above, where j is the scale and J denotes the largest or coarsest scale in the decomposition. Fig. 3 shows a representation of the multilevel wavelet decomposition at level 3.



Fig 3.     Level 3 Wavelet Decomposition [23].

*3) Squeeze net*: The highlight of most researchers exploring CNN has been how the field can achieve create precision on machine vision datasets. Several models have been built to achieve different levels of accuracy. Reference [24] have reviewed different CNN techniques to establish which model would require simpler parameter but attain a degree of precision that is similar to some other popular techniques. From his analysis, he concluded Squeeze Net has minimal parameters but still achieves a performance level that is comparable to common models. Fig. 4 and Fig. 5 shows the organization of convolution filters in the Fire module.

### Approach 1: Resizing the network 3x3 filters to 1x1 filters

This approach focused on decreasing the number of parameters resizing 3x3 filters to 1x1 filters. At the onset, this approach was confusing as one would assume that replacing 33x3 filters with 1x1 filters would provide less information and thereby, lower the accuracy of the model. Nevertheless, this was not the case. Naturally, 3x3 filters takes spatial data of pixels that are in close proximity. On the other hand, in 1x1 filters solely concentrates on one pixel and extracts the connection of channels within the pixel without focusing the adjacent one.



Fig 4.     Microarchitectural view: Organization of Convolution Filters in the Fire Module. In this Example, s1x1 = 3, e1x1 = 4, and e3x3 = 4. We Illustrate the Convolution Filters but not the Activations [24].



Fig 5.     Microarchitectural view: Organization of Convolution Filters in the Fire Module. s1x1 = 3, e1x1 = 4, and e3x3 = 4 [24].

Forrest used three main approaches when evaluating the different CNN architectures:

### Approach 2: Minimizing the number of inputs in the residual 3x3 filters

The approach focuses decreasing the parameter through cropping filters. This methodical strategy is actualized through feeding "squeeze" stratums into "expand" layers. Squeeze layers consist of 1x1 filters while "expand" layers have both 3x3 filters and 1x1 filters. Therefore, aggregate number of parameters are minimized through the lowering the "squeeze" stratum filters. This process is referred to as the "fire module" is used as the foundational block of creating the Squeezenet design.

### Approach 3: Downsample the end part of the network to create large feature maps in the convolution strata

Author in [24] propose that reducing the stride with the later convolution strata results in larger feature map, which enhances the accuracy of recognition. The manner in which the activation maps occur at the end of the network differentiates this design from well-known networks such as the VGG which is characterized with feature maps that shrink as one gets to end parts of the network.

The building block of Squeezenet is referred to as fire module. The fire module has two strata: squeeze layer and an expand layer. A Squeezenet has a pile of fire modules that are accompanied with minute pooling layers. Notably, both squeeze layer and the expand layer parallel activation map size. Nevertheless, the depth decreases in squeeze layer, while it surges the expand layer. Meanwhile, the neural models have a bottleneck layer and a trait of increasing. Besides, they have a tendency to generate a pattern of a surging depth, while minimizing activation map size with the intention of obtaining an elevated level abstract.

As illustrated in the chart above, the squeeze segment only has 1x1 filters. The latter implies that it operates as completely-attached layer dealing with feature points on the same area. The spatial abstract cannot be derived from the system. Notably, it helps in decreasing the depth of feature map, which simplifies and increases the speed of computation in the expand layer.

Reducing depth enables to achieve an effect when the number of computations performed by 3x3 filters in the expand layer decreases. The speed increases due to the fact that 3x3 filters required only as much computation as 1x1 filters need. Logically speaking, too much squeezing will decrease the exchange of data, while limited number of 3x3 filters is bound to inevitable limit space resolution. The SqueezeNet architecture enabled a 50X decrease in the size of the model when compared to AlexNet; it is worth noting that top-1 as well as top-5 accuracy of AlexNet was exceeded.

### B. Cross Validation

The suggested system uses SVM for classifying the different genders on the basis of train-test split approach that is one of techniques for cross-validation. The predictive patterns are computed through the division of the basic model to create a training set to train the system, and set test for evaluation.

SVM was previously proposed By Vapnik [25]. The support vector machine is among the maximum margin classifiers and based on the Structural Risk Minimization. SVM plan input silhouette to a top dimensional domain where the uppermost separating hyperplane is gathered. Linear support vector machine is primarily pronounced for binary classification.

Take for instance that the training data set as well as the labels $(x_n, y_n)$, $n=1,...,N$, $x_n \in \mathbb{R}^D$, $t_n \in \{-1, +1\}$, SVMs learning has the future controlled optimization:

$$\min_{w, \xi_n} = \frac{1}{2} w^T w + C \sum_{n=1}^{N} \xi_n$$
$$s.t. \quad w^T x_n t_n \geq 1 - \xi_n \qquad \forall_n \qquad (7)$$
$$\xi_n \geq 0 \qquad \forall_n$$

where $\xi_n$ are the slack variables, $w$ is the vector of coefficients, and $C$ is the capacity constant.

The unconstrained optimization problem as in (8) which is recognized as the primal form problem of L1-SVM:

$$\min_w = \frac{1}{2} w^T w + C \sum_{n=1}^{N} \max(1 - w^T x_n t_n, 0) \qquad (8)$$

Meanwhile L1-SVM is not differentiable, the L2-SVM is used which minimizes the squared hinge loss as in (9):

$$\min_w = \frac{1}{2} w^T w + C \sum_{n=1}^{N} \max(1 - w^T x_n t_n, 0)^2 \qquad (9)$$

To expect the class label of a test data x:

$$arg_t \max(w^T x) t \qquad (10)$$

To extend SVMs for multiclass problem, one-vs-rest approach is used. Representing the output of the *k-th* SVM as in (11):

$$a_k(x) = W^T x \qquad (11)$$

the forecast class is

$$arg_k \max a_k(x) \qquad (12)$$

### C. Score Fusion

A key question that should be addressed when handling information on combining systems is the nature of data that needs to be merged using the fusion module. There are several approaches that have been devised to help in the latter at different levels including the sensor level, feature level, rank level, as well as the decision level. In line with our recommended system, we will analyze combining systems considering it is the principal level of fusion at the match score level [26].

At the match score level, the relationship between input and prototype biometric feature silhouettes. Synthesis at the core level entails the successful combination of various match scores output to obtain a biometric recognition decision [26]. Some of the strategies used in consolidating the different biometric matchers include Majority Voting, and Weighted Majority Voting. Our proposed system recommends the use of Mode Voting Technique (MVT) when consolidating information at the decision level.

The method utilizes the standard class label values that are retrieved from the predicted label array obtained through the SVM discriminate classifier. MVT is used to establish the common non-repeated values in the predict label array X for the purpose of biometric identification.

$$Z = \text{mode}(X_{k,i}) \qquad (13)$$

where *Z* is the class label of the test image, *k* is the index of the test image, and *i* is the index of the descriptor.

## IV. EXPERIMENTAL RESULTS

As apparent in the methodology section, two databases (11k Database and CASIA Database) were used to obtain the relevant data records. The 11k database had hand images taken from individuals whose photos were stored in jpg formats.

The hand images displayed the palm/dorsal, and left and right images of the individuals. In total, images of 190 subjects, aged between 18 and 75, were retrieved from the database. As for the CASIA database, it had 5,503 palmprint images belonging to 312 persons and representing both the left and right palms.

The images of JPEG format each 8 bit gray-level. Slight adjustments were made to the images to facilitate CNN feature extraction. Specifically, they were modified to RGB and minimized to 227×227 pixels.

SVM was applied to discriminately sort the features. The features were retrieved from single-level 2D DWT through the application of Haar wavelet filter. Through this process, DWT provided an approximation coefficient matrix cA as well as detail coefficients matrices cH, cV, and cD (horizontal, vertical, and diagonal, respectively).

Two experiments are carried out over 11k database, one on all dataset, and the other on the dataset after extracting hand images with accessories. The predicted label arrays obtained from SVM as a classifier for the obtained features are fused using mode voting technique. Table I and Table II summarize the results obtained for both experiments.

The recognition rate for all dorsal dataset ranged from 97.57% to 99.71%, with mode voting technique on the resulting 8 predicted label arrays, the recognition rate reached 99.86% with processing time 0.27 sec for each image.

For all palmer dataset ranged from 96.86% to 98.57%, with mode voting technique on the resulting 8 predicted label arrays, the recognition rate reached 100% with processing time 0.28 sec for each image.

By Applying mode voting technique on the resulting 16 predicted label arrays for previous both dorsal and palmar images, the recognition rate reached 100% with processing time 0.55 sec for each image.

After excluding images with accessories, the recognition rate for dorsal dataset ranged from 97.27% to 99.64%, and for palmer dataset ranged from 97.27% to 99.27%. Also, after applying mode voting technique on the resulting 16 predicted label arrays, the recognition rate reached 100% with processing time 0.57 sec for each image.

For CASIA database, the experiment is carried out over the database which contains only palmprint images. Table III summarize the obtained results.

The recognition rate for palmer dataset ranged from 81% to 97.75%. After applying mode voting technique on the resulting 8 predicted label arrays (as the database only contain palmar side), the recognition rate reached 98% with processing time 0.25 sec for each image as in Table III.

TABLE I.    RESULTS OF 11K DATABASE FOR ALL IMAGES

| Descriptor | dwt2 | dorsal L | dorsal R | palmer L | palmer R |
|---|---|---|---|---|---|
| squeezenet fire9-concat | cA | 99.43 | 99.71 | 97.29 | 98.57 |
| | cH | 98.71 | 97.57 | 97.71 | 98.29 |
| | cV | 98.75 | 98.29 | 97.43 | 98.00 |
| | cD | 98.14 | 97.86 | 96.86 | 98.43 |
| Voting | | 99.86 | | 100 | |
| | | 100.00 | | | |
| Sys. Time/image in sec. | | 0.27 | | 0.28 | |
| | | 0.55 | | | |

TABLE II.    RESULTS OF 11K DATABASE EXCLUDING ACC. IMAGES

| Descriptor | dwt2 | dorsal L | dorsal R | palmar L | palmar R |
|---|---|---|---|---|---|
| squeezenet fire9-concat | cA | 98.73 | 99.64 | 99.27 | 98.73 |
| | cH | 98.00 | 97.82 | 98.00 | 97.27 |
| | cV | 98.18 | 98.00 | 98.00 | 98.00 |
| | cD | 97.45 | 97.27 | 97.27 | 97.82 |
| Voting | | 100 | | 100 | |
| | | 100.00 | | | |
| Sys. Time/image in sec. | | 0.57 | | | |

TABLE III.    RESULTS OF CASIA DATABASE

| Descriptor | dwt2 | palmar L | palmar R |
|---|---|---|---|
| squeezenet fire9-concat | cA | 97.75 | 95.25 |
| | cH | 92.25 | 82.50 |
| | cV | 87.75 | 88.50 |
| | cD | 82.00 | 81.00 |
| Voting | | 98.00 | |
| Sys. Time/image in sec. | | 0.25 | |

It is noted that, the high recognition rate is nearly obtained from the approximation coefficient matrix cA in both databases.

## V.    DISCUSSION

From the results, it is clear that the fusion at decision level using the mode voting technique guarantees an excellent recognition rate regardless of low recognition rate of some datasets. The mode voting technique ranks top of the list of SVM classifiers used for each database.

Table IV shows the performance comparison of proposed system CWNN using 11k and CASIA databases. It is clear that the performance is stable for both databases, the dorsal and palmar results in 11k database are nearly the same, so one of them can be used alone to reduce the system processing time. Also, the processing time for two databases is nearly the same for each image.

TABLE IV.    PERFORMANCE COMPARISON OF PROPOSED SYSTEM CWNN USING 11K AND CASIA DATABASES

| Database | 11k | | | CASIA |
|---|---|---|---|---|
| Hand | Dorsal L,R | Palmar L,R | All | Palmar L,R |
| No of Images | 5680 | 5396 | 11076 | 5502 |
| Feature extraction | Squeezenet & dwt2 | | | |
| Classifier | SVM | | | |
| No of predicted arrays | 8 | 8 | 16 | 8 |
| Result by voting | 99.86 | 100 | 100 | 98 |
| Sys. Time/image in sec. | 0.27 | 0.28 | 0.55 | 0.25 |

TABLE V. PERFORMANCE COMPARISON OF PROPOSED SYSTEM CWNN VS. EXISTING SYSTEMS USING 11K

| 11k | Afifi 2017 | | CWNN | |
|---|---|---|---|---|
| | Dorsal | Palmar | Dorsal | Palmar |
| All Images | | | 99.86 | 100 |
| no accessories | 97.30 | 94.20 | 100 | 100 |

Table V shows the performance comparison of the proposed systems CWNN vs. existing systems in the previous work that use 11k database.

## VI. CONCLUSION

This work reports a CWNN-Net: a new convolution wavelet neural network for gender classification using palm print system based on the mode voting technique, and compares the performance of the system using two datasets.

The novelty comes from using mode voting technique at decision level. Our experimental results show the efficiency of the suggested system. Using SqueezeNet and DWT in this system shows a promising result due to the advantage of SqueezeNet and wavelet decomposition.

We believe that the mode voting technique can serve as a step towards the construction of more accurate gender recognition and biometric identification systems. In the future, we will adopt this deep learning method in real-time palmprint recognition system and develop a more intelligent machine learning algorithm for feature extraction in palmprint recognition.

### REFERENCES

[1] G. Amayeh, G. Bebis and M. Nicolescu, "Gender classi_cation from hand shape", In: Computer Vision and Pattern Recognition Workshops. CVPRW'08. IEEE Computer Society Conference on, IEEE, 2008, pp. 1-7.

[2] M. Afifi, "11K Hands: Gender recognition and biometric identification using a large dataset of hand images." arXiv preprint arXiv:1711.04322, 2017.

[3] W. Ming and Y. Yubo, "Gender Classification Based on Geometry Features of Palm Image", the Scientific World Journal, vol. 2014, Article ID 734564, 2014.

[4] X. Font-Aragones, and, M. Faundez-Zanuy, "Hand-Based Gender Recognition Using Biometric Dispersion Matcher", in: Apolloni, B., Bassis, S., Esposito, A., & Morabito, F., (eds) Neural Nets and Surroundings. Smart Innovation, Systems and Technologies, vol 19. Springer, Berlin, Heidelberg, 2013.

[5] B. A. Golombn and D. T. Lawrence, "SEXNET: A neural network identifies sex from human faces", Advances in Neural Information Processing Systems, 1991, pp. 572–577.

[6] C. Yan, "Face Image Gender Recognition Based on Gabor Transform and SVM", In: Shen, G., & Huang, X., Advanced Research on Electronic Commerce, Web Application, and Communication. Communications in Computer and Information Science, vol 144. Springer, Berlin, Heidelberg, 2011.

[7] F. S. Iglesias, M. E. Buemi, D. Acevedo and J. Jacobo-Berlles, " Evaluation of Keypoint Descriptors for Gender Recognition", In: Bayro-Corrochano, E., & Hancock, E., (eds) Progress in Pattern Recognition, Image Analysis, Computer Vision, and Applications. Lecture Notes in Computer Science, vol. 8827. Springer, Cham, 2014.

[8] Y. Akbulut, A. Şengür, and S. Ekici, "Gender recognition from face images with deep learning", International Artificial Intelligence and Data Processing Symposium (IDAP), Malatya, 2017, pp. 1-4.

[9] J. Wang, W. Hu, Z. Wang and Z. Chen, "Human Identification and Gender Recognition from Boxing", in: Sun, Z., Lai, J., Chen, X., & Tan, T., (eds) Biometric Recognition. Lecture Notes in Computer Science, vol 7098. Springer, Berlin, Heidelberg, 2011.

[10] W. Yang, C. Chen, K. Ricanek and C, Sun, "Gender Classification via Global-Local Features Fusion", In: Sun, Z., Lai, J., Chen, X., & Tan, T. (eds) Biometric Recognition. Lecture Notes in Computer Science, vol. 7098. Springer, Berlin, Heidelberg, 2011.

[11] G. Azzopardi, A. Greco and M. Vento, "Gender Recognition from Face Images Using a Fusion of SVM Classifiers", In: Campilho A., Karray F. (eds) Image Analysis and Recognition. Lecture Notes in Computer Science, vol 9730. Springer, Cham, 2016.

[12] A. R. Ardakany and S. J. Louis, "Improving Gender Recognition Using Genetic Algorithms", In: Bui L.T., Ong Y.S., Hoai N.X., Ishibuchi H., Suganthan P.N. (eds) Simulated Evolution and Learning. Lecture Notes in Computer Science, vol 7673. Springer, Berlin, Heidelberg, 2012.

[13] C. Ng, Y. Tay and B. Goi, "A Convolutional Neural Network for Pedestrian Gender Recognition", In: Guo, C., Hou, ZG., & Zeng, Z., (eds) Advances in Neural Networks. Lecture Notes in Computer Science, vol 7951. Springer, Berlin, Heidelberg, 2013.

[14] M. Collins, P. Miller and J. Zhang, "Gait Based Gender Recognition Using Sparse Spatio Temporal Features", In: Gurrin, C., Hopfgartner, F., Hurst, W., Johansen, H., Lee, H., & O'Connor, N., (eds) MultiMedia Modeling. MMM 2014. Lecture Notes in Computer Science, vol 8326. Springer, Cham, 2014.

[15] B. Hatipoglu and C. Kose, "A gender recognition system from facial images using SURF based BoW method ", International Conference on Computer Science and Engineering (UBMK), Antalya, 2017, pp. 989-993.

[16] C. Nistor, A. C. Marina, A. S. Darabant and D. Borza , "Automatic gender recognition for "in the wild" facial images using convolutional neural networks", 13th IEEE International Conference on Intelligent Computer Communication and Processing (ICCP), Cluj-Napoca, 2017, pp. 287-291.

[17] M.V.M Cirne and H. Pedrini, "Gender recognition from face images using a geometric descriptor", IEEE International Conference on Systems, Man, and Cybernetics (SMC), Banff, AB, 2017, pp. 2006-2011.

[18] T. Liu, B. Sun, M. Chi and X. Zeng, "Gender recognition using dynamic gait energy image", IEEE 2nd Information Technology, Networking, Electronic and Automation Control Conference (ITNEC), Chengdu, 2017, pp. 1078-1081.

[19] G. Azzopardi, A. Greco, A. Saggese and M. Vento, "Fusion of domain-specific and trainable features for gender recognition from face images", IEEE Access, vol. PP, no. 99, 2018.

[20] CASIA Palmprint Database, http://biometrics.idealtest.org/

[21] G. Kumar and P. K. Bhatia, "A detailed review of feature extraction in image processing systems", In Advanced Computing & Communication Technologies (ACCT), Fourth International Conference, 2014, pp. 5-12.

[22] R. Haddadi, E. Abdelmounim, M. El Hanine and A. Belaguid, "Discrete Wavelet Transform Based Algorithm for Recognition of QRS Complexes", World of Computer Science & Information Technology Journal, 4(9), 2014.

[23] T. Williams and R. Li, "December. Advanced image classification using wavelets and convolutional neural networks", In Machine Learning and Applications (ICMLA), 15th IEEE International Conference, 2016), pp. 233-239.

[24] F. N. Iandola, S. Han, M. W. Moskewicz, K. Ashraf, W. J. Dally and, K. Keutzer, "Squeezenet: Alexnet-level accuracy with 50x fewer parameters and< 0.5 mb model size", 2016.

[25] V. Vapnik, "The Nature of Statistical Learning Theory ". NY: Springer-Verlag, 1995.

[26] A. A. Ross, K. Nandakumar and A. Jain, "Handbook of multibiometrics", Vol. 6, Springer Science & Business Media, 2006, pp. 73-82.

# The Method of Computer-Aided Design of a Bread Composition with Regard to Biomedical Requirements

Natalia A. Berezina[1], Andrey V. Artemov[2]
Department "Food Technology and Organization of Restaurant Business"
Orel State University, Orel, Russian Federation

Igor A. Nikitin[3]
Department "Technology of Grain Processing, Bakery,
Pasta and Confectionery Industries" K.G. Razumovsky
Moscow State University of Technologies and Management
(The First Cossack University) Moscow, Russian Federation

Alexander A. Budnik[4]
Department "Automated control systems"
K.G. Razumovsky Moscow State University of
Technologies and Management (The First Cossack
University) Moscow, Russian Federation

*Abstract*—**A method for efficient software implementation of bread optimized multicomponent mixtures has been developed. These polycomposite mixtures have a chemical composition that meets the modern physiological standards of nutrition for the elderly people. To implement the developed algorithm a high-level programming language Object Pascal was used using the IDE Borland Delphi 7.0. An unconventional raw material was selected, which allows to provide necessary requirements to the quality indicators of the finished bread in all modeled mixtures. Modeling the composition of flour mixtures for gerodietic nutrition using the software made it possible to obtain compositions with a specific ratio of prescripted components, balanced in accordance with the intended purpose.**

*Keywords—Modeling; polycomposite mixture; bread; gerodietic nutrition; quality*

## I. Introduction

The approach to creation of technologies of multicomponent products with a chemical composition [1] that is regulated in accordance with modern physiological nutritional standards predetermined the intensive development of research, united by the concept of "design" of consumer food products. The problem of creating food with a given level of nutritional adequacy has been formed for a long time [2] and is highly relevant at the present time. First of all, this is due to the increasing complexity of the composition of modern food products and the emergence of new knowledge about the effect of the components they contain on human health [3], expanding the range of new types of food ingredients.

The development of a new product range and technologies of bread with a complex raw material composition currently occupies an important place in the development of bakery production. The use of new and non-traditional raw materials has high prospects in the food industry, as it allows you to get bulk food with high nutritional value, preventive and therapeutic properties.

The principles of the food products composition design are based on standards and criteria for optimizing the quality of food, based on modern knowledge of biology, medicine and food chemistry. The solution of the problem of optimizing the compositions of polycomposite food products, including rye-wheat bread, is possible using indicators that can be described mathematically. As a result of research of domestic and foreign scientists in the field of biology and medicine [4], the concept of a balanced diet was formulated. According to the concept proportions of individual substances causing the sum of exchange reactions that underlie human life were determined [5].

The numerical values [6] of the optimal nutrient ratios, given in the balanced nutrition formula, make it possible to use them as formal optimization criteria for creating polycomposite mixtures with given levels of nutritional adequacy by analytical combining the main prescription component of breadmaking-flour in combination with new and unconventional raw materials.

To date, various software products have been developed and are being actively used to automate technological calculations of food recipes of various groups.

The most common instrument for creating such programs for calculating and optimizing recipes is the spreadsheet MS Excel. The initial data for calculations is entered into the corresponding cells of the spreadsheet, and the calculation formulas are entered into the others. The advantage of this shell is the prevalence [7], simplicity of calculations and a lot of methodological and reference literature on its use. In addition, it has a number of built-in modules to simplify routine procedures for finding a solution, correlation and regression analysis, etc. The disadvantage of using MS Excel is the lack of automated input of initial data and calculated dependencies, as well as obtaining a single solution for solving optimization problems using standard tools.

The program for calculating the chemical composition of food for catering in high school institutions [8] is designed to calculate the energy and nutritional value of a school food. The program allows you to evaluate the chemical composition of the school food ration according to 28 indicators, including the calculation of proteins and fats by origin, carbohydrates by molecular weight, vitamins A, B and E by their equivalents. The disadvantage of the program is the lack of possibility to optimize the diet for nutritional adequacy.

The optimization of parameters of the product under development by modeling the formulation using the integral balance criterion for a wide range of indicators was used in developing theoretical prerequisites for computer-aided design of food products for the elderly people. At the same time, a qualimetric multiplicative model was used [9], which allows to bring relative complex and simple individual quality indicators of different origin into one form, ensuring the independence of the properties of each indicator. This approach has a disadvantage inherent in combining many contradictory factors into one criterion. The resulting solution is unstable and very empirical.

A common drawback of existing software systems is the lack of a recipe optimization subsystem, based on a set of criteria for food, biological and mineral value, as well as finding a single solution that is optimal in terms of specified parameters. In addition, this approach does not allow to predict consumer properties of the product, which does not allow to speak about the possibility of its full launch on an industrial scale.

However, modeling of polycomposite mixtures for baking, based only on the analytical combination of the quantitative and qualitative component of the nutrients contained in them, has a very significant imperfection associated with the lack of guarantees of obtaining the final product of sufficient consumer dignity.

Combining the nutrient composition of mixtures for baking without taking into account the technical-functional properties of the final variants of mixtures makes the use of improvers a prerequisite for their use. It increases the cost of the final product, and also narrows the choice from a variety of analytical calculations of recipes at the raw set stage. As a result, the final choice of a polycomposite mixture is possible only after direct testing of mixtures and this fact greatly complicates and lengthens their development.

This shows the urgency of the task of developing time-efficient formalized methods for optimizing flour polycomposite mixtures for rye-wheat bread of high nutritional value with specified technological properties. Solving the problem is possible by modeling the technological and nutrient adequacy of baking mixtures. The approach is based on an analytical assessment of partial qualities of individual components of the mixture and designing a quantitative and qualitative nutrient composition of the polycomposite mixture for rye-wheat bread with the increased nutritional value. And their physical and chemical interaction with the main raw material ingredient (flour) must be taken into account.

The aim of this work is to develop a method of effective software implementation of the nutritional value and technological adequacy of rye-wheat bread, based on an innovative model of optimal composition and simulation algorithms.

## II. METHODOLOGICAL APPROACH OF QUANTITATIVE AND QUALITATIVE EVALUATION OF THE COMPOSITION OF A POLYCOMPOSITE MIXTURE FOR RYE-WHEAT BREAD

On the basis of the methodological approach of quantitative and qualitative assessment of the nutrient composition of the polycomposite mixture was the principle of separation as the key component of the protein part of the mixture. This is largely due to the fact that proteins, being an evolutionary-conditioned dominant of the diet, in general, determine the nature of nutrition. The satisfaction of a body with this component at a physiologically necessary level allows the body to display the functions of other nutrient components of food [10].

When considering the optimization of the amino acid composition of multicomponent products, the principles of Mitchell-Block [11] were recommended as principles for calculating quality, based on the interrelation of protein quality with its biological value.

Without a doubt, the computational method for determining the biological value of a baking mix composition has a number of drawbacks inherent in all computational methods [12], since it does not take into account the difference in protein digestibility in the various components that are included in the mixture. However, for practical purposes in computer design, the calculation method is currently only acceptable for the following reasons: biological methods are extremely complex [12], lengthy and expensive; when designing multicomponent mixtures, the priority is not the absolute value of the biological value, but the finding of such a ratio of prescripted components which ensures its maximum value [13].

The modeling of the nutrient composition of the model compositions of rye-wheat bread of a gerodietic orientation was carried out taking into account the basic medical and biological requirements for this group of products [14].

Modeling the technological adequacy of the flour mixture, which ensures a stable quality of the final product, was carried out by introducing the flour technological indicator called the "drop number", calculated using the Perten formula [15].

In order to develop a time-efficient and accurate algorithm for calculating the composition of the mixture in general, the task is formulated as follows: simulate the component composition of the mixture for gerodietic bread with a biological value of at least 60%, providing the "protein:fat:carbohydrate" ratio close to 1.0:0.8:3.5, "Ca:Mg:P" - 1:0.6:1.3, fiber - not less than 0.2 g per 100 g of finished bread, the "drop number" of the mixture is 200-240 s. The calculation algorithm is shown in Fig. 1.

Fig. 1. Algorithm for Calculating Polycomposite Mixtures for Rye-Wheat Bakery Products of High Nutritional Value with Specified Technical-Functional Properties.

On the basis of the proposed algorithm for calculating a polycomposite mixture in a high level programming language Object Pascal using IDE Borland Delphi 7.0 a software for calculating the optimized composition of a model mixture for

breads with a gerodietic orientation has been developed. An example of the software is shown in Fig. 2.

### III. MATERIALS AND METHODS

The chemical composition of raw materials and the "drop number" are shown in Table 1 as determined experimentally.

Analysis of the Table 1 data shows that selected raw ingredients have a rich chemical composition compared to bread flour, which will balance the composition of model mixtures for calcium, phosphorus and magnesium, increase fiber content and other components for balanced nutrition of the elderly. In addition, each of the ingredients has a special composition that in a mixture allows to obtain a product with a wide positive spectrum of health effects. Flax and sunflower seeds containing increased amount of fats will enrich model mixtures with polyunsaturated fatty acids, dry wheat gluten consisting mainly of protein will help balance the ratio of proteins and carbohydrates. The same role is played by ingredients with increased protein content - skimmed milk powder, soy flour, lentil flour, sunflower seeds. All types of raw materials have a higher fiber content than bread flour, which has a wide positive spectrum of effects on the body. Amaranth flour will allow to enrich model mixtures for bakery products with the strongest antioxidant - squalene, and buckwheat flour - with a routine that has anti-sclerotic and anti-hypertensive effects. Barley flour contains a large amount of β-glucan polysaccharide, which has a cholesterol-lowering effect and a lipid-lowering effect.

Practically all the ingredients will enrich model compositions with calcium, necessary for the prevention of osteoporosis, reduction of bone fragility and thyroid gland activity, and recommended for some muscle diseases. Calcium lactate, used as one of the components of the mixture, is widely used in the food industry not only as a source of calcium. It also has the functions of a flavor enhancer, flour and bread improver. The quality assessment of the calculated polycomposite mixtures was carried out by analyzing the chemical composition and balance of bread main nutrients from the mixtures by indicators of the mass fraction of moisture, specific volume, porosity, yield, organoleptic rate and Invitro digestibility using the enzyme pepsin [16].

TABLE I. THE CHARACTERISTIC OF RAW MATERIALS' CHEMICAL COMPOSITION

| Name of raw materials | proteins | lipids | Carbohydrates | | Calcium | Phosphorus | Magnesium | "drop number", sec |
| | | | mono- and disaccharides | cellulose | | | | |
| | g/100g | | | | mg/100 g | | | |
| rye flour | 9.9 | 1.7 | 73 | 0.3 | 34 | 189 | 60 | 150 |
| wheat flour | 12.5 | 1.2 | 70.8 | 0.2 | 32 | 184 | 73 | 210 |
| flax seeds | 18.9 | 42.16 | 28 | 27.3 | 255 | 642 | 392 | 450 |
| barley flour | 10.2 | 1.2 | 71 | 6.9 | 80 | 175 | 63 | 520 |
| amaranth flour | 9.5 | 3.9 | 67.8 | 1.1 | 179 | 620 | 229 | 320 |
| buckwheat flour | 12.4 | 3.2 | 73.7 | 1.9 | 23.9 | 264.3 | 147 | 320 |
| skimmed milk powder | 48.2 | 1.2 | 48.2 | 0 | 1155 | 920 | 160 | 60 |
| soy flour | 48 | 0.2 | 9.7 | 1.3 | 212 | 198 | 145 | 60 |
| lentil flour | 24 | 0.2 | 38 | 3.7 | 83 | 294 | 0 | 60 |
| sunflower seeds | 20.7 | 52 | 3.4 | 27.3 | 530 | 317 | 220 | 210 |
| dry wheat gluten | 78 | 0 | 15 | 0.6 | 142 | 260 | 25 | 450 |
| calcium lactate | 0 | 0 | 0 | 0 | 16000 | 0 | 0 | - |

Fig. 2. Software Interface Example.

## IV. RESULTS AND DISCUSSION

As a result of the software work (Fig. 2), more than 40 model flour mixtures were generated. After analyzing the calculated data, 6 mixtures were selected with a biological value of at least 75%. The research results are grouped in Table 2. For comparison, the characteristics of bread flour mixtures are given.

The data in Table 2 shows that modeling of the flour mixtures composition for gerodietic nutrition using the software allowed us to obtain compositions with a ratio of prescripted components balanced in accordance with the intended purpose.

The flour mixtures modeled by the software complex were used to make bakery products using the methods adopted in baking. The bread quality indicators are given in Table 3. Loaves of bread baked only from rye-wheat and wheat flour were control samples.

The determination of the digestibility of bakery products for gerodietic nutrition was carried out "in vitro" by incubating a 20% bread suspension in a solution of pepsin in glycine buffer. Bakery products from the trading network were used as control samples: the "Nareznoy" loaf made from wheat flour and the rye-wheat bread "Spassky". The research results are presented in Fig. 3.

It was determined that quality indicators of bread baked of model mixtures for gerodietic nutrition were not lower than

control samples, and thier "in vitro" digestibility even exceeded the digestibility seen in samples of "Spassky" bread and "Nareznoy" loaf. It is due to the large amount of water-soluble proteins in their composition, which are susceptible to proteolytic cleavage. This proves the high digestibility of the developed products. The result is important for the nutrition of older people, whose metabolism is slower than at a young age.

The calculation of the content of basic nutrients in 100 g of developed bakery products for gerodietic nutrition in accordance with methodological recommendations is made [17]. The results of the calculation are shown in Table 4. For comparison, the chemical composition of the "Nareznoy" loaf and the "Spassky" rye-wheat bread was calculated.



Fig. 3. Digestibility of Bakery Products From model Mixtures for Gerodietic Nutrition.

TABLE II.  QUANTITATIVE-QUALITATIVE CHARACTERISTICS OF FLOUR MIXTURES

| Component Name | Quantitative-qualitative characteristics of model mixtures | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | *A mixture of rye and wheat flour (control 1)* | *Rye flour (control 2)* | *1* | *2* | *3* | *4* | *5* | *6* |
| rye flour | 50 | 100 | - | - | - | 30.0 | 30.0 | 70.0 |
| wheat flour | 50 | - | 50.0 | 60.0 | 60.0 | 35.0 | 35.0 | - |
| buckwheat flour | - | - | 5.0 | 5.0 | 5.0 | 8.0 | 10.0 | - |
| lentil flour | - | - | 9.3 | 9.0 | 6.8 | 4.2 | 9.3 | 9.3 |
| amaranth flour | - | - | 9.5 | 8.3 | 10.0 | 9.1 | 4.3 | 8.0 |
| barley flour | - | - | 9.8 | - | - | - | - | - |
| flax seeds | - | - | 5.0 | - | - | - | - | - |
| soy flour | - | - | 10.0 | 5.0 | 5.0 | 10.0 | 7.7 | 9.0 |
| sunflower seeds | - | - | - | 10.0 | 6.0 | - | - | - |
| skimmed milk powder | - | - | - | - | 5.0 | - | - | - |
| dry wheat gluten | - | - | 0.6 | 2.2 | 1.8 | 3.0 | 3.0 | 3.0 |
| calcium lactate | - | - | 0.8 | 0.5 | 0.4 | 0.7 | 0.7 | 0.70 |
| The sum of the components of the mixture | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 |
| Biological value,% | 62.0 | 62.5 | 79.0 | 80.0 | 81.2 | 80.5 | 83.3 | 84.7 |
| Drop number, s | 200 | 150 | 210 | 220 | 230 | 230 | 200 | 200 |
| Calcium, mg/100g | 33.0 | 34.0 | 212.8 | 190.6 | 211.8 | 181.2 | 172.4 | 181.1 |
| Phosphorus, mg / 100 g | 186.5 | 189.0 | 280.0 | 251.8 | 288.2 | 240.3 | 226.3 | 234.8 |
| Magnesium, mg / 100 g | 66.5 | 60.0 | 102.7 | 92.16 | 95.15 | 86.85 | 75.46 | 74.12 |
| Cellulose, mg/100 g | 0.3 | 0.2 | 2.86 | 1.20 | 0.94 | 0.67 | 0.81 | 0.69 |
| Protein g / 100 g | 11.2 | 9.9 | 16.69 | 16.13 | 16.96 | 17.04 | 16.95 | 17.18 |
| Lipids, g / 100 g | 1.45 | 1.7 | 3.52 | 6.73 | 4.77 | 1.74 | 1.63 | 1.54 |
| Carbohydrates (mono-disaccha-rides), g / 100 g | 71.9 | 73 | 58.94 | 57.36 | 60.51 | 62.09 | 62.03 | 60.93 |
| The ratio of components | | | | | | | | |
| Calcium | 1.0 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 |
| Phosphorus | 5.7 | 5.6 | 1.32 | 1.32 | 1.36 | 1.33 | 1.31 | 1.30 |
| Magnesium | 2.0 | 1.8 | 0.48 | 0.48 | 0.45 | 0.48 | 0.44 | 0.41 |
| Protein | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 |
| Lipids | 0.1 | 0.2 | 0.21 | 0.42 | 0.28 | 0.10 | 0.10 | 0.09 |
| Carbohydrates | 6.4 | 7.4 | 3.53 | 3.56 | 3.57 | 3.64 | 3.66 | 3.55 |

TABLE III.  QUALITY INDICATORS OF BREAD MADE OF MODEL MIXTURES

| The name of indicators | Characteristics | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | *A mixture of rye and wheat flour (control 1)* | *Rye flour (control 2)* | *1* | *2* | *3* | *4* | *5* | *6* |
| Moisture content, % | 44.3 ±0.5 | 50.0 ±0.5 | 43.5±0.5 | 44.8±0.5 | 44.6±0.5 | 49.5±0.5 | 50.0±0.5 | 50.0±0.5 |
| Specific volume, g / cm | 1.5 ±0.1 | 1.5 ±0.1 | 1.5±0.1 | 1.6±0.1 | 1.6±0.1 | 1.4±0.1 | 1.6±0.1 | 1.5±0.1 |
| Porosity,% | 55.0 ±1.0 | 48.5 ±1.0 | 52.5±1.0 | 55.0±1.0 | 55.5±1.0 | 48.2±1.0 | 48.5±1.0 | 48.5±1.0 |
| Output, % | 148.2±0.5 | 149.9±0.5 | 148.2±0.5 | 147.7±0.5 | 148.4±0.5 | 148.8±0.5 | 147.9±0.5 | 150±0.5 |
| Organoleptic evaluation score, points | 70.5 ±2.0 | 70.5 ±2.0 | 70.5±2.0 | 65.0±2.0 | 70.5±2.0 | 72.0±2.0 | 71.5±2.0 | 69.5±2.0 |

TABLE IV.    CHEMICAL COMPOSITION OF BAKERY PRODUCTS FROM MODEL MIXTURES FOR GERODIETIC NUTRITION

| Name of food substances | Estimated composition | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | *"Nareznoy" loaf* | *"Spassky" rye-wheat bread* | *1* | *2* | *3* | *4* | *5* | *6* |
| Proteins, g | 7.1 | 7.6 | 11.3 | 10.9 | 11.5 | 11.5 | 11.5 | 11.6 |
| Fat, g | 2.7 | 1.0 | 9.0 | 8.7 | 9.2 | 9.2 | 9.2 | 9.3 |
| Carbohydrates, g | 49.1 | 48.6 | 39.8 | 38.8 | 40.9 | 42.0 | 41.9 | 41.2 |
| Dietary fiber, g | 0.2 | 0.2 | 1.9 | 0.8 | 0.6 | 0.5 | 0.5 | 0.5 |
| Calcium, mg | 18.3 | 22.3 | 143.8 | 128.8 | 143.1 | 122.5 | 116.5 | 122.4 |
| Phosphorus, mg | 105.1 | 126.0 | 189.2 | 170.2 | 194.7 | 162.4 | 152.9 | 158.7 |
| Magnesium mg | 41.7 | 44.9 | 69.4 | 62.3 | 64.3 | 58.7 | 51.0 | 50.1 |
| Potassium, mg | 101.7 | 110.5 | 277.4 | 225.9 | 245.1 | 226.2 | 209.0 | 239.4 |
| Sodium, mg | 5.7 | 9.1 | 31.0 | 37.6 | 51.4 | 26.9 | 22.5 | 29.0 |
| Omega-6, g | 2.2 | - | 4.2 | 1.7 | 3.4 | 5.6 | 5.7 | 5.8 |
| Vitamin E, mg | 2.2 | - | 4.1 | 4.4 | 4.8 | 5.3 | 5.3 | 5.2 |
| Energy value, kcal | 249.1 | 210.4 | 256.9 | 249.4 | 263.2 | 267.1 | 266.8 | 265.4 |
| Biological value,% | 66.5 | 62.2 | 79.0 | 80.0 | 81.2 | 80.5 | 83.3 | 84.7 |

Calculation of the food substances chemical composition of developed bakery products (100 g) made of model mixtures for gerodietic nutrition showed that the protein content in them is 1.5-1.6 times higher than that in the control samples - "Nareznoy" loaf and "Spassky" rye-wheat bread. The fat content is 3.2-3.4 times higher, the amount of carbohydrates decreased by 19.3-20.5%, the level of dietary fiber is 2.5-9.5 times higher, the amount of calcium is 2.5-9.7 times higher, phosphorus - 6.7-7.9 times higher, magnesium - 1.5-1.8 times higher, potassium - 1.2-1.7 times higher, sodium - 2.1-2.4 times higher, the content of omega-6 - 0.8-2.6 times higher, vitamin E - 1.9-2.4 times higher.

At the same time the developed bakery products have the increased biological value by 16.8-19.0 %. Proteins, fats and carbohydrates are in a ratio corresponding to the optimal absorption of these nutrients in gerodietic products - 1.0:0.8:3.5. The ratio of mineral substances Ca:Mg:P corresponds to 1:0.6:1.3, and the content of polyunsaturated fatty acids and fiber is at least 1.0 g and 0.2 g, respectively. Calculations show that the developed bakery products correspond well to requirements for products recommended for gerodietic nutrition.

## V. CONCLUSION

A software has been developed for calculating the optimized composition of a model mixture for gerodietic bread according to biological value. This was achieved by mathematical formalization of the problem, taking into account the "drop number" of the mixture as well as the limitations on the optimal content of unconventional raw materials and the number of amino acids that are deficient in bakery products, the ratio of proteins:fats:carbohydrates, Ca:Mg:P, fiber content.

Bakery products prepared from flour mixes modeled with the program complex differed by a higher protein content - 1.5-1.6 times, higher fats content - 3.2-3.4 times, the decrease of carbohydrates amount by 19.3-20.5%. At the same time, the bread had quality indicators that are traditional for this type of product.

As a further task for the study, it would be advisable to develop an integrated large-scale replenishable database of the chemical composition of raw ingredients for broader opportunities in the preparation of flour mixtures. It would also be important to be able to use different optimization criteria when composing flour mixes depending on the tasks.

### REFERENCES

[1] Friedman M. Nutritional value of protein from different food sources. J. Agric. Food Chem. 1996, no. 44 (1), pp. 6-29.

[2] McCollum E.V. A History of Nutrition: The Sequence of Ideas in Nutrition Investigations, Houghton Mifflin, 1957, 451p.

[3] Kodentsova V.M., Vrzhesinskaya O.A., Risnik D.V., Nikityuk D.B., Tutelyan V.A. Micronutrient status of population of the Russian federation and possibility of its correction. State of the problem. Voprosy pitaniia [Problems of Nutrition]. 2017, no. 86 (4) (in Russian).

[4] Bourn, D., Prescott, J. A comparison of the nutritional value, sensory qualities, and food safety of organically and conventionally produced foods. Critical reviews in food science and nutrition, no. 42 (1), p. 1-34.

[5] Tutelyan V.A. The laws of the nutrition science / V.A. Tutelyan // Modern medical technologies. - 2010. - № 4. - p. 98-99.

[6] Lisin, P.A. Computer modeling of multicomponent food products // Food industry, 2008. - № 11. - P. 60–61.

[7] Demina, I.A. Research and calculation of the composition of flour composite mixtures using MS Excel linear programming problems / I.A. Demina, V.S. Kubantsev // KASU Bulletin. - 2009. - №1.

[8] Certificate of official registration of computer program No. 2013613905. Calculation of the chemical composition of food products for catering in

educational institutions. [Text] / Z.I. Harisova, A.T. Zulkarnayeva, E.A. Povargo, I.I. Kurbangaliev. - The date of registration 04/18/2013.

[9] Dvoretsky S.I. Computer modeling and optimization of technological processes and equipment / S.I. Dvoretsky, A.F. Egorov, D.S. Dvoretsky // Textbook. Tambov: Publishing House Tambov state tech. university - 2003. - 224 p.

[10] Energy and protein requirement Report of a joint FAO/WHO ad hoc expert consultation/ WHO tech/ Ref. Ser. №724 – Geneva: WHO, 1985.

[11] Chernikov, M.N. The potential biological value of food proteins and the Mitchell principle [Text] / M.N. Chernikov // Vopr. medical chemistry. - 1989. - V24. - p.9.

[12] Molnar, P. Overall quality index for the evaluation of food products EOQC. Food Section Seminar on quality assurance in the food industry. Budapest, Hungary, 26-25.05.1986, pp. 123-142.

[13] Berezina N. A., Artemov A.V., Nikitin I.A, Zavalishin I.V., Ryazanov A.N. He Use of a Simplex Method with an Artificial basis in Modeling of Flour Mixtures for Bakery Products, (IJACSA) International Journal of Advanced Computer Science and Applications, 2017, vol. 8, no. 12, pp. 338-344.

[14] Methodical recommendations MR 2.3.1.2432-08. Norms of physiological energy and nutrient requirements for various groups of population of the Russian Federation, 2008. - 39 p.

[15] Methodical recommendations on the organization of nutrition of people of elderly and senile age (approved by the Ministry of Health of the USSR on February 19, 1975 No. 1225-75) - 31 p.

[16] Methodical recommendations on catering in institutions (departments) of social services for senior citizens and people with disabilities (Order of the Ministry of Health and Social Development of the Russian Federation of 04.06.2007 N 397), 2007. - 91 p.

[17] Control of raw materials, semi-finished products and finished bakery products: a teaching aid / N.V. Labutina, S.Ya. Koryachkina, N.A. Berezina, E.V. Khmeleva. - M .: DeLi, 2009 - 650 p.

# Quality of Service and Power Consumption Optimization on the IEEE 802.15.4 Pulse Sensor Node based on Internet of Things

Puput Dani Prasetyo Adi[1], Akio Kitagawa[2]

Micro Electronics Research Laboratory (MeRL)

Kanazawa University, Kanazawa, Ishikawa, Japan

*Abstract*—The Purpose of this research is to determine the Quality of Service (QoS) Zigbee or IEEE 802.15.4 sensor Node use the indicators, i.e. the Receiver Signal Strength and PathLoss (attenuation (-dB)) at the time of communication of the sensor node end device to the sensor Router node or Coordinator sensor node (sink). The factor power consumption sensor node is important to maintain the lifetime sensor node, The Sensor data in this research is the pulse sensor. The development of the Wireless Sensor Network communication system is in multi-hop communication, with efforts to obtain low power consumption on each sensor node. This study utilizes the Routing Protocol for Low Power and Lossy Network (RPL) method with position management on the sensor node on DODAGs consequently, that the average power consumption value for each sensor node is low. The benefit of the Sensor node is to send pulse sensor data from various nodes that are interconnected at different distances in multi-hop so that power consumption and Quality of Services (QoS) can be identified from the sensor node. From the research results, the average PathLoss value of IEEE 802.15.4 or Zigbee in free space is obtained by comparing the various simulation values and field experiments at a distance of 50 m at -75.4 dB and the Average Receiver Signal Strength (RSS) with a comparison of Equation and Experiments in the field with parameters The minimum Power Transmitter is 0 dBm and the Power Transmitter is maximum +20 dBm at a distance of 50m at - 66.6 dBm. Therefore, Pulse Sensor data will be displayed on the Web Page and stored in the MySQL database using Raspberry PI 3 as the Internet Gateway.

*Keywords—RPL; RSS; Pathloss; Zigbee; Pulse; DODAGs; IoT*

## I. INTRODUCTION

Heart rate is an important thing that can be one of the first indicators for patients as a doctor uses a stethoscope. The normal range for Pulse sensor between 60-100 beat per minute (bpm), there are two conditions for abnormalities in the heart i.e. Bradycardia and Tachycardia. If the heart conditions are abnormal, consequently, the blood flow will also be disrupted, this disorder is Arrhythmia, i.e. the heartbeat conditions are too slow heartbeat (Bradycardia), too fast heartbeat (Tachycardia) or irregular or variable conditions. The parameters that can be calculated are signal strength (RSSI) and Weakening signal (PathLoss) using ZigBee.

In this research, node sensors are sent using Pulse sensors through the Free Space area, there are multi-hop sensors that are interconnected and therefore, the closest sensors will send pulse data to the edge node, and edge nodes then send pulse sensor data to Raspberry Pi 3 as an internet gateway.

ZigBee or IEEE 802.15.4 Protocol technology was developed for the medical field, Zigbee specification is a light and comfortable prototype for patients, i.e. a pulse detection system using pulse sensors, the ability coordinator node to collect pulse sensors simultaneously without being discovered 3 (tree) nodes [1]. Zigbee has a frequency of 2.4 GHz and a data rate of 250 kbps, this is low bandwidth and is right for the Wireless Sensor Network, Zigbee has the ability to transmit data up to 100 meters.

Perception Layer is a ZigBee Communication System that is used i.e. Star, Tree, and Mesh Topology Fig. 4, from this communication system traffic data, can be analyzed from End Device Node, Router Node and Coordinator Node, this data traffic includes Quality Of Services (QoS) i.e. Throughput, Packet Loss and Delay, and Radio Frequency (RF) capabilities in sending data indicated 21`y Radio Signals Strength (RSS) or Radio Signals Strength Indicator (RSSI) [2, 3, 4, 5, 6]. at previous research testing the performance of Raspberry Pi 3 in capturing radio waves with the Master-Slave system using an RF Bluetooth RN-42 device in sending multi LM35 sensor data, accordingly in this research study the RSSI value is different from the environment can be monitored, from this study obtained the value of precise RSSI [7, 8, 9, 10], furthermore, the next research will be developed of sensor node using a multi-sensor (pulse and blood Pressure sensor). In addition to Zigbee, the wider Wireless Sensor Network device is GSM SIM900A [11] with Ultrasonic sensor data which is often referred to as WMAN.

Furthermore, the development of Wireless Sensor Network technology is applied to UnderWater, which is called the Underwater Wireless Sensor Network (UWSN). In research [12] the analysis was carried out by measuring RF strength with Path Loss Parameters, Velocity of Propagation, absorption loss and the rate of signal loss in a different underwater environment.

The Wireless Sensor Network on this research must be robust and have low power consumption for each node, because the Wireless Sensor Networks to be built are very numerous and interconnected, in this case, at the free Space. Therefore, routing techniques are needed on the Wireless Sensor Network architecture. Routing Protocol for Low Power and Lossy Network is routing that is used on IEEE 802.15.4 or

Zigbee protocol devices, Zigbee in this research is used on sensor nodes which will then function as senders of pulse sensor data. The pulse sensor type used is as in Fig. 1. Therefore, in this research there are two stages of research, 1) examine data transmission on sensor nodes with multi-hop networks. 2) Sending data to the coordinator node to the server at a certain distance with Pathloss and RSSI parameters in free space. The development of this research is towards an Internet of Things, because it is expected that this pulse sensor data can be an indicator for patients who want to check heart rate comfortably without having to go to the hospital, [7] in previous research the authors developed IoT using multi-sensor LM35 and Raspberry Pi 3.

Path Loss (PL) [13, 14, 15] is a weakening of Radio Frequency (RF) wave signals (in this research using ZigBee), which is caused by obstacles during the sending process from transmitter antenna to the receiving antenna, PathLoss (PL) [16, 17] is one of the indicators used to determine whether signal strength is good (dB) [18, 19] the external environment is influenced by Free-Space Path Loss (FSPL). Furthermore, Free-Space Path Loss (FSPL) [20, 21] measured using the Received Signal Strength Indicator (RSSI) technique using a module Zigbee. accordingly, there are two methods used to represent the results of PathLoss, i.e. in the condition of n = 2 or free space, PathLoss is called FSPL (Free Space PathLoss) in outdoor conditions (no space) or no interference, such as other radio waves, in conditions, when the signal is emitted with a known variable from the antenna height which has an impact on reflecting the ground or reflected waves, and the second method is sending data from the edge node to Internet Gateway and pulse data storage in MySQL Database using Python language and HTML and JASON to display graphs on web pages and about IoT Security has not been discussed in this research.

## II. RELATED STUDIES

Muhammad Niswar, Amil Ahmad Ilham, Elyas Palantei, Rhiza S. Sadjad, Andani Ahmad, Ansar Suyuti, Indra Bayu, Zaenab Muslimin, Tadjuddin Waris, and Puput Dani Prasetyo Adi [1], in their research entitled performance evaluation of ZigBee-based wireless sensor networks for patients' monitoring pulse status, this research is included in the Wireless Body Area Network (WBAN) other than Bluetooth, furthermore, from this research the performance of Zigbee as a sender of pulse sensor data is multiplexed or using the star network topology which causes a bottleneck in the coordinator node that functions as the receiver or antenna receiver. Finally, from the 5 sensor nodes that are sent together to the Coordinator node, they cause a loss when sending sensor data with 4 sensor nodes, so a maximum of only 3 sensor nodes or ZED can send data in full or maximum to the Coordinator node with a less than 30 distance meter.

Hana Mujlid, Ivica Kostanic [14], in their research entitled Propagation Path Loss Measurements for Wireless Sensor Networks in Sand and Dust Storms, in this study using empirical measurement data that is using Radio Frequency (RF) Zigbee or IEEE 802.15.4 which works on the 2.4 GHz frequency. The research taken is measuring Path Loss in

conditions of sandstorms and severe sandstorms that have an impact on the loss of propagation pathway from radio signals. From this research, the intercept values were 53 and 60 dB and the slope values were 28 and 37 dB / Dec.

Zhenran Gao, Weijing Li, Yan Zhu, Yongchao Tian, Fangrong Pang, Weixing Cao, and Jun Ni [22], in their research entitled Wireless Channel Propagation Characteristics and Modeling of Research in Rice Field Sensor Networks, in this study testing of Wireless Sensor Network devices in the location of agriculture or agriculture, is precisely in the rice field. So the result of sending signals from the Transmitter antenna to the Receiver antenna is affected by the growth of the rice field. testing was carried out at different altitudes (0.8 m, 1.2 m, 1.6 m, and 2.0 m) consequently, that analysis of the shipment would be found data with different node heights, the models used are free space models and two-Ray models.

Jose Vera-Pérez, David Todolí-Ferrandis, Salvador Santonja-Climent, Javier Silvestre-Blanes, and Víctor Sempere-Payá [23], The Routing Protocol for Low Power and Lossy Network (RPL) is considered for Low energy consumption in the Wireless Sensor Network i.e. IEEE 802.15.4 Network, to support the Internet of Things architecture. The parameters used in determining the quality of the Wireless sensor network i.e. Energy Consumption, Transmission Power, transmission receiver and Packet Size.

## III. METHODOLOGY

### A. Pulse Sensor

The pulse sensor used to detect Heart rate, there are 3 colors of the pin cable, e.g. the red cable for VCC equal to 3.3 - 5 Volts DC and Working Current is 4 mA, black cable for GND and purple color cable for Signal, if the Green LED in the middle of the circle pulse the sensor lights up indicates the Pulse Sensor is ON or active. Accordingly, the Pulse signal sensor is entered on the Analog A0 Pin Arduino pro mini. In this research, Waveform Heartbeat can be shown using an Integrated Processing Development of Environment (IDE) that is integrated with Arduino Pro mini-IDE. This electronic pulse sensor uses a green LED that continuously emits light to the skin. As long as the heart beats, consequently, blood waves flow into the skin capillaries so that light absorption occurs extra during that period. And the process of absorption of light is considered to represent beats heart. Fig. 1 shows three views of a pulse sensor, which is front, back and enlarged, therefore, this is an important part of the pulse sensor that needs to be highlighted in this research [24].

### B. Heart Beat Data Heart Beat Data

Pulse Wave or Photoplethysmogram (PPG) from Pulse sensor is shown in Fig. 2. It seems that a person's heart rate is 80 BPM, which means that the heartbeat of a normal person is a range of 60-100 BPM [1, 24, 25]. Point T is the starting point of forming a PPG signal, and point P is the peak point of the pulse wave amplitude. The percentage value of 25% and 50% shows the magnitude of the Amplitude value generated on the pulse sensor.

Fig. 1.    Pulse Sensor



Fig. 2.    Pulse Sensor Visualizer.

## C.  Heartbeat or Pulse

Arrhythmia is an abnormal condition in the heart, which is a slow pulse (Bradycardia), too fast (Tachycardia) or an ever-changing or irregular pulse. Accordingly, the heart rate is 60-100 BPM (beats per minute) [1]. Bradycardia is a condition of the heart rate below 60 BPM or below Normal, consequently, the condition causes a person to experience fatigue, weakness, dizziness, sweating, and fainting [1, 24, 25]. Conversely, Tachycardia is a condition of the heartbeat that is above 100 BPM or above normal, the condition of the heartbeat is too fast. Consequently, someone who has Tachycardia is dizzy, heart palpitations are very fast, fainting, mild headaches, shortness of breath, fatigue, in extreme situations sufferers can experience unconsciousness and Cardiac arrest. Moreover, the classification of heartbeat conditions is shown in Table I. therefore, this research Bradycardia and Tachycardia are important elements used as indicators on prototypes, where if BPM conditions show certain values such as Bradycardia then the LED indicator will show Yellow, conversely, if Tachycardia will show a red LED.

This pulse sensor data will be sent by the sensor node in the form of packet data (bytes) that can be monitored using X-CTU software. The different distance comparisons will be used as parameters for determining Quality of Service (QoS) from the sensor nodes in this research. The QoS analyzed include Throughput (Kbps), Packet Loss (Packet), Receive Signal Strength (RSSI) and Path Loss (dB) as a result of transmitting sensors data. Accordingly, the test area is Free Space value of Path Loss exponent (n) = 2.

## D.  Zigbee S2c

Zigbee works on the 2.4 GHz frequency and includes the IEEE 802.15.4 protocol, Zigbee is an RF device that uses a low voltage of 2.1-3.6 Volt, generally on the ZigBee module

is 3.3 volts, with Low current RX consumption is 18.8 mA and Tx is 17.4 mA, 9 dB of ZigBee gain and Average RF Power Transmitter (Pt) Zigbee is 14.77 dBm, with a data rate of 250 kbps. Zigbee has many types, e.g. Zigbee S1 and ZigBee S2, ZigBee S1 are compatible RF devices with a point to point and point to multipoint, commonly referred to as star networking, while ZigBee S2 (Fig. 3) is an RF device compatible with mesh networking, usually used as a router. On a mesh networking, ZigBee has the ability to communicate between routers. In this research, ZigBee S2c is compatible with the Mesh Network. ZigBee is a Radio Device that emits radio waves in all directions where there is a request from another radio in a collection of sensor nodes so that the ZigBee antenna belongs to the isotropic type. Furthermore, ZigBee is also included in the Ad-Hoc Networking category because ZigBee is a collection of mobile Nodes in accordance with the ZigBee capacity is compatible with the Mesh Network or only a tree or star network. Fig. 5 Clustering method is used to group nodes based on the closest distance and form clusters and lifts 1 cluster head (CH), therefore, Cluster Head is called Sink which is tasked with sending data of all node members to 1 cluster to Edge Router.

## E.  Routing Protocol for Low Power and Lossy Network

Routing Protocol for Low Power and Lossy Network (RPL) is a type of routing protocol based on the Destination Oriented Acyclic Graph (DODAG). DODAG is a special kind of DAG where each node wants to reach a single destination. The advantages of the Routing Protocol for Low Power and Lossy Network (RPL) are each embedded device or sensor node i.e. IEEE 802.15.4 or Low-Power Wi-Fi in large quantities can be connected effectively and optimally seen from limited power, memory, and processing resources. Accordingly, Fig. 6(a) and (b) are examples of DODAGs built using the Contiki cooja simulator.

The simulation in this research consists of Layer, Protocol and standards that are used to build simulation nodes, in detail can be seen in Table II [26]. The essence of RPL is an attempt to manage the energy or lifetime of a node.

TABLE I.    INDICATOR CLASSIFICATION ON THE PROTOTYPE

| beats per minute (bpm) | Classification | LED Indicator |
|---|---|---|
| >100 | Tachycardia | Red |
| 60  to 100 | Normal | Green |
| <60 | Bradycardia | Yellow |

TABLE II.    PROTOCOL STACK COMMUNICATION

| Layer | Protocol | Standard |
|---|---|---|
| Application | CoAP | IETF RFC 7252 |
| Transport | UDP | IETF RFC 768 |
| Network | IPv6 / RPL | IETF RFC 6550 |
| Adaptation | 6lowpan | IETF RFC 6282 |
| Data Link | IEEE 802.15.4 MAC (CSMA) | IEEE 802.15.4 |
| Radio Duty Cycling | ContikiMAC | - |
| Physical | IEEE 802.15.4 PHY | IEEE 802.15.4 |

Fig. 3.    Zigbee S2c.



(a.Zigbee Star Network)



(b.Zigbee Tree Network)



*(c.* Zigbee Mesh Network*)*

Fig. 4.    Zigbee Topology (Star, Tree Dan Mesh Network).



Fig. 5.    Clustering LowPAN Network.



Fig. 6.    (a)DAG (Directed Acyclic Graph) (b) Destination Oriented DAGs

### F.  Pulse Sensor Node on IoT Architecture

The architecture that will be built is shown in Fig. 7. There are two main analyzes, i.e. transmit sensor nodes on multi-hops which are indicated by setting the router node to EDGE router or sink. The first analyzes are the Quality of Services (QoS) this is an area to calculate the path loss and Receiver Signal Strength parameters. Furthermore, the next stage after the data arrives at the EDGE Router, the analysis is carried out on the Internet Gateway, setting up the MySQL database and display on the WEB so that any devices can see in real-time the Pulse Data sensor [27]. At this step is a need of an analysis of Reability and Energy Consumption in IoT.

### G.  Hardware

Zigbee Sensor nodes in this research shown in the following Fig. 8, built from a pulse, this sensor has been built from Arduino pro mini microcontroller and XBee S2c, this sensor node is ZED (ZigBee End Devices). In this node sensor prototype use 8x2 I2C LCDs and 3 LED colors with each 220 ohm resistor, the indicator espouse to read the patient's pulse condition. Furthermore, Logic Level Converter (LLC) is used as a regulator, LLC consists of several Pin High voltage i.e. 5 Volt and Low Voltage 3.3 Volt Pin. accordingly, which requires a 3.3 volt voltage is an 8x2 I2C and XBee S2c LCD, the pulse sensor can use 3.3 volts or 5 volts, furthermore, the battery used as a power supply is a 3.7 volt 1000 mAh battery type, while the FTDI 232 is used as an Arduino Pro mini board programmer from Arduino Integrated Development Environment (IDE) that uses C++ programming language, the library such as library pulse sensor, LCD and interrupt to calibrate the results of Pulse sensors, Precise material selection will make it convenient to use nodes for patients. Table III shows the hardware functionality of the sensor node.



Fig. 7.    Pulse Sensor Node on IoT Architecture.

Fig. 8. Node Sensor Schematic.

Table III shows the hardware and working functions on the sensor node, there are three types of ZigBee settings, e.g. Zigbee coordinator, Zigbee router and Zigbee end devices each have different functions, voltage usage is 3.3 volts and 5 volts on a sensor node, therefore, the voltage will be distributed, this is where the regulator works by using the High Voltage (HV) to Low Voltage (LV) and Low Voltage (LV) to High Voltage (HV). On Arduino pro mini bootloader programming is needed unlike the Arduino Uno version, so FTDI 232 is needed as a bootloader.

TABLE III. NODE HARDWARE

| No | Hardware | Notes |
|----|----------|-------|
| 1 | Microcontroller Arduino Pro mini | Processor, ADC, Data Serial Communication With a Base Station |
| 2 | XBee S2c End Device | Wireless sensor Network type to sending Pulse sensor data to Coordinator node |
| 3 | XBee S2c Coordinator | Wireless sensor Network type to receive Pulse sensor data from ZED or ZR to Base Station |
| 4 | XBee S2c Router | Wireless sensor Network type to sending Pulse sensor data to Coordinator node from ZED and Communicate between each Router at Mesh network |
| 5 | XBee USB Adapter | Setting Device to configuration the ZigBee S2c |
| 6 | Pulse Sensor | As Pulse Sensor data from patients |
| 7 | LCD 16X2 atau LCD 8X2 I2C | As Indicator |
| 8 | FTDI 232 | As Programming device from USB and Arduino IDE to Arduino Pro mini |
| 9 | Regulator 3.3 Volt to 5 Volt | Change the voltage to HV (High Voltage 5Volt ) or LV (Low Voltage 3.3 Volt) |
| 10 | Battery 3.7 Volt 1000 mAh | Supply power to the sensor node |

Actually the role is not important after the program upload process on Arduino pro mini, so it can be separated and the sensor node only uses the function of Arduino Pro Mini, but FTDI 232 is still useful besides being a bootloader, which is to supply 3.3 volts of power on ZigBee and 8X2 LCD. The role of supplying power to the node of this sensor is a 3.7 Volt 1000 mAh battery. The display of the prototype can be seen in Fig. 9. Accordingly, the position of the node or ZigBee end device is on it is necessary to see the color classification, there are three conditions, e.g. normal, represented in green (60-100 BPM), bradycardia (BPM >44 & <60 BPM) or (BPM >100 & <116) represented in yellow or a minor condition and Major or tachycardia (BPM<44 & >116 BPM) are represented in red, therefore, classification can take place easily from pulse indications.

Fig. 10 is a sensor node that will be used in mesh networks using the Routing Protocol for Low Power and Lossy Network (RPL). This RPL is done using the simulation software with the IEEE 802.15.4 protocol.



Fig. 9. Node Sensor.



Fig. 10. ZigBee end Devices and ZigBee Coordinator Connectivity.

Fig. 11. Battery Lifetime Zigbee Pulse Sensor Node.

With the specifications of the 3.7 Volt Battery node 1000 mAh, and load calculations e.g. XBee = 27 mA (TX Mode), LCD 8x2 = 2 mA, Arduino Pro mini 40 mA, FTDI232 is 50 mA and Pulse 5 mA, then total load is 124 mA. With a 5 Volt Voltage, then the Power is Voltage x Current = 5x0.124 is 0.62 Watt. Then I = 0.62 W / 3.7 Volt = 0.167 Ampere. Therefore, the power Consumption is 1000 mAH or 1AH / 0.167 A = 6 hours, this measurement is shown in Fig. 11.

### H. Equations of Received Signal Strength (dBm) and Path Loss (-dB) of IEEE 802.15.4 Zigbee

Zigbee is a Radio Frequency (RF) device of 2.4 GHz or 2400 MHz, so the wavelength ($\lambda$) of zigbee is 0.125 meters, this value is obtained from $\lambda = c / f$, where c is the speed of light which is $3 \times 10^8$. Power receiver (dB) will always weaken if distance (d) gets farther away equation (8), then the initial Power Receiver (Pr0) when d = 1 m, can be seen in equation (1), equation (2).[3]

The first test parameter is RSS (dBm) in the Free Space state, where the value of the exponent (n) of Free Space is 2. The value of the Receiver Signal Strength (RSS) results from the comparison of equation 10 and at the time of experiment or measurement in the field, while attenuation generated from equations 6, 7 and 8.

$$P0 = Ptx. Gt. Gr. \quad (\lambda/4\ \pi)^2 \tag{1}$$

$$Pr = \frac{P0}{d^2} \tag{2}$$

A logarithmic form (decibel scale), can be seen in equation (3).

$$10 \log_{10} Pr = 10 \log_{10} P0 - 20 \log_{10} d \tag{3}$$

*L Free Space* can generally be calculated according to the equation (4).

$$L\ Free\ Space = - (20 \log d + 20 \log f - 27,5) \tag{4}$$

C is a constant value that is responsible for wireless attenuation in the Free Space condition which is worth 27.5. So that if distance (d) gets bigger, the value of L Free Space (-dB) will increase. So that the frequency of Zigbee is 2.45 GHz, from equation 4 becomes equation (5).

$$L\ Free\ Space\_2450\ MHz = - (20 \log d + 40.3) \tag{5}$$

In Zigbee Datasheet Zigbee, attenuation in free space has a different value in equation (6).

$$L\ Free\ Space\_2450\ MHz = - (33 \log d/8 + 58.5) \tag{6}$$

In Fig. 15 shows the attenuation of Zigbee based on the value of d which is getting bigger, so that the value of L Free Space (-dB) is getting bigger as well. In fact, on experiment 1 and experiment 2 Attenuation is affected by noise when transferring data on the Zigbee transmitter sensor to zigbee receiver sensor so that the data is not so smooth, but in outline decreases (-dB) or PathLoss value gets bigger. In Fig. 10 at a distance of 1 m Pathloss shows an average value of -35 dB and at a distance of 50 m Pathloss average value of - 75.4 dB. So that after the attenuation value (-dB) has been found, the Pr value (Power Receiver (dBm) can be determined by equation (7).

$$Pr = Pt + Gr + Gt + L \tag{7}$$

Assuming the antenna gain for both transmitter and receiver is 0 dB (1 mW) on the Zigbee specification, on free space propagation is formulated in equation (8). And for indoor propagation condition is shown in equation (9).

$$Pr = Pt - (33 \log d/8 + 58.5) \tag{8}$$

$$Pr = Pt - (40 \log d/8 + 50.3) \tag{9}$$

The specification of the Zigbee Power Transmitter datasheet shows the value of 0 dBm for Low Power and + 20 dBm for High Power Transmitters so that you can make a comparison chart. An example of sending data from the TX Zigbee pulse sensor Node to this RX Zigbee sensor node is shown in Fig. 12.



Fig. 12. Zigbee Transmitter and Receiver at Free Space Propagation.

*I. DODAGs of Sensor Nodes Power Consumption Analyze*

In this research using Routing Protocol for Low Power and Lossy Network (RPL) analysis [7], so that the value of traffic sensor nodes can be generated in detail. During the simulation, the Sensor node or IEEE 802.15.4 parameter needs to be included in the simulated parameters, so the results can be compared with the experiment. In accordance with the Zigbee S2c specifications used in this research, the ZigBee 2.4 GHz frequency, the ZigBee data rate is 250 kbps, and the transmission rate is 10-100 meters on Line-of-sight area. in this simulation, the transmission rate (Tx) range is 50 meters. Please note that when the sensor data reaches the sink, the analysis process continues at the next point. At the sensor node, when on the Tx range other sensor nodes are worth 100%, meaning that they are included in the sink range or other sensor nodes that are the router, and the maximum distance is 50 meters. If the distance is > 50 meters, the sensor node is not detected and cannot be analyzed, consequently, the sensor node is required to enter in the range of 50 meters. Therefore, in the simulation of sensor nodes that enters the range of 50 meters indicated by a value of 100%. The Power Consumption test parameters on the sensor node used include CPU, LPM, Radio Transmit (Tx) and Radio Listen (Rx).

---

***Header File Code to represent the Value of Power Consumpsion each sensor node***

```
public static final long TICKS_PER_SECOND = 4096L;
public static final double VOLTAGE = 3;
public static final double POWER_CPU = 1.800 * VOLTAGE;      /* mW */
public static final double POWER_LPM = 0.0545 * VOLTAGE;      /* mW */
public static final double POWER_TRANSMIT = 17.7 * VOLTAGE;   /* mW */
public static final double POWER_LISTEN = 20.0 * VOLTAGE;     /* mW */
```

*-------Power Consumption Spesification of Sensor nodes --------- (10)*

---

***Equation the energy Consumption of Sensor node and relationship with time***

*CPU energy:* $cp = c*1.8/tm$
*LPM energy:* $lp = l*0.0545/tm$
*Transmit energy:* $lt = t*17.7/tm$
*Listen energy:* $lr = r*20/tm$
*Total energy:* $n = cp+lp+lt+lr$
$tm = c+l$
where:
*tm* is the total time,
*c* is the time that the CPU was used.
*l* is the the time that the sensor was in Low Power Mode (LPM)
*t* is the transmit time and
*r* is the Listen time

*-------Energy Consumption equation of Sensor nodes --------- (11)*

---

In this research, the approach is carried out if there are many nodes that are interconnected with node specifications that have specifications in equation (10) and equation (11). With the DAGs and DODAGs method by setting the Sink number in the Clustering system, furthermore, how does the power consumptions affect and the analysis of changes in the position of nodes on the Routing Protocol for Low Power and Lossy Network (RPL) topology, the aim is to get a low Power Consumption on each node, furthermore, that the Average Average Power Consumption (mW) is obtained. Consequently, the test is done by setting the node position in

DODAG's Directed Acyclic Graph (DAG) and Destination (DAG).

In Fig. 13, the router node (node 3) has a different task, namely as a router node for all nodes joined in the DAGs topology. In this case, the Router node must have more Power to receive all data multiplexing from the data end nodes.

While in Fig. 14, divide the Router node into two nodes, namely Router node (2) and node (11) with the same number of nodes in Fig. 13 and Fig. 14, i.e. 10 sender nodes and 1 node coordinator (sink). This decreases the Power needed by the Router node to send sensor data to the coordinator node (sink). Fig. 17 and Fig. 18 show that Power Consumption is needed for each node, so that conclusions from 2 DODAG topology can be obtained which is the most appropriate to be applied when sending pulse sensor nodes in this research.

In Fig. 17 and Fig. 18, there is a significant difference between Topology DODAG 1 (Fig. 13) and DODAG 2 (Fig. 14), this difference is seen from the energy consumption required by the sensor node in communicating and sending data to other nodes (neighboard). In Fig. 17, the power consumption required by node 3 is much higher than the other 10 nodes, this is because node 3 is the directed node or the router node that sends all nodes to the sink or coordinator node. Consequently, Average Power consumption in Fig. 17 can be seen in Fig. 19, and the Average in the Power Consumption of 2 Router node can be seen in Fig. 20.



Fig. 13. Multi DAGs on DODAG.



Fig. 14. DAGs with a Division Router.

## IV. RESULTS AND ANALYSIS

Fig. 15 shows the PathLoss (-dB) value for Free Space compared to Equation 4, 5 and 6 and 9 with - (40 log d / 8 + 50.3) this is indoor position. The average of the 4 equations at a distance of 1m is -30.9 dBm and at a distance of 50m is - 78.9 dBm.

Fig. 16 shows the value of the Power Receiver (dBm) based on the distance and value of the Power Transmitter if using the minimum Pr value (0 dBm) and Pr value (+20 dBm) at different distances. Fig. 16 is Comparing Graph from the Calculation value in equation (7) and the Power Receiver (dBm) generated from the experiment in the field using the Digi X-CTU Zigbee software.

Fig. 19 and Fig. 20 are Average Power Consumption on 1 Router node and 2 Router nodes; sampling is for 10 minutes in both experiments on 1 and 2 Router nodes. From the experiment it can be concluded that using 1 Router Node will use Power Consumption greater than using 2 router nodes, with 1 router node that handles 10 sender nodes requires 1.35 mW Power Consumption while with 2 routers it can be reduced to an average 1 mW Power Consumption.



Fig. 15. ZigBee PathLoss (*-dB*) on Free Space Area.



Fig. 16. Received Signal Strength of Zigbee (*dBm*) on Free Space Area.



Fig. 17. Power Consumption of 1 Router Node.



Fig. 18. Power Consumption of 2 Router Node.



Fig. 19. Average Power Consumption of 1 Router Node.

Fig. 20. Average Power Consumption of 2 Router Node.

In Fig. 21, Nodes 22, 23, 24 and 26 are Router nodes whose task is to send data from end devices to communicate with the coordinator node.

The conclusion that can be generated from analysis at Fig. 21 is that the more sensor nodes are in the position of end devices with a large number of hops that will not affect power consumption, however a power consumption is influenced by the number of end nodes handled by the node, for example, nodes 15 and 26 are nodes most in the DODAGs topology which has the highest number of nodes than the node handled by other Router nodes, see in Fig. 22. What is different is that node 15 has a number of hops 2, meaning that its position is not as a router but as a node rank = 2, see in Fig. 23.

In this section, the Pulse sensor data is in the Node Coordinator and is ready to be transferred to the MySQL database via the Python Language platform on the Internet Gateway. The internet gateway used is Raspberry Pi 3. Python language sends pulse sensor data to the MySQL database, then the data is sent from the MySQL database to HTML and JASON so that the pulse sensor data can be seen in real-time on the Internet, Fig. 24, according to the Internet of things architecture in Fig. 7.



Fig. 21. DAGs with a 4 Division Router Node.



Fig. 22. Average Power Consumption of 26 Nodes.



Fig. 23. Network Hops.



Fig. 24. Realtime Data Pulse (BPM) at Web Page HTML JASON.

## V. CONCLUSIONS

The Quality of Services (QoS) from a Zigbee pulse sensor node can be seen from the parameters i.e. Pathloss (Attenuation (-dB)) and Power Receiver (Receiver Power Strength (dBm)), in this case, ZigBee PathLoss (-dB) on Free Space area will be represented by 1-50m distances, consequently, the average of the 4 equations at a distance of

1m is -30.9 dBm and at a distance of 50m is -78.9 dBm. Furthermore, the average Power Receiver (dBm) on Free Space in 1-15m distances for Zigbee is -66,671 dBm. In the IEEE 802.15.4 routing sensor, the sensor node shows with 1 router node and the 10 sender nodes handle require 1.35mW Power Consumption, this can be concluded that the more router nodes the sensor is Power Consumption will be lower, this is because the load given is not based on just one router node. Sending data from Sink node or EDGE node to Internet Gateway has been successful. Sensor Pulse Data is successfully stored in the MySQL database and displayed on the Web in real-time.

## VI. DISCUSSION AND SUGGESTIONS

This research will be applied to the Monitoring Health Patient Status, as long as there is internet connection, the Patient data e.g. Pulse will be very easy to obtain. In this research discuss about the performance of Zigbee devices at a distance of 50 m with the Coordinator node or receiver connected to the Raspberry Pi 3 as an internet gateway. Analysis of future research will also discuss the comparison of topology in ZigBee in the free space area so that the data pulse will be analyzed based on details of QoS based on the ZigBee topology and its impact on the ZigBee coordinator node. Furthermore, the data pulse for the ZigBee coordinator node will be sent via Raspberry Pi 3 as an internet gateway using Python Programming. This research still has to be developed with the analysis of the Internet of Things Protocol and Security Methods of the Internet of Things. Comparison of 6LOWPANs, Zigbee and LoRAWAN can be done in the next research. Development of the Routing Algorithm needs to be done to obtain a low Energy Consumption value on the Sensor node.

## ACKNOWLEDGMENT

## REFERENCES

[1] Muhammad Niswar, Amil Ahmad Ilham, Elyas Palantei, Rhiza S. Sadjad, Andani Ahmad, Ansar Suyuti, Indrabayu, Zaenab Muslimin, Tadjuddin Waris, Puput Dani Prasetyo Adi, Performance evaluation of ZigBee-based wireless sensor network for monitoring patients' pulse status, 2013 International Conference on Information Technology and Electrical Engineering (ICITEE) DOI: doi/10.1109/ICITEED.2013.6676255

[2] Jinze Du, Jean-François Diouris and Yide Wang, A RSSI-based parameter tracking strategy for constrained position localization, Du et al. EURASIP Journal on Advances in Signal Processing (2017) 2017:77, DOI 10.1186/s13634-017-0512-x

[3] Jiyan Huang, Peng Liu, Wei Lin and Guan Gui, RSS-Based Method for Sensor Localization with Unknown Transmit Power and Uncertainty in Path Loss Exponent, Sensors 2016, 16, 1452; doi:10.3390/s16091452

[4] Jungang Zheng, Yue Liu, Xufeng Fan and Feng Li, The Study of RSSI in Wireless Sensor Networks, Advances in Intelligent Systems Research, volume 133, 2nd International Conference on Artificial Intelligence and Industrial Engineering (AIIE2016) Copyright © 2016, the Authors. Published by Atlantis Press.

[5] P. K. Dutta, O. P. Mishra, and M. K. Naskar, Analysis of dynamic path loss based on the RSSI model for rupture location analysis in underground wireless sensor networks and its implications for Earthquake Early Warning System (EEWS), International Journal of Automation and Smart Technology (AUSMT), DOI: 10.5875/ausmt.v5i3.858

[6] Pooyan Abouzar, David G. Michelson, and Maziyar Hamdi, RSSI-Based Distributed Self-Localization for Wireless Sensor Networks used in Precision Agriculture, arXiv:1509.02400v1 [cs.DC] 21 Aug 2015, https://arxiv.org/pdf/1509.02400.pdf

[7] Puput Dani Prasetyo Adi and Akio Kitagawa, "Performance Evaluation WPAN of RN-42 Bluetooth based (802.15.1) for Sending the Multi-Sensor LM35 Data Temperature and RaspBerry Pi 3 Model B for the Database and Internet Gateway" International Journal of Advanced Computer Science and Applications (IJACSA), 9(12), 2018. DOI: dx.doi/10.14569/IJACSA.2018.091285

[8] Pushan Dutta, O.P.Mishra, M.K.Naskar, Analysis of dynamic path loss based on the RSSI model for rupture location analysis in underground wireless sensor networks and its implications for Earthquake Early Warning System (EEWS), September 2015 International Journal of Automation and Smart Technology 5(3), DOI: 10.5875/ausmt.v5i3.858

[9] Ranjan Kumar Mahapatra, N. S. V. Shet, Localization Based on RSSI Exploiting Gaussian and Averaging Filter in Wireless Sensor Network, Arabian Journal for Science and Engineering, August 2018, Volume 43, Issue 8, pp 4145–4159, DOI: doi/10.1007/s13369-017-2826-2

[10] Yu Xiaoqing, Zhang Zenglin, Han Wenting, Experiment Measurements of RSSI for Wireless Underground Sensor Network in Soil, IAENG International Journal of Computer Science, 45:2, IJCS_45_2_02, Advance online publication: 28 May 2018

[11] Puput Dani Prasetyo Adi and Rahman Arifuddin, Design of Tsunami Detector Based Sort Message Service Using Arduino and SIM900A to GSM/GPRS Module, JEEMECS (Journal of Electrical Engineering, Mechatronic and Computer Science) Volume 1, No.1. 2018, DOI: doi/10.26905/jeemecs.v1i1.1982

[12] Umair Mujtaba Qureshi, Faisal Karim Shaikh, Zuneera Aziz, Syed M. Zafi S. Shah, Adil A. Sheikh, Emad Felemban and Saad Bin Qaisar, RF Path and Absorption Loss Estimation for Underwater Wireless Sensor Networks in Different Water Environments, Sensors 2016, 16(6), 890; https://doi.org/10.3390/s16060890

[13] Daihua Wang, Linli Song, Xiangshan Kong, and Zhijie Zhang, Near-Ground Path Loss Measurements and Modeling for Wireless Sensor Networks at 2.4 GHz, Hindawi Publishing Corporation International Journal of Distributed Sensor Networks Volume 2012, Article ID 969712, 10 pages doi:10.1155/2012/969712

[14] Hana Mujlid, Ivica Kostanic, Propagation Path Loss Measurements for Wireless Sensor Networks in Sand and Dust Storms, Frontiers in Sensors (FS) Volume 4, 2016, DOI: doi/10.14355/fs.2016.04.004 www.seipub.org/fs

[15] Hristos T. Anastassiu, Stavros Vougioukas,Theodoros Fronimos, Christian Regen, Loukas Petrou, Manuela Zude 4 and Jana Käthner, A Computational Model for Path Loss in Wireless Sensor Networks in Orchard Environments, Sensors 2014, 14, 5118-5135; doi:10.3390/s140305118

[16] J. Miranda, R. Abrishambaf, T. Gomes, P. Gonçalves, J. Cabral, A. Tavares and J. Monteiro, Path Loss Exponent Analysis in Wireless Sensor Networks: Experimental Evaluation, https://www.researchgate.net/publication/256733582, Conference Paper July 2013, DOI: 10.1109/INDIN.2013.6622857

[17] Michael Cheffena and Marshed Mohamed, Empirical Path Loss Models for Wireless Sensor Network Deployment in Snowy Environments, IEEE Antennas and Wireless Propagation Letter (Volume:16), 11 September 2017, DOI: 10.1109/LAWP.2017.2751079

[18] Naseer Sabri, S A Aljunid, M S Salim, R Kamaruddin, R B Ahmad, M F Malek, Path Loss Analysis of WSN Wave Propagation in Vegetation, Journal of Physics: Conference Series 423 (2013) 012063, doi:10.1088/1742-6596/423/1/012063, ScieTech 2013, IOP Publishing.

[19] Sinant Kurt and Bulent Tavli, Path Loss Modeling for Wireless Sensor Network : Review of Models and Comparative Evaluations, IEEE Antennas and Propagation Magazines, July 2016, DOI:10.1109/MAP.2016.2630035

[20] Tajudeen O. Olasupo, Carlos E. Otero, Kehinde O. Olasupo, Ivica Kostanic, Empirical Path Loss Models for Wireless Sensor Network Deployments in Short and Tall Natural Grass Environments, IEEE Transactions on Antennas & Propagation, 2016, Manuscript ID is AP1512-1931.R2, DOI 10.1109/TAP.2016.2583507, IEEE Transactions on Antennas and Propagation

[21] Xiaoqing Yu, Wenting Han, Zenglin Zhang, Path Loss Estimation for Wireless Underground Sensor Network in Agricultural Application, Agric Res (2017) 6: 97. DOI: doi/10.1007/s40003-016-0239-1

[22] Zhenran Gao, Weijing Li, Yan Zhu, Yongchao Tian, Fangrong Pang, Weixing Cao, and Jun Ni, Wireless Channel Propagation Characteristics and Modeling of Research in Rice Field Sensor Networks,Sensors 2018, 18, 3116; doi/10.3390/s18093116

[23] Jose Vera-Pérez, David Todolí-Ferrandis, Salvador Santonja-Climent, Javier Silvestre-Blanes, and Víctor Sempere-Payá, "A Joining Procedure and Synchronization for TSCH-RPL Wireless Sensor Networks", Sensors 2018, 18, 3556; doi:10.3390/s18103556

[24] Puput Dani Prasetyo Adi, Analisis kinerja jaringan sensor nirkabel untuk monitoring denyut nadi pasien, April 2018, DOI: 10.13140/RG.2.2.29145.83040

[25] Pradana, H., & Adi, P. (2018). Monitoring Detak Jantung Dan Sistem Implementasi Telemetri Pada Pelaksanaan Lari. SinarFe7, 1(2), 552-556. Retrieved from http://ejournal.fortei7.org/index.php /SinarFe7 /article/view/187

[26] Widya Cahyadi, Muhammad Arief Wahyudi, dan Catur Suko Sarwono, "Analisis Perbandingan Konsumsi Energi dan Masa Hidup Jaringan pada Protokol LEACH, HEED, dan PEGASIS di Wireless Sensor Network", Jurnal Rekayasa Elektrika, VOLUME 14 NOMOR 2, AGUSTUS 2018

[27] Muhammad Ateeq, Farruh Ishmanov, Muhammad Khalil Afzal, and Muhammad Naeem" Multi-Parametric Analysis of Reliability and Energy Consumption in IoT: A Deep Learning Approach", Sensors 2019, 19, 309; doi:10.3390/s19020309

# Segmentation of Touching Arabic Characters in Handwritten Documents by Overlapping Set Theory and Contour Tracing

Inam Ullah[1], Mohd Sanusi Azmi[2]
Mohamad Ishak Desa[3]
Faculty of information technology and Communication
Universiti Teknikal Malaysia (UTeM)
Melaka, Malaysia

Yazan M. Alomari[4]
Department of Management Information Systems College of
Applied Studies and Community Services
Imam Abdulrahman Bin Faisal University
Dammam, Saudi Arabia

*Abstract*—Segmentation of handwritten words into characters is one of the challenging problem in the field of OCR. In presence of touching characters, make this problem more difficult and challenging. There are many obstacles/challenges in segmentation of touching Arabic handwritten text. Although researches are busy in solving the problem of segmentation of these touching characters but still there exist unsolved problems of segmentation of touching offline Arabic handwritten characters. This is due to large variety of characters and their shapes. So in this research, a new method for segmentation of touching Arabic Handwritten character has been developed. The main idea of the proposed method is to segment the touching characters by identifying the touching point by overlapping set theory and ending points of the Arabic word by applying some standard morphology operation methods. After identifying all the points, segmentation method is applied to trace the boundaries of characters to separate these touching characters. Experiments were conducted on touching characters taken from different data sets. The results show the accuracy of the proposed method.

*Keywords*—*Offline handwritten characters; touching characters; segmentation; overlapping set theory; morphological operation*

## I. INTRODUCTION

Modern age, also called the age of information technology because the computer has importance in every field of life. Computer is considered as essential part of human life. Although it is true that computer have not much intelligent compare to human. Human can recognize any type of text image from historical and degraded documents lying in the libraries while computer can't directly understand these text images [1]. Many efforts are required to convert these historical documents to machine understandable format [2][3] because it is not sufficient to store the information in image format [4]. Researchers are busy in developing new algorithms for segmentation of touching Arabic characters, but still this is long standing problem for conversion of handwritten images to electronic form [5][6]. The problem becomes more serious, when dealing with touching handwritten Arabic words because still there is a gap between human and machine abilities in reading handwriting text under noisy conditions especially for overlapped Arabic manuscripts. This due to the nature of font and style of the Arabic characters, which is written from right to left and is always cursive in both machine printed and handwritten text [7][8][9]. Numerous attempts have been made for the recognition/segmentation of overlapped words in Arabic and other languages as well but these overlapped characters still exist gaps [10]. All these efforts emerge the idea of Optical Character Recognition (OCR).

OCR is a technology that is used to convert paper scanned or other types of images to editable format [11][12]. But before converting these images to editable images it needs image segmentation. Image segmentation is one of the important step in OCR because segmentation subdivides an image and distinguish the area of interest and ignore unwanted information [13]. Although image segmentation is not directly related to recognize the segmented images but they are closely related to each other. Image segmentation is important basis for image recognition [14]. If image segmentation is accurate recognition rate is high otherwise recognition ratio is low.

Segmentation which is used to break the text into lines, words and characters of handwritten text is still a challenging task because handwriting is natural and differs from person to person; therefore, many researchers are investigating solutions to solve the problem and some of them have made remarkable achievements, still more research is needed to improve the performance of already developed systems. To discuss all developed methods in this paper is not possible but research done by, address the issues of touching Arabic handwritten characters.

The rest of the paper is organized as follows. Section II covers basic background about the properties of Arabic language with touching types in the Arabic handwritten documents. Section III describes the related works. Algorithm details are in Section IV. Experimental results are reported in Section V. Conclusion and future work is discussed in Section VI.

## II. PROPERTIES OF ARABIC LANGUAGE

### A. Arabic Alphabets

In Arabic language consists of 28 alphabets[15][16] shown in Fig. 1.

Fig. 1.    Arabic Alphabets [15] [16].

### B. Shape of Characetrs in Words

Shape of alphabets according to location is shown in Table I. Some letters have the same shape but numbers, location (above and below) of dots and strokes differ one alphabets from other e.g. Ba ( ب ), Ta ( ت ), and Tha (ث).

### C. Touching Character and its Types

Characters may be joined with the character of other word or with in the same words to form simple touching [18] as shown in Fig. 2.

TABLE I.        ARABIC CHARACTER FORMS [17]

| Characters | Isolated | Beginning | Middle | End |
|---|---|---|---|---|
| Alif | ا | - | - | ـا |
| Ba | ب | بـ | ـبـ | ـب |
| Ta | ت | تـ | ـتـ | ـت |
| Tha | ث | ثـ | ـثـ | ـث |
| Jeem | ج | جـ | ـجـ | ـج |
| Ha | ح | حـ | ـحـ | ـح |
| Kha | خ | خـ | ـخـ | ـخ |
| Dal | د | - | - | ـد |
| Thal | ذ | - | - | ـذ |
| Ra | ر | - | - | ـر |
| Zai | ز | - | - | ـز |
| Seen | س | سـ | ـسـ | ـس |
| Sheen | ش | شـ | ـشـ | ـش |
| Swad | ص | صـ | ـصـ | ـص |
| Dwad | ض | ضـ | ـضـ | ـض |
| Tah | ط | طـ | ـطـ | ـط |
| Dha(Zwa) | ظ | ظـ | ـظـ | ـظ |
| Ain | ع | عـ | ـعـ | ـع |
| Ghain | غ | غـ | ـغـ | ـغ |
| Fa | ف | فـ | ـفـ | ـف |
| Qaf | ق | قـ | ـقـ | ـق |
| Kaf | ك | كـ | ـكـ | ـك |
| Lam | ل | لـ | ـلـ | ـل |
| Meem | م | مـ | ـمـ | ـم |
| Noon | ن | نـ | ـنـ | ـن |
| Ha | ه | هـ | ـهـ | ـه |
| Waw | و | و | ـو | ـو |
| Ya | ي | يـ | ـيـ | ـي |



Fig. 2.    Touching Character.

Characters can be joined in such a way to form even complicated touching, such as writing in Arabic calligraphy shown in Fig. 3.

Touching of characters taken from handwritten Arabic AHDB dataset is shown in Fig. 4.

Thus, the possible types of touching that normally exists between characters are shown in Table II.



Fig. 3.    Complicated Touching Characters.



Fig. 4.    Touching Characters in AHDB Dataset.

TABLE II.        GENERAL TOUCHING TYPES [19]

| Type | Letters | Sample |
|---|---|---|
| A | Top:[ر, ز, س, ش, ص, ض, ن, ق, و, ى, ي]<br>Bottom:        [ا, ط, ظ, ك, ل] |  |
| B | Top:        [ر, ز, م, و]<br>Bottom:   [ص, ض, ة] |  |
| C | Top:   [ج, ح, خ, ع, غ]<br>Bottom:  [ا, ط, ظ, ك, ل] |  |
| D | Top:        [ج, ح, خ, ع, غ]<br>Bottom:        [ه] |  |

### III. RELATED WORK

By exploring the published literature related to segmentation of touching handwritten Arabic characters, there are number of methods proposed for segmentation of handwritten and printed touching text. Some of them are.

In [19], they proposed a new method on the basis of morphological analysis for segmentation of touching lines in Arabic handwritten document. Four types of touching types in Arabic document. Image is converted into skeleton image and the curved is traced by angular variance starting from the starting point. The purpose is to trace the skeleton image in right direction. This method is working for selected touching but not for all touching types because of varieties in handwritten Arabic characters.

In [20], they proposed a new method called template based segmentation method for segmentation of touching characters. In this method created a dictionary file, which contains template of all possible touching with necessary detail. Then compare with input image and select template from the dictionary file. This method works well some image but can't segment all touching characters.

Problems identified in this method:

*1)* Tedious method because of creation of template file.

*2)* Template dependent method, will not segment all touching characters.

*3)* Time consuming, because for simple touching characters will search whole dictionary file.

In [21], proposed a new segmentation method for segmentation of touching handwritten digits. First finding the end-points by applying some standard method. Now by tracing the boundary from end-point collecting all co-ordinates and segmentation point. Thus segmented the touching digits with this recognition ratio 71.43%. But problems identified in the proposed method are given below:

*1)* Proposed method is only for segmentation of numeric digit not used for character data.

*2)* Proposed method is only for simple touching for multiple touching occurs over-segmentation.

In [22], proposed a new segmentation technique for solving the problem touching handwritten Arabic characters, touching is between two characters in same word or other line. This method is template based and these touching images are compare with the template file already create before. Problems in the proposed method:

*1)* Performance of proposed method depended on template.

*2)* As in handwritten text there is a lot of variation even between same writer. Any incorrect template selection affects the recognition rate.

*3)* Sometime produce broken characters during segmentation.

In [23], proposed method is basically for segmentation and recognition of Arabic touching characters. The concept is to trace the boundary of character. Normalize center of gravity of

region. Then calculate minimum distance and apply horizontal distance between connected region. Segment the connected characters. Problem identified:

*1)* Segmentation of ligature is not explained neither give the results or comparison to clear the methodology.

*2)* Authors claimed the small modification can be used for segmentation of various documents including Othman script of Al-Quran. But Quran of Othman Script have multiple touching points and also have some broken characters are possible, especially in handwritten script. No solution is given for merging of broken characters

In [14], proposed method for segmentation of touching handwritten and printed Latin characters. Basically this method is a combination of three already developed methods taken from the literature. During database selection for testing, highlighted the problem of unavailability of standard datasets. Although they claim high success segmentation ratio but still this method is only applicable to certain situation or document, which is the drawback of the proposed method.

### IV. PROPOSED METHOD

In this section, proposed method for segmentation of touching handwritten Arabic character with mathematical background is explained. The model of the proposed method is illustrated in Fig. 5.

In the first stage of the model, input image is converted to binary image and also to enhance the quality of the image by removing unwanted information.

Step 1. Find all endpoints of touching image. In Fig. 6 Four endpoints in the image $E_1$, $E_2$, $E_3$ and $E_4$.

Step 2. Find coordinate between any two Endpoints by tracing boundary of image. Suppose between $E_1$, $E_2$ and $E_3$, $E_4$. Thus
Set A = {Coordinates between $E_1$, $E_2$ }
Set B = {Coordinates between $E_3$, $E_4$ }

Step 3. Apply overlapping set theory on Set A and B
Set C = Set A $\cap$ Set B
If IsEmpy Set C
No touching character
Else
Touching character and element of Set C is the touching point.

Step 4. END

In the second stage of the model, End-points, touching and neighboring point are detected in the input image. For endpoints and neighboring points of thinned image applied standard method but for touching point applied overlapping set theory. Steps of overlapping set theory are to find the junction or touching point.

After finding the touching point, next is the neighbor point in all direction of the image near touching point. Touching and neighboring points are shown in Fig. 6 and Fig. 7 below.

Fig. 5.    Proposed Method Model.



Fig. 6.    Endpoints and Touching Point.



Fig. 7.    Neighboring Points.

In the third stage (Formulate direction), correct boundary of touching character is identified by neighbor points. The purpose of this step is to trace the boundary of touching character in correct direction because at the location of touching point there are many paths or curves. The coordinates of neighboring points $C_{p+1}$, $C_{p+2}$, $C_{p+3}$ and $C_{p+4}$ help in selection of right direction. Fig. 8 shows the selection of curve near touching point.



Fig. 8.    Formulate the Direction.

The proposed method is based on contour tracing starting from one endpoint and continue tracing the boundary of character to the touching point($C_p$). At this point, there are three possible directions (see Fig. 8) towards endpoint $E_2$, $E_3$ and $E_4$. Here the neighbor points play an important role for direction. Therefore tracing character boundary from $E_1$ will follow the curve towards endpoint $E_2$ because by comparing the coordinates of neighboring point $C_{p+1}$ with other neighboring points $C_{p+2}$, $C_{p+3}$ and $C_{p+4}$, there is a abrupt change towards point $C_{p+2}$ and $C_{p+4}$. While normal change towards neighboring point $C_{p+3}$. Fig. 8 shows three curves near touching point.

At this stage, endpoints, touching point, neighboring points and also curve direction of the touching image are identified. Next stage is segmentation of touching character image into separate characters. Fig. 9 explains the segmentation algorithm flowchart.



Fig. 9.    Flowchart.

The Explanation of the flowchart in Fig. 9 can be detailed as follow:

```
Algorithm : Segmentation of touching
   characters

Start
Step 1:    Input an image
   Image  is  handwritten  or  printed
   touching   or   without   touching
   characters.

Step 2: Perform Pre-processing
   Very  important  step  to  enhance the
   quality   of   image   and   remove
   unwanted information.

Step 3: Find End-Points and Touching
   point
   Using   overlapping   set  theory  to
   find the touching.

Step 4: Find Neighboring points near
   Touching point
   Here  find  neighbor  points  in  all
   direction,  near  the  junction  or
   common point

Step 5: Formulate the direction
   The  purpose  of  this  step  is  to
   trace  the  boundary  of  touching
   character  in  correct  direction
   because   at   the   location   of
   touching  point  there  are  many
   paths or curves.

Step 6: Segmentation of the touching
   characters
   At  this  step  touching  character
   image  is  segmented  to  separate
   characters.

End
```

## V. RESULTS AND DISCUSSION

Basically, this proposed method is for segmentation of touching handwritten Arabic characters. As due to lack of standard dataset for handwritten data. Touching handwritten Arabic data are collected from different datasets, while some are converting to touching by doing some manually work. In section A list of datasets for testing the proposed method. Section B is for comparison with others two selected methods, which are closely related to segmentation of touching characters.

### A. List of Dataset for this Research

Data collection is very important to test the performance of the proposed method. Due to unavailability of standard touching character datasets. Selected number of dataset, details of the collected data are shown in Table III below.

TABLE III. DATA COLLECTION

| Database | Purpose |
|---|---|
| AHDB | Off-line handwriting |
| IFN/ENIT | Tunisian city names Off-line handwriting |
| Arabic handwritten 1.0 | Off-line handwriting |
| IBN SINA | Arabic Manuscript |
| IAM | English Handwritten. Total number of writers, 657 collected handwriting samples. Number of isolated and labelled word 115320. |
| NIST | There are 150,000 handwritten binary digits number of broken digits 2600 |

Touching character were selected manually from the datasets in Table III, especially from AHDB as shown in Fig. 4. While in some cases did some manual work in order to get some challenging types of two touching characters to test the performance of the proposed method.

### B. Analysis

In this section presented results of the experiments, which were conducted to prove the performance of the proposed method. Comparison of proposed with that of other methods are shown in Table IV.

It shows from the experiments that proposed method is very flexible and efficient for segmentation of touching Arabic handwritten characters. Samples results are given in Fig. 10 below.

TABLE IV. ANALYSIS

| Touching Characters | Proposed Method | | |
|---|---|---|---|
| | Total Selected Images | Segmented Images | Segmentation Percentage |
| Other Method 1 [22] inter-word touching. | 220 | -- | 94% |
| Other Method 2 [19] inter-word touching. | 622 | 620 | 96.88% |
| Proposed method | 220 | 214 | 97.27% |



Fig. 10. Samples of Results.

Based on Fig. 10, the proposed method is correctly segmented touching handwritten Arabic character in an efficient way. Total number of images selected 220 of almost every type of touching character and out of these 220 images correctly segmented 214 images. Only 6 images (2.73%) either over or under segmented because of varieties in Arabic handwritten data.

## VI. CONCLUSION

Touching handwritten character images normally exist in old historical documents or writing text in Arabic calligraphy. These touching characters extensively happen in English, numbering and Arabic handwritten historical materials. However, for this research will considered only Arabic handwritten characters and numeric data. Thus proposed method is only for segmentation of touching Arabic handwritten characters to solve the longstanding and unsolved problem of segmentation of touching characters. In final conclusion, to say that proposed method is very flexible and if improved further can be used for multiple touching Arabic characters while in future work can be used for other languages that are similar to Arabic such as Urdu, Pashto and Farsi languages.

## ACKNOWLEDGMENT

### REFERENCES

[1] S. A. Malik, M. Maqsood, and F. Aadil, Advances in Information and Communication, vol. 70. Springer International Publishing, 2020.

[2] H. Modi and M. C., "A Review on Optical Character Recognition Techniques," Int. J. Comput. Appl., vol. 160, no. 6, pp. 20–24, 2017.

[3] C. S. Lwin and W. Xiangqian, "Image Purification Technique for Myanmar OCR Applying Skew Angle Detection and Free Skew," Int. J. Sci. Res. Sci. Technol., vol. 6, no. 1, pp. 186–203, 2019.

[4] N. Aouadi, S. Amiri, and A. K. Echi, "Segmentation of Connected Components in Arabic Handwritten Documents," Procedia Technol., vol. 10, pp. 738–746, 2014.

[5] T. Saba and A. Rehman, "Character Segmentation in Overlapped Script using Benchmark Database," pp. 140–143.

[6] A. Gattal, Y. Chibani, and B. Hadjadji, "Segmentation and recognition system for unknown-length handwritten digit strings," Pattern Anal. Appl., vol. 20, no. 2, pp. 307–323, 2017.

[7] J. H. Alkhateeb, "A Database for Arabic Handwritten Character Recognition," Procedia Comput. Sci., vol. 65, no. Iccmit, pp. 556–561, 2015.

[8] A. Amin, "Off-line Arabic character recognition," Pattern Recognit., vol. 31, no. 5, pp. 517–530, 2002.

[9] S. Ahmed, S. Naz, M. Razzak, and R. Yusof, "Arabic Cursive Text Recognition from Natural Scene Images," Appl. Sci., vol. 9, no. 2, p. 236, 2019.

[10] M. S. Deshmukh and S. R. Kolhe, "A Hybrid Character Segmentation Approach for Cursive Unconstrained Handwritten Historical Modi Script Documents," pp. 967–978, 2019.

[11] N. Vincent and J. M. Ogier, "Shall deep learning be the mandatory future of document analysis problems?," Pattern Recognit., vol. 86, pp. 281–289, 2019.

[12] M. Ayesh, K. Mohammad, A. Qaroush, S. Agaian, and M. Washha, "A Robust Line Segmentation Algorithm for Arabic Printed Text with Diacritics," Electron. Imaging, vol. 2017, no. 13, pp. 42–47, 2017.

[13] S. Eskenazi, P. Gomez-Krämer, and J. M. Ogier, "A comprehensive survey of mostly textual document segmentation algorithms since 2008," Pattern Recognit., vol. 64, no. October 2016, pp. 1–14, 2017.

[14] G. A. Farulla, N. Murru, and R. Rossini, "A fuzzy approach to segment touching characters," Expert Syst. Appl., vol. 88, pp. 1–13, 2017.

[15] S. Khan, H. Ali, Z. Ullah, N. Minallah, S. Maqsood, and A. Hafeez, "KNN and ANN-based Recognition of Handwritten Pashto Letters using Zoning Features," Int. J. Adv. Comput. Sci. Appl., vol. 9, no. 10, 2018.

[16] S. Wshah, Z. Shi, and V. Govindaraju, "Segmentation of Arabic handwriting based on both contour and skeleton segmentation," Proc. Int. Conf. Doc. Anal. Recognition, ICDAR, no. January, pp. 793–797, 2009.

[17] Y. M. Alginahi, "A survey on Arabic character segmentation," Int. J. Doc. Anal. Recognit., vol. 16, no. 2, pp. 105–126, 2013.

[18] K. Anwar, Adiwijaya, and H. Nugroho, "A segmentation scheme of Arabic words with harakat," 4th IEEE Conf. Commun. Networks Satell. COMNESTAT 2015 - Proc., pp. 111–114, 2016.

[19] N. Ouwayed and A. Belaïd, "Separation of overlapping and touching lines within handwritten arabic documents," Lect. Notes Comput. Sci. (including Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinformatics), vol. 5702 LNCS, pp. 237–244, 2009.

[20] L. Kang and D. Doermann, "Template based segmentation of touching components in handwritten text lines," Proc. Int. Conf. Doc. Anal. Recognition, ICDAR, pp. 569–573, 2011.

[21] A. N. G. Lopes Filho and C. A. B. Mello, "Segmentation of Overlapping Digits through the Emulation of a Hypothetical Ball and Physical Forces," pp. 223–226, 2015.

[22] N. Aouadi and A. Kacem, "A proposal for touching component segmentation in Arabic manuscripts," Pattern Anal. Appl., pp. 1–23, 2016.

[23] Z. Saber, A. Q. Sabri, A. Kamsin, and S. Hakak, "Efficient Approach to Segment Ligatures and Open Characters in offline Arabic Text," Int. J. Comput. Commun. Instrum. Eng., vol. 4, no. 1, pp. 40–44, 2017.

# Pedestrian Safety with Eye-Contact between Autonomous Car and Pedestrian

Kohei Arai[1], Akihito Yamashita[2], Hiroshi Okumura[3]

Dept of Information Science, Saga University, Saga City, Japan

*Abstract*—Method for eye-contact between autonomous car and pedestrian is proposed for pedestrian safety. The method allows to detect the pedestrian who would like to across a street through eye-contact between an autonomous driving car and the pedestrian. Through experiment, it is found that the proposed method does work well for finding such pedestrians and for noticing a sign to the pedestrians from the autonomous driving car with a comprehensive representation of face image displayed onto front glass of the car.

*Keywords*—*Autonomous driving car; eye-contact; self driving car; pedestrian safety; Yolo; OpenCV; GazeRecorder*

## I. Introduction

In recent years, the world has become more interested in autonomous driving of cars. However, the realization of fully automatic driving has various problems in terms of safety, and can be mentioned as difficult problems. The International Transport Forum at the OECD is an intergovernmental organization with 54 member countries. It acts as a think tank for transport policy and organizes the Annual Summit of transport ministers. ITF is the only global body that covers all transport modes [1]. Centralized raw data fusion, neural networks for machine learning is developed for autonomous driving car [2]. 360-degree real-time perception of vehicle surroundings is also developed. The developed software system runs on a flexible yet automotive-grade hardware platform with state-of-the-art FPGA, SoC and safety processors.

An autonomous vehicle driving control system carries a large number of benefits, whether for the engineering industry or engineering education. The system discussed here provides the measurements obtained from vision namely, offset from the centerline at some look-ahead distance and the angle between the road tangent and the orientation of the vehicle at some look-ahead distance and these are directly used for control [3].

The hardware capabilities are already approaching the levels needed for well-optimized Autonomous Vehicle software to run smoothly. Current technology should achieve the required levels of computational power—both for graphics processing units (GPUs) and central processing units (CPUs)—very soon [4]. Cameras for sensors have the required range, resolution, and field of vision but face significant limitations in bad weather conditions. Radar is technologically ready and represents the best option for detection in rough weather and road conditions. Lidar systems, offering the best field of vision, can cover 360 degrees with high levels of granularity. Although these devices are currently pricey and

too large, a number of commercially viable, small, and inexpensive ones should hit the market in the next year or two.

At present, connected vehicles equipped with a system that can communicate by IoT have appeared in order to ensure safety between vehicles. However, it does not hold between cars and pedestrians. In this study, we focus on eye contact with cars and pedestrians, which are not considered in the current autonomous driving, and aim to avoid danger and secure certain safety.

In this paper, from the background of the research in Section 1, we have presented the problems raised and the purpose of this research. In Section 2, we described trends in autonomous driving of automobile-related companies. The following section describes the system proposed in this paper. Then the following section gives an overview of this system followed by the conditions of eye contact. After that, image recognition, specifically pedestrian detection, and face and eye detection, as well as eye gaze estimation are described.

The experiments conducted in this study are described in the third chapter. Object detection using YOLO and the results are described. Then, the OpenCV face / eye detection procedure and results are followed together with the procedure and results of gaze estimation using GazeRecorder. After that, the discussion of the experiment is described. Final chapter gives the conclusions with some discussions.

## II. Proposed Method

### A. Proposed Method Overview

This section outlines the proposed system for securing pedestrian safety by eye contact detection in autonomous driving. This system projects the face of a person on the windshield of a car and sets eye contact with pedestrians. Therefore, we use an onboard camera to monitor pedestrians, detect pedestrians from the camera images, and confirm eye contact to avoid danger. Based on these, simulation is performed to verify whether safety can be maintained.

### B. Eye-Contact

In this study, we define the conditions under which eye contact between a car and a pedestrian can be made as follows, and the situation is shown in Fig. 1 [1].

*1)* A pedestrian can be recognized from the on-board camera image (Fig. 1(a)).
*2)* It turns out that a pedestrian is going to cross (Fig. 1(b)).
*3)* Noticed the car and confirmed the driver (Fig. 1(c)).

(a) Condition#1


(b) Condition#2


(c) Condition#3

Fig. 1.  Conditions for Eye Contact.

The pedestrian can cross the street safely by satisfying the above conditions and giving an eye contact to the pedestrian from the car and giving an intention to stop (Fig. 2).



Fig. 2.  Notifying a Sign to the Pedestrians from the Autonomous Driving Car with a Comprehensive Representation of Face Image Displayed onto front Glass of the Car.

## C.  Pedestrian Finding

In order to perform pedestrian recognition, we use an image recognition algorithm YOLO (You Only Look Once) [2] [3] that can be processed in real time by Joseph Redmon et al. In this research, we use Darknet [4] to implement YOLO.

## D.  Face and Eye Detection

We use the Haar-like feature classifier implemented in OpenCV (Open Source Computer Vision Library) [5] for face and eye detection. The features of various objects are learned in advance, and the xml file list [6] of the created feature data is shown in Table I.

In this research, we use "haarcascade_frontalface_alt.xml" and "haarcascade_eye.xml" for face detection and eye detection.

## E.  Gaze Detection

In this research, safety of pedestrians is secured by eye contact detection. Therefore, in order to perform eye contact, it is considered essential to estimate the line of sight after detecting the face and eye from detection of the pedestrian. This time, we use a system called GazeRecorder [7] for eye gaze estimation. The main features of this system are shown below:

*1)* Gaze tracking using a webcam
*2)* There is no physical contact for gaze estimation
*3)* Real-time gaze tracking is possible
*4)* Dynamic heat map can be generated using an adaptive time window

Since real-time processing is important in this system, GazeRecorder is adopted for gaze estimation.

TABLE I.      XML File List of Haar-Like Feature Classifier Implemented in Open CV

| File Name |
| --- |
| haarcascade_upperbody.xml |
| haarcascade_smile.xml |
| haarcascade_russian_plate_number.xml |
| haarcascade_righteye_2splits.xml |
| haarcascade_profileface.xml |
| haarcascade_lowerbody.xml |
| haarcascade_licence_plate_rus_16stages.xml |
| haarcascade_lefteye_2splits.xml |
| haarcascade_fullbody.xml |
| haarcascade_frontalface_default.xml |
| haarcascade_frontalface_alt2.xml |
| haarcascade_frontalface_alt_tree.xml |
| haarcascade_frontalface_alt.xml |
| haarcascade_frontalcatface_extended.xml |
| haarcascade_frontalcatface.xml |
| haarcascade_eye_tree_eyeglasses.xml |
| haarcascade_eye.xml |

## III. EXPERIMENT

### A. *Object Detection using YOLO*

We implemented YOLO under the following environment.

- OS: Windows 10 (64 bit)
- CPU: Intel Core i5-5275U
- Memory: 8.00 GB
- Visual Studio 2015
- OpenCV 3.2
- Darknet

For the construction of the development environment, we referred to "Darknet[1] installed on windows 10 to perform object recognition and object discrimination" [8]. Although it affects the processing speed, this research does not use the GPU[2], and CUDA[3] as well as cuDNN[4] are not used because the processing is performed by the CPU.

### B. *Experimental Results for Pedestrian Detection*

Fig. 3 shows the execution screen of YOLO in one scene of still image and moving image. Fig. 3(a) shows a screen shot of the process result of YOLO software operation while Fig. 3(b) shows an example of screen shot of the moving picture of object detection by YOLO. When processing with a still image, various objects other than human beings were detected. And even those who turned out to look backwards or from the image could be detected as humans. Also, it was found that the processing in the case of using the image in Fig. 3(c) takes approximately 2 seconds. When processing was performed using moving pictures, people and cars could be detected even though the image quality was poor. However, because processing was performed only by the CPU, it took more than 40 minutes to complete all processing even for a 30-second movie.

### C. *Face and Eye Detection by Python OpenCV*

In this experiment, we use Anaconda [9], a Python distribution provided by Continuum Analytics, in order to use Python and, "Face Detection using Haar Cascades" [10] was referred to for face and eye detection. Fig. 4 shows the results of face and eye detection. The face and detected parts are enclosed in a yellow frame, and the eye and detected parts are enclosed in a light blue frame.

### D. *Gaze Estimation by GazeRecorder*

When GazeRecorder starts up, the screen shown in Fig. 5 is shown. Here, it is also possible to set execution results and eye tracking. GazeRecorder is downloadable software [11]. After the object detection [12], gaze is detected with GazeRecorder. After all, learning process on gaze detection can be done with anaconda [13].

When you click the icon that says "Start Cam", the webcam starts. After that, the screen "Look at dot" is

displayed, so watch the white point at the center as shown in Fig. 6. Perform calibration by focusing on this point. And start tracking eye gaze.


(a) Still picture#1(Process result)


(b) Still picture#2


(c) Moving picture

Fig. 3. Pedestrian Detection by YOLO.

---

[1] https://pjreddie.com/darknet/

[2] https://ja.wikipedia.org/wiki/Graphics_Processing_Unit

[3] https://ja.wikipedia.org/wiki/CUDA

[4] https://developer.nvidia.com/cudnn

(a) Lena        (b) Co-author

Fig. 4. Face and Eye Detection.



(a) Star-up Screen Display.



(b) Execution Result Setting Screen Display.



(c) Result Setting Screen Display.

Fig. 5. Parameter Setting Screen Display of the Gaze Recorder.



Fig. 6. Screen Display for Calibration.

Clicking on the "Rec" icon will start recording, and clicking on the "Stop Rec" icon will stop recording. The videos taken here are stored in an arbitrary directory, and a heat map is automatically generated.

As shown in Fig. 7, GazeRecorder shows that the face is captured in three dimensions. The light mesh area represents the entire eyeball, and the red circle indicates the pupil and the iris. Also, the light blue lines extending from both eyes indicate the pupil and the iris.

In the heat map, it shows whether a human is looking at the screen part, and also the place where the user is gazing is divided according to the degree of color as shown in Fig. 8, and the red part is the place that is most often seen. The purple point from the center of the screen to just above is what is actually tracking the eye gaze.



Fig. 7. Example of the Detected Line of Sight Vector and Gaze using Gaze Recorder.



Fig. 8. Example of the Heat Map Derived from Gaze Recorder.

When performing eye tracking, the PC specifications described in the previous section enable real time eye tracking with almost no time lag.

Thus, the pedestrian gaze can be recognized. When the gaze to point to the autonomous car, then the car recognizes pedestrian's intention to across the street. Therefore, the car stops for the pedestrian with the sign shown in Fig. 2 for notifying the pedestrians with a comprehensive representation of face image displayed onto front glass of the car

## IV. Conclusion

Method for eye-contact between autonomous car and pedestrian is proposed for pedestrian safety. The method allows to detect the pedestrian who would like to across a street through eye-contact between an autonomous driving car and the pedestrian. Through experiment, it is found that the proposed method does work well for finding such pedestrians and for noticing a sign to the pedestrians from the autonomous driving car with a comprehensive representation of face image displayed onto front glass of the car.

It is found that when performing pedestrian detection with YOLO and eye tracking with GazeRecorder together with face and eye detection with OpenCV, the PC specifications described in the previous section enable real time eye tracking with almost no time lag.

Thus, the pedestrian gaze can be recognized. When the gaze to point to the autonomous car, then the car recognizes pedestrian's intention to across the street. Therefore, the car stops for the pedestrian with the sign shown in Fig. 2 for notifying the pedestrians with a comprehensive representation of face image displayed onto front glass of the car.

Future research works are as follows:

In this experiment, just visible camera is used for pedestrian detection and face/eye detection as well as gaze estimation. Further experiments have to be done with Near Infrared: NIR cameras and NIR Light Emission Diode: LED for same objectives in a bad weather condition and in night time operations.

## Acknowledgments

## References

[1] International Transport Forum 2 rue André Pascal 75775 Paris Cedex 16 Franc https://cyberlaw.stanford.edu/files/publication/files/15CPB Auto nomousDriving.pdf.

[2] Marcel Zolg, Techniocal University of Munich, https://www. mentor.com/mentor-automotive/autonomous

[3] KHALID BIN ISA Department of Computer Engineering, Faculty of Electric and Electronic Engineering, University College of Technology Tun Hussein Onn, 86400 Parit Raja https://www.ijee.ie/articles/Vol21-5/Ijee1674.pdf.

[4] By Kersten Heineke, Philipp Kampshoff, Armen Mkrtchyan, and Emily Shao https://www.mckinsey.com/industries/automotive-and-assembly/our-insights/self-driving-car-technology-when-will-the-robots-hit-the-road

[5] https://www.youtube.com/watch?v=QdYMOZRGdMw

[6] Joseph Redmon, Santosh Divvala, Ross Girshick, Ali Farhadi, "You Only Look Once: Unified, Real-Time Object Detection"

[7] YOLO | Joseph Redmon, https://pjreddie.com/darknet/yolo/

[8] Darknet | Joseph Redmon, https://pjreddie.com/darknet/

[9] OpenCV (Open Source Computer Vision Library), https://opencv.org/

[10] OpenCV Face detection（Haar-like feature extraction）|Qiita, https://qiita.com/hitomatagi/items/04b1b26c1bc2e8081427

[11] GazeRecorder | GazeRecorder, http://gazerecorder.christiaanboersma.com/gazerecorder/

[12] Object detection, ryuji shirata , "Darknet on windows10, http://weekendproject9.hatenablog.com/about.

[13] Anaconda and JetBrains Join Forces to Launch 'PyCharm for Anaconda' https://www.anaconda.com/.

## Authors 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 and 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 an Adjunct Professor of University of Arizona, USA since 1998. He also is Vice Chairman of the Commission-A of ICSU/COSPAR since 2008. He received Science and Engineering Award of the year 2014 from the minister of the ministry of Science Education of Japan and also received the Bset Paper Award of the year 2012 of IJACSA from Science and Information Organization: SAI. In 2016, he also received Vikram Sarabhai Medal of ICSU/COSPAR and also received 20 awards. He wrote 34 books and published 520 journal papers. He is Editor-in-Chief of International Journal of Advanced Computer Science and Applications as well as International Journal of Intelligent Systsems and Applications. http://teagis.ip.is.saga-u.ac.jp/

# Tennis Player Training Support System based on Sport Vision

Kohei Arai[1], Toshiki Nishimura[2], Hiroshi Okumura[3]

Dept. of Information Science
Saga University
Saga City, Japan

*Abstract*—**Sports vision based tennis player training support system is proposed. In sports, gaze, dynamic visual acuity, eye movement and viewing place are important. In sports vision, Static eyesight, Dynamic visual acuity, Contrast sensitivity, Eye movement, Deep vision, Instant vision, Cooperative action of eye, hand and foot, and Peripheral field are have to be treated. In particular for the tennis, all of the items are very important. Furthermore, trajectory of gaze location and tennis racket stroke gives some instructions for skill-up of tennis play. Therefore, sports vision based tennis player training system is proposed. Through experiment, it is found that the proposed system does work well for improvement of tennis players' skills.**

*Keywords—Sport vision; static eyesight; dynamic visual acuity; contrast sensitivity; eye movement; deep vision; instant vision; cooperative action of eye; hand and foot; peripheral field*

## I. INTRODUCTION

Sports vision based tennis player training system is proposed. There is Sport Vision Association in Japan [1]. In sports, gaze, dynamic visual acuity, eye movement and viewing place are important. In sports vision, Static eyesight, Dynamic visual acuity, Contrast sensitivity, Eye movement, Deep vision, Instant vision, Cooperative action of eye, hand and foot, and Peripheral field are have to be treated. In particular for the tennis, all of the items are very important.

Face detection from the acquired imagery data is common and popular technology. Computer input system based on viewing vector estimation with iris center detection from face image acquired with web camera allowing users' movement is proposed [1]. Also, method for face identification with Facial Action Coding System: FACS based on eigen value decomposition is proposed [2]. The first book of its kind devoted to the emerging field of computer vision in sports is published [3]. Quite recently, the 5th International Workshop on Computer Vision in Sports (CVsports) at CVPR 2019 is held in Long Beach California, U.S.A [4]

A new Olympic vault dataset is proposed [5] and present three frameworks for action quality assessment which improve upon published results: C3D-SVR, C3D-LSTM and C3D-LSTM-SVR.The frameworks mainly differ in the way they aggregate clip-level C3D features to get a video-level description. This video-level description is expressive about the quality of the action. The task of detecting swimming strokes in the wild is demonstrated [6]. However, without

modifying the model architecture or training method, the process is also shown to work equally well on detecting tennis strokes, implying that this is a general process. The outputs of the system are surprisingly smooth signals that predict an arbitrary event at least as accurately as humans (manually evaluated from a sample of negative results).

A convolutional neural network (CNN) has been designed to interpret player actions in ice hockey video [7]. The hourglass network is employed as the base to generate player pose estimation and layers are added to this network to produce action recognition. As such, the unified architecture is referred to as action recognition hourglass network, or ARHN. ARHN has three components. Group activity recognition in sports is often challenging due to the complex dynamics and interaction among the players. In this paper, we propose a recurrent neural network to classify puck possession events in ice hockey. Our method extracts features from the whole frame and appearances of the players using a pre-trained convolutional neural network. In this way, our model captures the context information, individual attributes and interaction among the players [8].

In order for action recognition to be useful in sports analytics a finer-grained action classification is needed. For this reason we focus on the fine-grained action recognition in tennis and explore the capabilities of deep neural networks for this task. In our model, videos are represented as sequences of features, extracted using the well-known Inception neural network, trained on an independent dataset. Then a 3-layered LSTM network is trained for the classification [9]. A learning-based framework that takes steps towards assessing how well people perform actions in videos is proposed [10]. The approach works by training a regression model from spatiotemporal pose features to scores obtained from expert judges. Moreover, the approach can provide interpretable feedback on how people can improve their action.

Using the line-of-sight camera, acquire eye-gaze video and analyze the difference between the beginner and the experienced person or the gaze at good and bad times. Also, I would like to support technical improvement by proposed system. Using head set of work camera, players' gaze is estimated together with ball trajectory for prediction. Then, some instructions can be provided to the player in concern. This is what I intend to do. Moreover, it is possible to make a simulated experience by Virtual Reality: VR from the line-of-sight image.

---

[1] http://www.sports-vision.jp/deta.htm

In the next section, the proposed system for improving tennis players' skills is described. Then, preliminary experiment is followed. After that, conclusion is described together with some discussions. Finally, future works are followed.

## II. PROPOSED SYSTEM

### A. Design Concept

The proposed system provides expected position of tennis ball immediately after the ball is hit by the offensive player, as well as instructions and week points to tennis players. Fig. 1 shows design concept of the proposed system. There are two cameras, one is observe the tennis court and the other one is attached to the tennis player's head in the system.

From the video data acquired #1 camera, tennis ball is looked at and always calculate the trajectory and predict the expected tennis ball position touch down on the court. Therefore, instruction can be provided to the tennis player. Thus, the tennis player can take a next action so quickly. On the other hand, #2 camera is always looking forward forehead direction, head pose direction. Therefore, it is understand which direction is tennis player is looking at. Thus, week points can be provided when the tennis player is not looking at the appropriate direction after the play. That will help to improve their skills.



Fig 1.    Design Concept

Tennis player wears GoPro of camera and play. Tennis ball can be extracted from the acquired moving picture using OpenCV. The trajectory of the extracted tennis ball and gaze locations are plotted and displayed onto computer screen. Then, valuable instructions are given to the players.

### B. Software and Hardware

Tennis ball can be extracted from the video and analyze where the players are looking together with gaze direction. There are some actions in tennis plays, Serve, Stroke, Volley, Smash, Match (singles). These actions are to be identified from the acquired videos. The identified action types can be provided to the tennis player through wireless communication network.

Major hardware and software used are as follows:

- GoPro HERO6   (Camera)
- OpenCV 3.1.0 for image analysis
- Python 2.7 or Visual Studio2017 C++
- Major specification of GoPro HERO6 is as follows,
- Weight: 117 g

Action camera

Image stabilization

Hi-Vision

Waterproof function

Field of view: "Wide angle", "fisheye"

Frame rate 4K / 60fps, 1080p / 240fps

Fig. 2 shows outlook of the GoPro HERO6.



(a)Back view                    (b)Front view

Fig 2.    Outlook of the GoPro HERO6.

### C. Major Characteristics of Tennis Player

From data by SONY smart tennis sensor [2] (SONYSmartTennisSensor), major characteristics of tennis player are as follows:

Male average speed in 20's

Forehand stroke: 90 km / h (25 m / s)

Serve: 115 km / h (31 m / s)

The vertical length of the tennis court (the distance between the players): 23.77 m (about 25 m)

Until the hit ball reaches the opponent

Forehand stroke: about 1.0 second

Serve: about 0.8 seconds

### D. OpenCV

In order to analyze the acquired imagery data, OpenCV is used. The main functions that can be done are as follows:

- Filtering
- Matrix operation
- Object tracking
- Area segmentation
- Camera calibration
- Feature point extraction
- Object Recognition
- GUI
- Machine learning
- Panorama synthesis

---

[2] http://tennisblog.smartsports.sony.net/entry/worlddata-japan-detail-stroke

- Computational photography

Color based method is used for tennis ball detection. Detection of tennis ball is performed not only by using circle detection but also by color extraction as follows:

*1)* Hue, Saturation, Value: HSV conversion of the original image
*2)* Extract only yellow color
*3)* Extraction of contour (extraction of maximum contour only)
*4)* Draw the minimum circumscribed circle

Fig. 3 shows examples of detected circles using color information based method. The proposed tennis ball detecting method is based on the combination between shape based and color based methods. Therefore, tennis ball can be detected even if the tennis ball is not circle shape and even if the tennis ball is not yellow color.



Fig 3.    Example of yellow colored object detection.

*E. Tennis Ball and Gaze Location Trajectories as well as the Acquired Image of the One Forehand Stroke*

Tennis ball and gaze location trajectories can be analyzed in the 1080 by 1920 pixels of computer screen from the acquired moving pictures. One of the examples is shown in Fig. 4. Also, tennis ball location trajectory is analyzed and displayed onto computer screen as shown in Fig. 5. Furthermore, Fig. 6(a) and (b) show the acquired image of the one forehand stoke image of an expert and that of a beginner.



Fig 4.    Example of gaze location trajectory for one forehand stroke.



Fig 5.    Example of tennis ball location trajectory.



(a)Expert



(b)Beginner

Fig 6.    Acquired images with GoPro camera attached to an expert and a beginner of tennis players.

Thus, valuable instructions can be made available for tennis players based on the trajectories of ball and gaze locations.

## III.   EXPERIMENT

Six experts and six beginners of tennis players are participated to the experiment. Trajectories of the gaze location are analyzed for five forehand strokes. Fig. 7(a) to (f) shows the trajectories of experts while Fig. 7(g) to (l) shows those of beginners.

Where frame rate is 59 f/s which corresponds to about 0.0179 sec. per frame while time duration between two dots in the gaze trajectory figure is 0.0179 sec. From the trajectories, hitting point and stroke position and its blurring are recognized. It is quite obvious that head pose is stable, hitting

point is also stable, pulling jaw, tennis ball is caught with eyes for experts while opposite for beginners. In other word from the sport vision of point of view, experts are good at the followings:

Peripheral vision
kva motion vision
Deep vision
Fixed gaze



(a)Expert#1



(b)Expert#2



(c)Expert#3



(d)Expert#4



(e)Expert#5



(f)Expert#6



(g)Beginnert#1

(h)Beginner#2


(i)Beginner#3


(j)Beginner#4


(k)Beginner#5


(l)Beginner#6

Fig 7.   Gaze trajectories for 5 forehand strokes of experts and beginners.

Also, expert views opponent, catches the entire tennis court, catches the tennis ball properly as shown in Fig. 8(a) while beginner cannot view opponent, does not catch tennis ball properly, does not catch tennis court as shown in Fig. 8(b). This is the results from the acquired image of experts and beginners at the hitting point.


(a)Expert


(b)Beginner

Fig 8.   Acquired image of experts and beginners at the hitting point.

## IV. Conclusion

Sports vision based tennis player training support system is proposed. In sports, gaze, dynamic visual acuity, eye movement and viewing place are important. In sports vision, Static eyesight, Dynamic visual acuity, Contrast sensitivity, Eye movement, Deep vision, Instant vision, Cooperative action of eye, hand and foot, and Peripheral field are have to be treated. In particular for the tennis, all of the items are very important. Furthermore, trajectory of gaze location and tennis racket stroke gives some instructions for skill-up of tennis play. Therefore, sports vision based tennis player training system is proposed. Through experiment, it is found that the proposed system does work well for improvement of tennis players' skills.

It is found that expert views opponent, catches the entire tennis court, catches the tennis ball properly while beginner cannot view opponent, does not catch tennis ball properly, does not catch tennis court

For future research works, real time instruction has to make available for tennis players.

*1)* Points looking at for expertized tennis player. This fact will be clarified.

*2)* Then, action classification from the acquired moving pictures will be conducted. If it is possible to identified just after the initial action, then next action can be recognized. After that, some preparatory action can be done with the recognized actions.

## Acknowledgment

## References

[1] Kohei Arai and Hiromi Uwataki, Computer input system based on viewing vector estimation with iris center detection from face image acquired with web camera allowing users' movement, Electronics and Communication in Japan, 92, 5, 31-40, John Wiley and Sons Inc.,2009.

[2] Kohei Arai, Method for face identification with Facial Action Coding System: FACS based on eigen value decomposition, International Journal of Advanced Research in Artificial Intelligence, 1, 9, 34-38, 2012.

[3] https://www.springer.com/us/book/9783319093956

[4] http://www.vap.aau.dk/cvsports/

[5] https://arxiv.org/pdf/1611.05125.pdf

[6] http://homepage.cs.latrobe.edu.au/zhe/files/StrokeDetection.pdf

[7] http://openaccess.thecvf.com/content_cvpr_2017_workshops/w2/papers/Fani_Hockey_Action_Recognition_CVPR_2017_paper.pdf

[8] http://openaccess.thecvf.com/content_cvpr_2017_workshops/w2/papers/Tora_Classification_of_Puck_CVPR_2017_paper.pdf

[9] http://www.doc.ic.ac.uk/~wjk/publications/vinyes-knottenbelt-cvsports-2017.pdf

[10] https://www.csee.umbc.edu/~hpirsiav/quality.html

## Authors 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 and 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 an Adjunct Professor of University of Arizona, USA since 1998. He also is Vice Chairman of the Commission-A of ICSU/COSPAR since 2008. He received Science and Engineering Award of the year 2014 from the minister of the ministry of Science Education of Japan and also received the Bset Paper Award of the year 2012 of IJACSA from Science and Information Organization: SAI. In 2016, he also received Vikram Sarabhai Medal of ICSU/COSPAR and also received 20 awards. He wrote 34 books and published 520 journal papers. He is Editor-in-Chief of International Journal of Advanced Computer Science and Applications as well as International Journal of Intelligent Systsems and Applications. http://teagis.ip.is.saga-u.ac.jp/

# A Low Power Consuming Model of Parallel Programming for HPC System

Mohammed Nawaf Altouri[1], Abdullah M. Algarni[2]

Department of Computer Science, King Abdulaziz University (KAU)

P.O. Box 80221, Jeddah 21589, Saudi Arabia

*Abstract*—**For most of the past five decades, the growing computational power of supercomputers has come primarily from a doubling of clock frequency every 18 months. Over this time period, the clock rate increased by six orders of magnitude, while the number of processors increased by three orders of magnitude. The major challenge caused by the increasing scale and complexity of HPC systems is the massive power consumption. Due to constraints on heat and the power requirements of today's microprocessors, vendors have shifted to putting multiple processors (cores) on a chip. The number of cores per chip is expected to continue increasing exponentially over the next decade. One expected strategy is the correct usage of parallel programming models that decrease power consumption and increase system performance through massive parallelism (concurrency). In the current study, we have proposed a Hybrid MVAPICH-2 + CUDA (HMC) parallel programming model that outperformed other state-of-the-art dual and tri hierarchy level approaches with respect to power consumption and execution time. Moreover, the HMC model was evaluated by implementing the matrix multiplication benchmarking application. Consequently, it can be considered a leading model for the emerging Exascale computing system.**

*Keywords—HPC; parallel computation; power consumption; hybrid programming; MVAPIC2; CUDA*

## I. INTRODUCTION

In the next decade, an extreme level computing system called Exascale computing is anticipated to revolutionize computational science and engineering by providing 1018 FLOPS operations per second, which will be comprise hundreds of thousands of heterogeneous compute nodes linked by complex networks [1]. A projection from the world's most powerful system with the capability of handling Petaflops per second developed in the recent past (2014,) creates the possibility of producing Exascale systems deployed in the 2020 timeframe [1][2]. For this Ultra-scale computing system, an extensive change in node architectures is expected, replacing the current trend of increasing clock speed by doubling the number of cores in a system [3][4]. However, a prominent level of computation for the Exascale system has some valid limitations, such as energy consumption (20MW), time of delivery (2020), number of cores (100 million) and cost ($200M) [3][5]. According to the US Department of Energy (DOE), energy consumption per flop must be less than 20 Pico-Joules (PJs) [6]. Under these hard limitations, the development community must rethink its use of existing technologies and expand the co-design space to find better solutions, new applications, algorithms, better technology, and performance.

In attaining the emerging supercomputing goal, one faces number of effective challenges (such as massive power consumption, extraordinary parallelism, memory, bandwidth, latency, hierarchy level, communication and synchronization mechanism, resilience and heterogeneity) that must be determined by introducing new hardware, general-purpose multi-cores and special-purpose accelerators, frameworks, programming models and algorithms. In the last two decades, the microprocessor-based single CPU has increased system performance and decreased cost; however, because of heat dissipation and energy consumption issues, this approach reached a limit [7], [8]. The solution was switched to a model in which the microprocessor has multiple processing units known as "cores" [9], [10]. Therefore, two approaches were introduced and are currently being used "multi-core" [11] integrates a few cores into a single microprocessor while with "many-core" [12] , a large number of cores are integrated into a single device, called a GPU (Graphical Processing Unit) or GPGPU (General Purpose Computing Graphical Processing Unit). Therefore, the aim is to minimize power consumption in the system by increasing system performance through a combination of course-grain and fine-grain parallelism using heterogeneous parallel programming models [13].

In the current study, our primary objective was to introduce a new parallel programming model that can reduce power consumption in the system to deal with several linear/non- linear HPC benchmarking applications. In terms of introducing such a model, we have conducted an empirical study in which we investigated the existing parallel programming models being considered for emerging supercomputing systems. In terms of consequences, we proposed a new parallel programming model: Hybrid MVAPICH-2 + CUDA (HMC), which section III describes in detail.

The rest of the content is organized follows. Section II demonstrates the related work of parallel computing software and hardware technologies and renowned approaches. Section III describes the proposed HMC model. Section IV describes the empirical study, consequences and comparative analysis.

## II. RELATED WORK

In [14], [15], [16] and [17], the authors investigated the primary challenges for emerging Exascale computing systems. According to the authors, supercomputer architectures have gone from 1000 processors to 100,000 processors in the last five years while next-generation systems will have more than one million processors. The rate of growth of parallelism is in fact accelerating, it will likely exceed one hundred million when

Exascale systems appear. Some estimates even predict that the need for multiple threads to cover main memory and communication latency means that scientific codes will contain billions of threads. However, they determined that the major challenge caused by the increasing scale and complexity of HPC system's the massive power consumption. One expected strategy is the correct usage of parallel programming models that decrease power consumption and increase system performance through massive parallelism (concurrency). The increase of concurrency from hundreds of thousands to hundreds of millions is also a tremendous challenge for system software to manage; in addition, it is a challenge with respect to application's ability to achieve good performance at this level of parallelism.

Further, Carretero. J, et al. [8] said that an Exascale architecture should be both energy-efficient and power proportional. The subject of reducing the energy consumption in computing brought up two main research directions. The first direction is concerned with power-aware and thermal- aware hardware design, including low-power techniques on all levels. The second research direction is based on the development of power-aware software for the entire software stack, including operating systems, compilers, applications and algorithms. The authors investigated several energy- saving mechanisms in hardware and software including DVFS (Dynamic voltage frequency scaling), clock Gating, power gating, and coprocessors or accelerators. Clock gating reduces power consumption by disabling the clock in those parts of the circuit that are idle or, as in the case of flip-flops, maintain a steady state that does not need to be refreshed. Similarly, power gating achieves a better power reduction than does clock gating. This had led emerging technologies- such as GPU (Graphics Processing Unit), MIC (Many Integrated Cores) and FPGA (Field Programmable Gate Array) becoming the most promising technologies to overcome the power consumption challenge.

Maglione et al., [18] investigated the Advanced Configuration and Power Interface (ACPI), which is an open standard for device power management co-developed by Hewlett-Packard, Intel, Microsoft, Phoenix, and Toshiba. Other advanced tools including Memscale and PGCapping were also investigated. Memscale includes a control algorithm that minimizes the overall system energy based on performance counter monitoring while PGCapping integrates power gating with DVFS for chip multiprocessors.

Feng and Xixhou [19] analyzed different modern architectures and common applications and illustrated that some system components such as caches and network links-disproportionately consume extensive power for common HPC applications. They demonstrated that a large percentage of power consumed in caches and networks can be saved using our approach automatically. Regarding energy optimization, energy spent on cooling accounts for about 40% of total energy consumption in a data center. The author's focus was to extend energy optimization work beyond machine energy savings so that cooling energy could be reduced. They worked on a runtime system which used Dynamic Voltage and Frequency Scaling (DVFS) to minimize the occurrence of hotspots by keeping core temperatures in check. The consequences showed that we can save considerable cooling energy using this temperature aware load balancing.

Aniruddha, et al. [20] discussed performance optimization under power consumption constraints. According to the authors, the power consumption of Intel's Sandy Bridge family of processors can be user-controlled through the Running Average Power Limit (RAPL) library. It has been shown that an increase in the power allowed for the processor (and/or memory) does not yield a proportional increase in the application's performance. Consequently, for a given power budget, it could be better to run an application on a larger number of nodes with each node capped at lower power rather than fewer nodes, with each running at its TDP.

Geist, Al- and Daniel A. Reed [21] conducted a survey about primary high-performance computing challenges. They explored the state of the art and reflected on some of the challenges likely to be faced during the building of trans- Petascale computing systems. The energy required for Petascale clusters is now measured in megawatts; commercial cloud data centers consume 25–60 megawatts. Putative Exascale systems would consume hundreds of megawatts using current technologies. For future technologies, two architectural paths are emerging to address the three key challenges including reliability, energy consumption and software complexity. Consequently, energy consumption is a major driver in the emergence of the above two architectural designs. For large-scale heterogeneous system, energy efficiency can be obtained through energy constrained scheduling and adaptive parallelism. Energy consumption and power costs must be managed with as much care as performance and resilience. By contrast, software complexity must be managed to decrease software development costs.

According to [22], hybrid techniques are solutions that allow for emerging HPC computing systems to deal with the primary issues of power consumption and enhancing performance. The authors proposed different hybrid techniques of MPI+X where X is any parallel programming model to compute GPU. Conventionally this X is considered CUDA or Open ACC.

Further, the authors in [23] proposed a tri-level hybrid of MPI+OMP+CUDA (MOC) to achieve massive parallelism. The primary purpose of this model was to achieve all levels of parallelism including coarse grain, fine grain and finer granularity in the system during the execution of any benchmarking application over a large cluster system. Therefore, the proposed MOC model was implemented on different HPC systems and evaluated the performance as well as power consumption. The MOC model reduced the system's power consumption by 20 MW overall by enhancing the performance by 20 PFLOPS in HPC systems.

## III. PROPOSED HYBRID PARALLEL PROGRAMMING MODEL

In this section, we present the proposed dual-level parallel programming model for the high-end computing system. The proposed approach is a hybrid of MVAPICH-2 and CUDA, named the Hybrid MVAPICH-2+CUDA (HMC) Model.

### A. Selected Model for Enhancement

Indeed, as more and more compute cores become available on a single node, the expectation is that communication of the local node will play an increasingly important role in the overall performance of parallel applications such as MPI applications. Therefore, it is crucial to optimize intra-node communication

paths utilized by MPI libraries. As much communication of data between processors as it will consume more power, so by reducing that communications overhead between processors we can increase that performance and limit that power consumption through intra-node communication. We have decided that MPI+CUDA is the best hybrid model. The main reason for that is that it has less communication overhead; other reasons are described in detail in Section IV. This is done by conducting an empirical study and analysis for all the parallel programming model that shown in Table I and Table II.

Many versions of MPI can deliver the best performance. One way is to use MVAPICH-2, an open source implementation of the Message Passing Interface that can deliver the best scalability, performance and fault tolerance for high-end computing systems and servers using InfiniBand, which is used for interconnect communication and was first popular in high performance computing environments, Internet Wide Area RDMA Protocol (IWARP) and RDMA Over Converged Ethernet (ROCE) networking technologies. It facilitates the task of porting MPI applications to run on clusters with NVIDIA GPUs by supporting standard MPI calls from GPU device memory. Furthermore, it optimizes the data movement between host and device and between GPUs in the best way possible while requiring no effort from the application developer.

Random Direct Memory Access (RDMA) is hardware architecture we used to implement our hybrid model. It is especially useful in massively parallel computer clusters. After the initialization of MPI, a global communicator will contain all processors on that library. Therefore, unprecedented scalability, resiliency, and overhead limitations will be on MPI application. However, MVAPICH-2 has more directives that can decrease power consumption and deliver the best performance. This library will have MPI Sessions- a fundamental change in how we address and organize MPI processes that remove the known scalability barriers by no longer requiring the inclusion of all possible communication peers on the global communicator.

TABLE I.    EXECUTION TIME (PERFORMANCE)

| Matrix Size | MPI + OMP | OMP + CUDA | MPI + Open ACC | MPI + CUDA | MPI+ OMP+ CUDA |
|---|---|---|---|---|---|
| 1000^2 | 4.5 | 3.22 | 3.79 | 3.01 | 5.01 |
| 2000^2 | 14.55 | 8.55 | 5.99 | 4.81 | 10.81 |
| 3000^2 | 45.99 | 37.83 | 38.21 | 34.2 | 32.2 |
| 4000^2 | 107.06 | 85.18 | 60.86 | 52.29 | 46.29 |
| 5000^2 | 202.39 | 145.62 | 91.39 | 88.74 | 70.74 |
| 6000^2 | 341.56 | 189.59 | 128.22 | 109.94 | 107 |

TABLE II.    POWER CONSUMPTION

| Matrix Size | MPI + OMP | OMP + CUDA | MPI + Open ACC | MPI + CUDA | MPI + OMP+ CUDA |
|---|---|---|---|---|---|
| 1000^2 | 311 | 187.7 | 198.92 | 195.4 | 235.67 |
| 2000^2 | 334.7 | 200 | 207.07 | 205.13 | 248.91 |
| 3000^2 | 329.6 | 207 | 226.58 | 213.05 | 256.34 |
| 4000^2 | 332 | 245.82 | 243.41 | 233.21 | 263.56 |
| 5000^2 | 323 | 263.26 | 262 | 255.84 | 268.29 |
| 6000^2 | 326 | 271.44 | 279.52 | 271.72 | 273.603 |

Session facilitates these efforts with two key contributions:

- A scalable representation of communication groups.

- A tighter integration of MPI applications with the underlying runtime system.

### B. HYBRID MVAPICH-2+CUDA (HMC) Model

As shown in Fig. 1 below, MPI initializes by creating a session with a specific run time and get the session name. Therefore, any named set of processes that are exposed by that session can be converted into a group 'MPI group' and define the ranks. These groups will communicate through a parent communicator" which is used by MPI to orchestrate (matching) the communication needed to create a new communicator". Thereafter, it will initialize the data to be computed and every processor will retrieve with a rank number. If the rank does not master the processor these data distributed over processors otherwise, its release and exit. So here will call for the CUDA kernel which means it will parallelize the data from the host (CPU) to device (GPU) cores for actual execution of DMM.

In CUDA, the GPU block dispatcher will schedule the grid by assigning each thread to one of computational core and these threads will be synchronized by self-cooperation. Each block has its own shared memory so that the thread will process the data using this shared memory within the block then return the result to the scheduler .Therefore, this processed data will be on GPU global memory that is visible to the host or CPU memory so it will copy the processed data from Device to the host memory. Finally, if the rank is a master processor, it will receive the processed data from all ranks then it will finalize and print the result.



Fig. 1.   HMC Architecture.

### IV. EXPERIMENTAL RESULTS

This section present the results of experiments conducted for the proposed study. All the experiments were performed on HPC Xeon Phi with GPU 1070-ti. To quantify the primary factors (including performance and power consumption) in the proposed HMC model, we firstly carried out an empirical study in which we computed the different datasets of matrix multiplication and Implemented multiple parallel programming models, including single [24], dual [22], [25] and tri-hybrid [23]. Then we executed the same datasets on HMC model. Leading to quantification factors, we measured the system's performance by quantifying

execution time against each matrix size. By contrast, the second factor of power consumption was measured using the TechpowerUp GPU-z 2.6.0 software. A running screen of TechpowerUp GPU-z is shown in Fig. 2.

According to existing state-of-the-art models and strategies, single level parallel programming models are not individually applicable for emerging HPC systems. Therefore, it must be in hybrid to achieve massive parallelism through a cluster system with multiple nodes. However, we excluded single level models from our comparative analysis, and followed dual- and Tri-hybrid parallel programming models. Leading to dual-level models, we compared the results with HMC in matrix multiplication with different size executions in the range of 1000x1000 to 6000x6000 matrix size, as shown in Fig. 3.

As per Fig. 3, we noticed that at an initial execution 1000 and 2000 matrix sizes, no big difference was seen in all executions, through a tremendous change was seen when matrix size was increased. We can see a major difference in execution time at matrix size 4000x4000 where heterogeneous Computation over cluster system outperformed the homogeneous/single computing system. From this position, the graph changed gradually onward, where the peak reading was noted for the hybrid of MPI+OMP with 342 number of seconds for 6000 matrix size execution. By contrast, the heterogeneous model on single node executed in less time, as OMP+CUDA took 190 seconds .We noticed that other heterogeneous models with the hybrid of MPI+ Open ACC and MPI+CUDA executed the same matrix size in 125 and 110 seconds, respectively, as MPI version 3 improved in terms of enhancing performance under power consumption limitations. Therefore, we integrated MPI-3 (MVAPICH-2) with CUDA which is the proposed model of current study and executed all these datasets in similar way. We discovered a drastic improvement in performance which was 80 second execution time against 6000x6000 matrix size. This is because on MPI-3 it

has more directives, like sessions and groups, and it will be scalable. These directives facilitate an improvement in the performance of our HMC model.

According to [23], the proposed MOC model was considered promising for emerging HPC systems to attain a massive performance. Also, the MOC model introduced three level of parallelism, including coarse grain, fine grain and finer grain. Therefore, we also executed the similar matrix multiplication dataset on MOC and compared it with the proposed HMC model as shown in Fig. 4.

It was observed that HMC outperformed throughout the execution as compared to MOC in all datasets. We observed a 20-second difference in the execution for a large dataset, which is the improved performance in HMC.

Further, to quantify the second objective of study, we observed the power consumption in all the selected models and compared it to proposed HMC model. The analysis mechanisms were similar to those used for execution time. We firstly evaluated dual-level homogeneous and heterogeneous models and compared them to HMC with respect to power consumption, as shown in Fig. 5.

MVAPICH-2 introduced the new directives, which are used to optimize the communication cost among the processors. Consequently, it causes a reduction in the system's power consumption during execution. Based on these improvements, we observed that HMC decreased the power consumption throughout the executions. For small computations, up to 50 watts' power consumption was observed that was reduced in HMC by comparing the consumed power measured in the best model MPI+CUDA from existing state-of-the-art mechanisms. As shown in Fig. 6, the power consumption was evaluated in the MOC model to conduct a comparative analysis with HMC model.



Fig. 2.    TechpowerUp GPU-z Running.

Fig. 3.    (Performance) Execution Time for Matrix Multiplication in Dual Level Models vs HMC.



Fig. 4.    (Performance) Execution Time for Matrix Multiplication in MOC Model vs HMC.



Fig. 5.    (Power consumption) Dual Level Models vs HMC.



Fig. 6.    (Power Consumption) MOC Tri-Hybrid vs HMC Model.

This time we observed a vital difference in power consumption throughout the executions. For small dataset, the observed power consumption was up to 80 watts, which is big difference and an achievement in terms of study's second objective. The same ratio was discovered for all other executions. Finally, HMC consumed 240 watts for 6000x6000 matrix size. We critically noted that this level of power consumption was observed in 2000x2000 for the MOC model, which was computed in a very small time. However, a drastic change in both objectives was achieved in proposed HMC model, which is a big achievement with respect to satisfying the requirements of emerging HPC systems.

## V.    DISCUSSION

The proposed study was primarily concerned with emerging High-Performance Computing and its perspectives objectives, which are majorly concerned to enhancing performance under the power consumption limitations. These concerns are vital challenges now a day for current and future ICT. According to research communities, there are two solutions to these primary challenges, increasing number of cores in the system to achieve massive performance in the system. This approach is not feasible, as it will increase the power consumption in the system, there another solution is required which is achieving massive parallelism in the system to reduce the execution time that will eventually decrease the power consumption in the system. Leading to the second option, this study proposed a new parallel programming model called Hybrid MVAPICH-2 and CUDA (HMC).

HMC is fundamentally an extension of the dual level model of MPI and CUDA. The issue in a hybrid of MPI+CUDA was similar that it could not fulfill the demand of HPC systems. MPI-3 (MVAPIC2) it has more directives like sessions and groups and it will be scalable. We observed that the quantified execution in HMC was 10% less as compare the other parallel programming models. Further, we also noted that the power consumption was much consumed in dual and tri-level hybrid models but HMC consume 50 to 60 WATTs less comparatively.

## VI.    CONCLUSIONS

Emerging HPC technologies are experiencing more priority and demand in all scientific fields. It has been anticipated that Exascale HPC systems will be introduced at the end of 2020. According to current state-of-the-art technologies, Exascale systems face two vital challenges, including Power consumption when increasing system performance to achieve Exa-flops level of calculations. Various research communities have taken initiatives to address these challenges. With respect to these objectives, the current study proposed a new model named Hybrid MVAPICH-2+CUDA (HMC) to address such challenges. The HMC model implemented in matrix multiplication benchmarking application and compared the quantified performance and power consumption with existing dual- and tri-hybrid parallel programming models. We observed that the HMC model outperformed in all cases when we compared it to other dual- and tri-level parallel programming models. HMC reduced the power consumption up to 80 watts with the same dataset execution within 70 sec less time,

comparatively. These improvements can serve as the foundation of an initiative to consider HMC as a leading model in the era of HPC systems.

From future perspectives, HMC is required to implement a large cluster system through which we can quantify the said attributes on different platforms. Moreover, we must implement HMC on different benchmarking applications to observe the behavior of the proposed model when we change the benchmark.

REFERENCES

[1] Perarnau, Swann, Rinku Gupta, and Pete Beckman. "Argo: An Exascale Operating System and Runtime." (2015).

[2] Carretero, Jesus, et al. "Energy-efficient Algorithms for Ultrascale Systems."Supercomputing frontiers and innovations 2.2 (2015): 77-104.

[3] Shalf, John, Sudip Dosanjh, and John Morrison. "Exascale computing technology challenges." High-Performance Computing for Computational Science–VECPAR 2010. Springer Berlin Heidelberg, 2010. 1-25.

[4] Abraham, Erika, et al. "Challenges and Recommendations for Preparing HPC Applications for Exascale." arXiv preprint arXiv:1503.06974 (2015).

[5] Hall, Mary, et al. "ASCR Programming Challenges for Exascale Computing." (2011).

[6] Tolentino, Matthew, and Kirk W. Cameron. "The optimist, the pessimist, and the global race to exascale in 20 megawatts." Computer 45.1 (2012): 0095-97.

[7] Diaz, Javier, Camelia Munoz-Caro, and Alfonso Nino. "A survey of parallel programming models and tools in the multi and many-core era." Parallel and Distributed Systems, IEEE Transactions on 23.8 (2012): 1369-1386.

[8] Rajamony, Ramakrishnan, and Alan L. Cox. "Parallel programming tools." Wiley Encyclopedia of Electrical and Electronics Engineering (1998).

[9] W. Hwu, K. Keutzer, and T.G. Mattson, "The concurrency challenge," IEEE Design and Test of Computers, vol. 25, no. 4, pp. 312-320, July 2008.

[10] Macedonia, Michael. "The GPU enters computing's mainstream." Computer 36.10 (2003): 106-108.

[11] Geer, David. "Chipmakers turn to multicore processors." Computer 38.5 (2005): 11-13.

[12] Satish, Nadathur, Mark Harris, and Michael Garland. "Designing efficient sorting algorithms for manycore GPUs." Parallel & Distributed Processing, 2009. IPDPS 2009. IEEE International Symposium on. IEEE, 2009.

[13] Nakajima, Kengo. "Hybrid vs. flat mpi on the earth simulator: Parallel iterative solvers for finite-element method." Applied Numerical Mathematics 54.2 (2005): 237-255.

[14] Geist, Al, and Robert Lucas. "Major computer science challenges at exascale." The International Journal of High- Performance Computing Applications 23.4 (2009): 427- 436.

[15] Bergman, Keren, et al. "Exascale computing study: Technology challenges in achieving exascale systems." Defense Advanced Research Projects Agency Information Processing Techniques Office (DARPA IPTO), Tech. Rep 15 (2008).

[16] Reed, Daniel A., and Jack Dongarra. "Exascale computing and big data." Communications of the ACM 58.7 (2015): 56-68.

[17] Bergman, Keren, et al. "Exascale computing study: Technology challenges in achieving exascale systems peter kogge, editor & study lead." (2008).

[18] Maglione, Stephen C., and Edward Stanley Suffern. "Power control of servers using advanced configuration and power interface (ACPI) states." U.S. Patent No. 8,250,382. 21 Aug. 2012.

[19] Feng, Xixhou, Rong Ge, and Kirk W. Cameron. "Power and energy profiling of scientific applications on distributed systems." Parallel and Distributed Processing Symposium, 2005. Proceedings. 19th IEEE International. IEEE, 2005.

[20] Marathe, Aniruddha, et al. "A run-time system for power- constrained HPC applications." International conference on high-performance computing. Springer, Cham, 2015.

[21] Geist, Al, and Daniel A. Reed. "A survey of high- performance computing scaling challenges." The International Journal of High-Performance Computing Applications 31.1 (2017): 104-113.

[22] Ashraf, M. Usman, Fathy Alboraei Eassa, and Aiiad Ahmad Albeshri. "High-performance 2-D Laplace equation solver through massive hybrid parallelism." 2017 8th International Conference on Information Technology (ICIT). IEEE, 2017.

[23] Ashraf, M. Usman, et al. "Performance and power efficient massive parallel computational model for HPC heterogeneous Exascale systems." IEEE Access 6 (2018): 23095-23107.

[24] Ashraf, Muhammad Usman, Fadi Fouz, and Fathy Alboraei Eassa. "Empirical Analysis of HPC Using Different Programming Models." International Journal of Modern Education & Computer Science 8.6 (2016).

[25] Ashraf, Muhammad Usman, and Fathy Elbouraey Eassa. "Hybrid model-based testing tool architecture for an exascale computing system." International Journal of Computer Science and Security (IJCSS) 9.5 (2015): 245.

# Effective Framework of Pedagogy

## Using the Plan-Do-Check-Act Improvement Cycle

Tallat Naz[1], Momeen Khan[2], Khalid Mahmood[3]
Department of Computer Science, IIC University of Technology
Phnom Penh, Cambodia

*Abstract*—Learning paths drive learners to proficiency by using a selected sequence of training activities under time constraints. Therefore, learners can regulate learning and give feedback for pedagogy improvements. Studying learning path evaluation provides a useful conceptual reference to enhance pedagogically. This paper proposes an approach based on the Plan-Do-Check-Act improvement cycle to systematically evaluate learning paths in learning management systems. The framework is a valuable resource that consolidates existing practices in learning management evaluation. Our approach integrates learning styles, learning profile, along with cognitive activities. The proposed framework was compared with current learning path methods. Results were competitive compared with related works.

*Keywords—Learning management system; learning styles; learning path; Plan-Do-Check-Act (PDCA)*

## I. INTRODUCTION

A learning management system (LMS) provides functions of course creation tools, tracking, assessment, collaboration, and reporting to learners and administrators. An LMS provides content managers with the tools to create learning paths for specific course outcome. A learning path is a guided set of educational activities that aim to increase learner proficiency; thus, it is considered an organizational asset. A learning path should consider the current learner skillset, the required job skills, and the future learner skillset to fulfill his career goals.

A useful LMS should provide facilities to create and track the learning paths of prospective employees. Therefore, it aligns job skillset with resources; consequently, the organization LMS value is increased. However, building a learning path without taking personal preferences, cognitive activities, and learning style into consideration could lead to improper resources required to carry out the LMS project tasks. Without considering the overall value of the learning path to an organization, the learning path assessments remain provisional. Therefore, a non-achieved course or a seminar outcome could lead to organization goal failure.

Although there are many works on learning path creation and optimization [1]–[4], little work exists on learning path evaluation from the viewpoint of teachers. An assessed learning item of a learning path provides little feedback on current courses content or course items sequence. Therefore, the problem learning path evaluation is seen as a personal activity of learners, leaving the teacher assessment out of scope.

In practice, a learning path evaluation process should consider multidimensional data including learning profile, preferences, and behavior actions. A useful learning profile helps learner strive to be inquirers, knowledgeable, communicators and reflective. Learner preferences are based on their expected learning targets, timeframes, and budget constraints. Behavior actions include history and log files, learner recommendation, and learner usability records of an LMS.

We claim that enhancing pedagogy is affected by Porter's five forces [5] in the context of e-learning systems. The competitive rivalry of other active learning paths; the threat of adding extra, unrelated, or improper learning items sequences; the threat of substituting new learning seminars; and the bargaining power between learner and the instructor.

This paper applies Deming's plan-do-check-act (PDCA) cycle [6] (as shown in Fig. 1) to improve and evaluate the learning paths of LMSs. At the planning phase, the course outcomes are linked to a set of processes toward creating new or improved course content. Then the set of activities in the planning phase is executed in the "do" phase. The "check" phase evaluates the learning path based on available data. Consequently, the "act" phase fixes issues and recommends actions that need human involvement. The PDCA iterates until satisfaction determined by the evaluator. Table I. shows a list of the PDCA actions under PDCA phase. In each phase of the PDCA, a set of actions is proposed that systematically combine learning profile, the existing log file of LMS, and learning styles. Therefore, the statistics of learners and teachers are combined.

The paper is structured as follows. Section 2 summarizes related work. Section 3 explains the proposed model, while Section 4 evaluates the proposed model. Section 5 provides conclusions and future research.

Fig. 1. Learning Path Plan-do-Check-Act Cycle.

TABLE I. THE PDCA STAGES IN THE LEARNING PATH EVALUATION PROCESS

| Phase | Action |
|---|---|
| **Plan** | develop a learner personnel profile identify the learning style apply cognitive activities process LMS repository datasets acquire organization assets |
| **Do** | the previous step is enacted |
| **Check** | course objectives coverage sequencing of learning activities course learning and practice time relationship or synthesis within learning provided by other sources. |
| **Act** | extend or improve course objectives alter sequencing of learning activities modify required practice time address learning of knowledge within the tasks or activities of learning path identify gaps in the learning process compared with other sources. |

## II. RELATED WORK

### A. Learning Management System Evaluation

In the market, there are more than 300 active LMSs that provide basic features like content options, course creation tools, assessments, collaboration, reporting, and skill tracking. Many LMSs provide mobile learning, certification management, gamification, and social learning. One major feature of an LMS is the Shareable Content Object Reference Model (SCORM). SCORM defines a way of constructing training content so it can be shared with other SCORM compliant systems. LMSs can be in three main categories: content preparation systems like Moodle and Blackboard; corporate training systems like Litmos, and Zoho recruit; and school management systems like Edmodo, and Schoology. A useful LMS should be scalable, user-friendly with simple reporting features. The progress tracking and evaluation of learning goals ensure effective pedagogy.

There are many approaches to evaluate LMSs. Gartner uses a data-driven approach known as the magic quadrant. The quadrant has two dimensions the capability and value. The capability focuses on an LMS capability such as functionality and system integration while the value provides dimensions of

satisfaction and market price [7][8]. The quadrant dimension identifies learners, masters, pacesetters, and contenders.

The adoption of an LMS has been further studied [9] to identify technology and pedagogy effect on the learning process. The study identified a lack of development in LMS usage and pedagogy. The predictors of behavioral intention to use an LMS are identified as perceived usefulness, perceived ease of use, enjoyment, subjective norm, satisfaction, and interactivity and control with the validated structural model [10]. Author in [11] identified that tracking, deficiency of audit trail, and insufficient reporting are the main gaps in current LMS systems.

The adoption of an LMS depends on technology, system usefulness, and organizational constraints. Technology issues such perceived use and usefulness are the main factors in an LMS adoption [12]. However one of the most widely used models, the technology acceptance model (TAM) and its variants are problematic with conflicting results [12]. A unified theory of acceptance and use of technology (UTAUT) model was proposed to handle the limitation of a single technology model [13].

Other approaches try to evaluate an LMS system from the utility function. The reference [14] handled the satisfaction and usability factors of LMS using the Analytical Hierarchy Process (AHP) approach to determine satisfaction factors (accuracy, timeless, ease of use and format). They found that the design of an LMS interface will enrich the effectiveness of pedagogy. They identified interface problems with their root cause using the ISO 9241-11 model [15]. Other approaches evaluate the usability of software, in general, using software comments [16], [17]. Teachers believe that an LMS is useful in communication and improving learning while students barely believe that an LMS usage improves teaching [18]. Therefore, organizational support is needed to provide pedagogy. Reported results support the hypotheses that organizational support plays a primary role in enhancing the faculty's LMS self-efficacy and technical support [19].

### B. Learning Styles

Learning styles provide common ways of learning. Some people prefer using pictures while others prefer to use music. VARK is one of the leading models of learning styles that refers to the learning styles: visual, auditory, reading/writing preference, and kinesthetic [20]. A right LMS should provide cognitive and self-directing skills [4]. Personality evaluation could be used for guiding the students according to their preferences [21]. There is a direct relationship between learning styles and personality traits [21]. The reference [22] proposed a neural network model to detect learning styles. The reference [23] focused on recommending learning contents with an adaptive user interface. The course interface was changed using the Felder-Silverman Learning Style Model (FSLSM) model and generic adaptive rules.

### C. Learning Paths

One possible way to evaluate the learning path is by using a checklist guide. System reports about collected data of learner activities of keystrokes and mouse clicks is another approach. However, other approaches are based on previous knowledge

and log file of online activities [1]; therefore, the seminar taken adequately by learner interacts with individual preference [1]. The reference [2] proposed a genetic algorithm to sequence learning path based on learner preferences. Learning analytics provides tools for collecting analyzing data about learners to optimize learning. The learning processes could be understood and managed using learning analytics approach [3]. Therefore, learner preferred learning modes and modality preferences influence their learning and capabilities. The preferred learning modes should be matched with an appropriate learning path.

### III. PROPOSED FRAMEWORK

Learners can self-regulate their learning based on cognitive and behavioral dimensions. However, learners are influenced by external factors of seminar content style (teacher style) compared to his (internal) learning style. Therefore, the proposed model takes into consideration the learner style and his profile.

Fig. 2 depicts the proposed model. The model is composed of two significant steps the "plan-do" and "check- act" of the learner path cycle. The "develop profile" process captures information from the LMS repository including previous courses taken, learner actions, and the LMS user behavior. The process also is subject to organizational process assets that could be part of the LMS. Organizational process assets include policies, procedures, learning process, lesson learned, and other knowledge that is deemed appropriate for the organization such as an HR system.

Therefore, taking the learning outcome as the critical part of the learning path evaluation, then we express the learning path effectiveness as shown in equation (1).

$$Effectivness = \frac{\sum weight_i . Acc_i}{\sum weight_i}$$

The second process, "identify learning style", is crucial to pedagogy and an active learning path. Learners should be comfortable with learning contents that matches with their cognitive abilities. In this process, using available learning styles such as VARK and FSLSM, the leaner should be categorized to his preferred content. For example, a visual learner style implicates the need for visual course contents. Visual learning retains twice as much information as those in the auditory condition [24]. The identification of a leaner style is affected by factors such as organization, cognitive abilities, and technical constraints [25]. Several lite tests should be carried out to identify the learner style.

The primary critical step of the proposed framework is the "evaluate learning path" shown in the second phase of the framework. The process has several sub-processes as stated in Fig. 3(a) to Fig. 3(d). Table II summarizes these sub-processes and the "act" that we propose to use.

Where $weight_i$ is the weight of the outcome $i$, $Acc_i$ is the accuracy of a completed outcome $i$. The weight of an outcome is calculated based on the outcome importance or time needed to accomplish the outcome. A teacher should score each outcome value to quantify pedagogy. However, the terms in

equation (1) depends on quantifying the prescriptions and actions in Table III which is outside the scope of this paper.



Fig. 2. Proposed Model.

TABLE II. LEARNING PATH EVALUATION (THE "CHECK-ACT" OF THE PDCA)

| Figure | Description | Action |
|---|---|---|
| 3(a) | Was the course expected outcomes covered in the learning path? | If the average learners' score for the learning path is less than 50%, then prompt an immediate update to course outcomes or contents. Start a process to address gained learning knowledge. |
| 3(b) | Check learning path sequence of activities (curriculum)? | If average learners score of a subset of learning path is below 50%, then consider better sequencing or improving course content. |
| 3(c) | Check the alignment between course contents and allotted time? | If the average number of learners fails to complete the course item in allotted time, then consider either extending time or splitting Course item. |
| 3(d) | Does the learning path do better than other competitors? | If the number of enrollments in the current learning path of an LMS is less than a competitor LMS, then consider a revision of sequence, content, price, Enrollment process and course features. |



Fig. 3. Evaluating Learning Path Process (Check-Act).

## IV. EVALUATION

Although many works have been adopted to evaluate a learning path, most works target learning path design. To our knowledge, no complete learning path framework quantify learning path from learner and teacher perspectives. The proposed framework is an overarching approach intended to measure the learning path effectiveness from the learner and the organization (teacher) point view.

We carry out a comparison against the following list of learning path evaluation features that are either extracted from the literature review or suggested by this paper. The suggested criteria are as follows:

*1) Agility:* The ability to be flexible and able to deal with unseen changes in technology, environment, and pedagogy.

*2) Overall performance:* The ability to measure the overall performance of a specific learning path.

*3) Identify course content gaps:* The framework should be able to detect if learning material should be changed

*4) Learner profile:* An overall framework should be able to take into consideration the learner profile including the cognitive and learning styles.

Table III shows the comparison of our approach and a list of selected works from the literature. The ☑ indicates that the criterion is fully available while ☒ means a non-applicable criterion.

From the table, we deduce that the proposed model and selected works are agile, since they are not linked to specific learning management platform. All the studied evaluation methods do not provide the ability to evaluate the learning path from the organization and learner perspectives. The works [1] and [3] did not provide feedback on learning course items or feedback back to expected learning outcomes. The lowest performance was in LMS reports that only use log files of learner actions. The proposed approach reported the highest performance due to its abstractness and richness of extra knowledge extracted from learning styles, preferences, and cognitive activities.

The proposed framework was evaluated experimentally by expert judgment. Teachers liked the idea, but they were looking for a real-life system while learners were happy with the feature of learning style detection. Moreover, one organization that adopts an LMS likes the idea of linking organization assets with the LMS learning path.

TABLE III.    COMPARISON OF OUR APPROACH AND RELATED APPROACHES

| Criterion | Our approach | Checklist | LMS reports | [1] | [2] | [3] |
|-----------|-------------|-----------|-------------|-----|-----|-----|
| A | ☑ | ☒ | ☑ | ☑ | ☑ | ☑ |
| B | ☑ | ☒ | ☒ | ☒ | ☒ | ☒ |
| C | ☑ | ☑ | ☒ | ☒ | ☑ | ☒ |
| D | ☑ | ☑ | ☒ | ☑ | ☑ | ☑ |

## V. DISCUSSION

This study extends previous work on learning path evaluation by providing a theoretical framework that specifies the unidirectional feedback between learners and teachers. It also offers several important implications for research and practice, and thus should help in the design, evaluation and widespread adoption in LMSs.

Despite the significant contribution of the research findings, this study involved limitations that should be considered when interpreting the results. First, the features of influencing learning styles evaluations are primarily identified based on a limited number of literature reviews from 2010-2018. This limitation may restrain the generalizability of the findings, and hence, future researchers should interpret the current findings with caution. Second, the proposed approach has not been implemented in an LMS. Therefore, results probably only reflect the general concept of a learning path evaluation. Future work should consider the application of machine learning approaches to analyze and categorize collected data, before proceeding to implementation.

## VI. CONCLUSIONS

This paper proposes a learning path framework that can measure the effectiveness of a learning path. The proposed model evaluates the learning path from the viewpoint of learners, teachers, and organization. The model is based on learner style, profile, and cognitive activities. The evaluation was carried out by comparison of the proposed model with a list of related works. Results show that the framework is deemed useful in the context of state of the art LMSs. In the future, we plan to implement the framework in an open source LMS such as Moodle.

### REFERENCE

[1] T. Lerche and E. Kiel, "Predicting student achievement in learning management systems by log data analysis," Comput. Human Behav., vol. 89, pp. 367–372, 2018.

[2] P. Dwivedi, V. Kant, and K. K. Bharadwaj, "Learning path recommendation based on modified variable length genetic algorithm," Educ. Inf. Technol., vol. 23, no. 2, pp. 819–836, 2018.

[3] C. Schumacher and D. Ifenthaler, "Features students really expect from learning analytics," Comput. Human Behav., vol. 78, pp. 397–407, 2018.

[4] M. Rani, K. V. Srivastava, and O. P. Vyas, "An ontological learning management system," Comput. Appl. Eng. Educ., vol. 24, no. 5, pp. 706–722, 2016.

[5] M. E. Porter, "The five competitive forces that shape strategy," Harv. Bus. Rev., vol. 86, no. 1, pp. 25–40, 2008.

[6] R. D. Moen and C. L. Norman, "Circling back: Clearing up myths about the Deming cycle and seeing how it keeps evolving," Qual. Prog., vol. 43, no. 11, p. 22, 2010.

[7] Gartner® Inc, "FrontRunners Methodology," 2018. [Online]. Available: https://www.saimgs.com/imglib/other_pages/FrontRunners/Meth odologyOverview.pdf. [Accessed: 28-Feb-2019].

[8] Gartner, "Magic Quadrant for Business Intelligence and Analytics Platforms, Retrieved on 15 October, 2017 from https://www.gartner.com/ doc/reprints?id=1- 3TYE0CD&ct=170221&st=sb," 2017.

[9] J. Sinclair and A.-M. Aho, "Experts on super innovators: understanding staff adoption of learning management systems," High. Educ. Res. Dev., vol. 37, no. 1, pp. 158–172, 2018.

[10] D. Findik-Co\cskunçay, N. Alki\cs, and S. Özkan-Yildirim, "A Structural Model for Students' Adoption of Learning Management Systems: An Empirical Investigation in the Higher Education Context," J. Educ. Technol. Soc., vol. 21, no. 2, pp. 13–27, 2018.

[11] "Functionality Gaps in the Design of Learning Management Systems," Int. J. Adv. Comput. Sci. Appl.

[12] R. R., "Why some teachers easily learn to use a new virtual learning environment: A technology acceptance perspective," Interact. Learn. Environ., vol. 24, p. 539, 2014.

[13] V. Venkatesh, J. Y. L. Thong, and X. Xu, "Consumer acceptance and use of information technology: extending the unified theory of acceptance and use of technology," MIS Q., vol. 36, no. 1, pp. 157– 178, 2012.

[14] Y. Tjong, L. Sugandi, A. Nurshafita, Y. Magdalena, C. Evelyn, and N. S. Yosieto, "User Satisfaction Factors on Learning Management Systems Usage," in International Conference on Information Management and Technology (ICIMTech 2018), 2018, pp. 11–14.

[15] N. Phongphaew and A. Jiamsanguanwong, "Usability evaluation on learning management system," in International Conference on Applied Human Factors and Ergonomics, 2017, pp. 39–48.

[16] I. Atoum, "A Novel Framework for Measuring Software Quality- in-use based on Semantic Similarity and Sentiment Analysis of Software Reviews," J. King Saud Univ. - Comput. Inf. Sci., p. , 2018.

[17] I. Atoum and C. H. Bong, "A Framework to Predict Software 'Quality in Use' from Software Reviews," in Proceedings of the First International Conference on Advanced Data and Information Engineering (DaEng-2013), vol. 285, J. Herawan, Tutut and Deris, Mustafa Mat and Abawajy, Ed. Kuala Lumpur: Springer Science+Buisness Media Singapore, 2014, pp. 429–436.

[18] T. S. D., "Saving time or innovating practice: Investigating perceptions and uses of learning management systems," Comput. Educ., vol. 53, p. 686, 2009.

[19] Y. Zheng, J. Wang, W. Doll, X. Deng, and M. Williams, "The impact of organisational support, technical support, and self- efficacy on faculty perceived benefits of using learning management system," Behav. Inf. Technol., vol. 37, no. 4, pp. 311– 319, 2018.

[20] N. D. Fleming, "I'm different; not dumb. Modes of presentation (VARK) in the tertiary classroom," in Research and development in higher education, Proceedings of the 1995 Annual Conference of the Higher Education and Research Development Society of Australasia (HERDSA), HERDSA, 1995, vol. 18, pp. 308–313.

[21] A. Kamal and S. Radhakrishnan, "Individual learning preferences based on personality traits in an E-learning scenario," Educ. Inf. Technol., vol. 24, no. 1, pp. 407–435, 2019.

[22] M. S. Hasibuan, L. E. Nugroho, and P. I. Santosa, "Model detecting learning styles with artificial neural network," J. Technol. Sci. Educ., vol. 9, no. 1, pp. 85–95, 2019.

[23] S. V. Kolekar, R. M. Pai, and M. P. Manohara, "Rule based adaptive user interface for adaptive E-learning system," Educ. Inf. Technol., 2018.

[24] J. Cuevas and B. L. Dawson, "A test of two alternative cognitive processing models: Learning styles and dual coding," Theory Res. Educ., vol. 16, no. 1, pp. 40–64, 2018.

[25] S. Graf, T.-C. Liu, and others, "Identifying learning styles in learning management systems by using indications from students' behaviour," in Eighth IEEE international conference on advanced learning technologies, 2008, pp. 482–486.

# Comparative Performance Analysis of RPL for Low Power and Lossy Networks based on Different Objective Functions

Mah Zaib Jamil[1], Danista Khan[2], Adeel Saleem[3], Kashif Mehmood[4], Atif Iqbal[5, *]

Department of Electrical Engineering, University of Engineering and Technology, Lahore, 54000, Pakistan[1]
Department of Electrical Engineering, The University of Lahore, Lahore, 54000, Pakistan[2, 3, 4]
School of Electrical & Electronic Engineering, North China Electric Power University, Beijing, 102206, China[3]
School of Electrical Engineering, Southeast University, Nanjing, 210096, China[4]
School of Renewable Energy & Clean Power, North China Electric Power University, Beijing, 102206, China[5*]

*Abstract*—**The Internet of Things (IoT) is an extensive network between people-people, people-things and things-things. With the overgrown opportunities, then it also comes with a lot of challenges proportional to the number of connected things to the network. The IPv6 allows us to connect a huge number of things. For resource-constrained IoT devices, the routing issues are very thought-provoking and for this purpose an IPv6 Routing Protocol for Low-Power and Lossy Networks (RPL) is proposed. There are multi-HOP paths connecting nodes to the root node. Destination Oriented Directed Acyclic Graph (DODAG) is created taking into account different parameters such as link costs, nodes attribute and objective functions. RPL is flexible and it can be tuned as per application demands, therefore, the network can be optimized by using different objective functions. This paper presents a novel energy efficient analysis of RPL by performing a set of simulations in COOJA simulator. The performance evaluation of RPL is compared by introducing different Objective functions (OF) with multiple metrics for the network.**

*Keywords*—*ETX; ELT; HOP; internet of things; IP; networks; network performance; routing; RPL*

## I. INTRODUCTION

Radio-frequency identification (RFID) group generally define Internet of Things (IoT) as a network in which objects are globally interconnected and accessible by their unique addresses based on standard communication protocols [1]. The vision and applications of Internet of Things have been promoted by the intercommunication of simple embedded wireless sensing devices as depicted in [2]. The data communication in wireless networks takes place over multiple wireless links i.e. multi-HOP communication. Routing protocol is responsible for finding the best path towards the destination with minimum cost.

For delay resistive communication, routing protocol is responsible for minimizing the overhead time required to find a path and ensuring an energy efficient communication. For Low power and Lossy Networks (LLNs), a light-weight, energy and memory efficient routing protocol for resource constrained IoT named as the Routing Protocol for Low-Power and Lossy

Networks (RPL) is proposed by IETF ROLL group [3]. In RPL, the objective function is either maximized or minimized depending upon the node's metrics which are shared among the nodes. The RPL routing protocol is designed to be highly flexible, which mandates its tuning for application specific requirements. Recent studies on RPL target the techniques to improve the energy efficiency of RPL. For this reason, many objective functions and metrics are proposed in prior literature. The preferred parent nodes are overburdened due to more than one child node connected to each parent node, thus resulting in weakening of these nodes [4]. To cope up with this, the greedy approaches are applied to choose the best possible parent in RPL which result in frequent path changes, maintenance overhead and instability of significant routing protocol. These issues affect the energy efficiency and have a crucial impact on the lifetime of the nodes, resulting in the disconnection of a part of the network and therefore the reconstruction of the Destination Oriented Directed Acyclic Graph (DODAG) [4]. Thus, these approaches decrease energy efficiency.

This paper highlights this problem by carrying out the performance analysis of RPL based on its stability and energy efficiency. Different objective functions are used for RPL's performance evaluation and the feasibility investigation of the objective functions. The three objective functions are HOP count (HOP), expected transmissions (ETX) and expected lifetime (ELT). The performance of these objective functions is analyzed and compared on the basis of the distances of nodes from the root node. The HOP count and ETX are standardized objective functions in RPL [4]. We have proposed a new novel implementation of ELT objective function, to minimize the overburdening of nodes thus prolonging the node's lifetime. The results of ELT are then compared with the performance results of two standard objective functions of RPL.

The rest of the paper focuses on the analysis of RPL based on different metrics of objective function. Section 2 explains the background of this work. A brief introduction of the RPL is given in Section 3. Algorithm of the paper is discussed in Section 4 whereas the simulation model is presented in Section 5. Performance results are discussed in Section 6. Future modifications and conclusions are elaborated in Section 7.

---

*Corresponding Author

## II. LITERATURE REVIEW

Many research studies have been performed related to the optimization of the RPL in different conditions and within the application's scope. No standardized parameter's configuration is provided by the IETF ROLL, so RPL metrics can be modified according to the application demand. The IETF ROLL working group standardized only two different objective functions, namely, the Objective Function zero (OF0) [5] and the minimum Rank with Hysteresis Objective Function (MRHOF) [3, 6]. In [7, 8], the authors have shown the effectiveness of RPL pertaining to exiguous delay, quick configuration, loop-free topology and self-healing. Efficient energy consumption in RPL majorly depends upon the control messages which are responsible for the creation and selection of the paths. The control traffic overhead is minimized using trickle timer [9]. At the initial stages of the network creation, the number of control messages is high and it decreases as the network reaches towards stability [10]. A 20-nodes based network is simulated and the packet error rate observes to be 1% and the control traffic overhead fluctuate around 25%. Furthermore, the control traffic overhead increases as the number of nodes is increased in the network and is recorded almost 75% for 100 nodes [11].

Control traffic overhead must be reduced to save energy utilization of the nodes and this can be done by adjusting different RPL parameters like trickle timer [12]. In [12], it is also stated that if the routes are maintained by optimizing the energy metric alone then it may result in more packet loss. There should be a compromise between energy metric and the link quality metric to increase the performance of RPL. Duty cycle aware routing scheme is proposed in [13] to save energy. The message packets are transmitted to all the neighboring nodes and the earliest active node forwards the packet. The scheme proposed in [14] is based on multiple node metrics for RPL routing. Factors impacting the energy consumption are also targeted such as the distances of the nodes from the root node. A technique to save energy is the selection of cost e effective path and further reduction in energy utilization can be done by using a probing technique [15]. Recently, in [16] authors reported a serious problem of link repair issue in RPL using a Contiki-RPL implementation. The authors reported that if a link goes down, then the RPL repair process takes at least 200 seconds to find an alternate path.

In [17], the authors reported the instability of RPL and proposed an estimation of control traffic for RPL. However, the proposed scheme mandates that the wireless links are bidirectional and symmetric, which is not a realistic assumption. In [18] a hybrid routing metric is proposed which considers the link conditions and the path lengths based on two dynamic threshold levels. To decrease the energy consumption and increase node's lifetime, the radio transceiver should sleep for a longer duration. Therefore, in [19] a scheme is proposed which considers the link quality and node metrics to improve the objective function for the route selection [19]. In [7, 20] the performance of network for two different objective functions is analyzed using a Contiki-RPL simulator that results in similar outcomes when computing the rank. Several researchers have also tried different methods to optimize routing metrics and OFs for RPL to meet different requirements in specific application scenarios [21, 22]. The topology of the network for the two standard objective functions (HOP count and ETX) is such that the overburdened bottle neck nodes drain more energy as compared to the other nodes which can be chosen as parent nodes. This results in the network failure. The authors proposed a new objective function to balance the overburdened nodes by selecting the other preferred parents to maximize the network lifetime [4].

## III. PROTOCOL OVERVIEW

### A. Routing Protocol for LLNs (RPL)

RPL develops a Directed Acyclic Graph (DAG) topology. The nodes which are preferred parent to the child nodes acts like sink nodes. Hence the child nodes forward packets to their preferred parent until the packet reaches its destination.

The nodes are connected in a network through multiple routes down the tree and if any link failure occurs an alternate path is determined. Thus, path selection plays a vital role in RPL. For this purpose, different metrics are introduced in objective functions as determining the path selection criteria which is communicated in the form of control packets among the network nodes using standard TCP/IP protocols.

### B. Network Topology in RPL

RPL maintains a tree like graph topology in wireless networks. A typical RPL topology is shown in Fig. 1. The top node is the root node, the nodes at the edges are the leaf nodes and rests of the nodes are intermediate parent nodes that connect leaf nodes to the root node. The DAG topology of the RPL maintains direct routes and allows network traffic to flow in both upward and downward directions [23].

### C. Routing Process in RPL

*1) DODAG:* RPL organizes DODAG topology [23]. Usually the DODAG contains only one root node and it keeps the information of the paths from the nodes towards the root node. Each node has a set of preferred parents and any one of the preferred parents is selected by the node for packet forwarding to the destination node. The node obtains the information required for the selection of the parent from the DODAG's data it receives. A rank is allotted to each node in the network. The rank of the root node is always 1 and the next node closest to the root is assigned the incremented rank.

*2) Control messages:* RPL needs control messages for the creation and the maintenance of the network. These control messages are of the following three types: [19, 24] DODAG Information Object (DIO) message contains the information about the RPL instance, the node parameters like rank and the information needed to identify a DODAG. DODAG Information Solicitation (DIS) is a request message from a node which do not receive DIO message from its neighbor nodes after a specific DIO advertising time interval. It is a method to poke nearby DODAGs [23]. DODAG Destination Advertisement Object (DAO) message carries the destination information in the upward direction to the root node. This message is used to establish downward paths.

*3) Control messages flow and DODAG formation:* The root node sends DIO messages to the nearby nodes. The nearby nodes upon processing the DIO select the root as their parent due to the lowest rank and sends DIO messages to the neighbor nodes. The nodes with a lower rank than the node which sends DIO ignore the message and the rest of the nodes process the DIO message for parent selection. If any node lies in more than one neighbor list of the nodes, then it receives more than one DIO message and parent selection is governed by the information in the objective function. This process repeats for all nodes. When all nodes join DODAG then child nodes start sending DAO messages in the upward direction to their respective parents for downward communication paths. This process repeats for every child node until the root node is reached. Meanwhile, any node upon not receiving DIO message, within a specified time interval, sends DIS messages to the neighbor nodes for seeking DIO message to join the DODAG [25].

*4) Objective functions for routing:* The Objective function starts playing its role when a node receives more than one DIO message from nearby preferred parents. The objective function is defined as a criterion for the route selection [25]. The computation of rank and the selection of parents are also governed by the objective function. It starts calculating and optimizing the path cost of every possible path choice.

According to the application demand, the objective function is either maximized or minimized. Generally, there are two objective functions which are implemented in different operating systems, namely Objective Function Zero (OF0) and Expected Transmissions Count (ETX). The optimization criterion for OF0 is HOP count and for ETX criteria is the expected number of transmission attempts. OF0 minimizes the number of the nodes between the sender node and the root node while ETX ensures minimum time for the packet to be delivered successfully to the root node [26]. The objective function is based on routing metrics and dictates the network operation.

*5) Metric of objective function for routing:* Metric of the objective function is quantified as the path cost and is responsible in the selection of the path. For example, routes can be selected based on link metrics, i.e. ETX, or based on node metrics i.e. battery status. The routes which do not match the metric criteria in the objective function are omitted.



Fig. 1.   RPL Topology.

## IV. PERFORMANCE ANALYSIS TECHNIQUES

### A. RPL Objective Function

The OF along with a set of metrics are used for the selection of DODAG, computing the ranks of the nodes, listing the preferred parents and selection of routes based on the metrics [23]. In RPL routing protocol, every child has to choose a preferred parent among the neighboring nodes towards the root node, depending upon the information available in objective function [25].

### B. Minimum Rank with Hysteresis Objective Function (MRHOF)

MRHOF is an objective function based on the minimization of the metric defined in the objective function. It omits the routes which do not fulfill the criteria of the objective function. Meanwhile the response to small metric changes is reduced using hysteresis. The metrics used in MRHOF are additive along a route. The two most commonly used routing metrics in RPL are ETX (expected number of transmission attempts per packet) and HOP (HOP count) [18, 26].

MRHOF19 with the ETX metric selects the paths of nodes to the root node with minimum number of retransmission value and MRHOF with HOP count metric select the path based on the minimum number of HOP counts of nodes towards the root node [18, 26]. ETX metric is the wireless link quality and is approximated by averaging the number of successful transmissions of the data frames [18]. HOP Count metric gives the rank of the node which is the distance of the node from the root node in terms of number of HOPs [26].

### C. ELT

In this paper, three different routing metrics in the objective function are used to evaluate the performance of the RPL. These routing metrics are HOP Count, ETX and ELT. The former two are already implemented and the latter one is now proposed for the comparative analysis of RPL. The ELT is comparatively new metric [27] which is expected lifetime metric of the node based on the energy consumed of the nodes in active mode (Rx and Tx).

To implement ELT of a node, it is calculated using the formula proposed by the authors in [28].

$$ELT = E_{res} / E_{use} \qquad (1)$$

Equation (1) shows that the expected life on a node is the ratio of the residual energy ($E_{res}$) to the energy used ($E_{use}$). The following approximation is used to simplify it

$$E_{use} = Time \times P_{tx} \qquad (2)$$

Combining (1) and (2), ELT becomes as shown in (3).

$$ELT = E_{res} / (Time \times P_{tx}) \qquad (3)$$

$P_{tx}$ is a transmission power of a particular mote with a constant value. Time can be approximated as (4) shown below

$$Time = (TotTx \times PacketSize)/DataRate \qquad (4)$$

$TotTx$ is the total number of transmissions. Now using (4) and (3), it becomes

$$ELT = Eres / (((TotTx \times PacketSize) /DataRate) \times Ptx) \quad (5)$$

So, (5) is the final form of the ELT.

In ELT based objective function, the child node chooses the parent among the neighboring nodes depending upon the longest lifetime of the nodes. ELT ensures that the network lifetime is maximized. The goal of ELT is to maximize the life as in (6) or we can minimize the time as in (7).

$$max\ ELT = max\ Eres / (Time \times Ptx) \quad (6)$$

$$max\ ELT = Eres / min\ (Time \times Ptx) \quad (7)$$

Equation (7) shows that maximizing the ELT is equivalent to the minimization of the Time. From (4), for a fixed data rate for all the transmissions, the minimization of Time is proportional to the minimization of the number of transmissions as in (8).

$$max\ ELT =$$

$$Eres / (min ((TotTx \times PacketSize) /\Box DataRate) \times Ptx) \quad (8)$$

So, maximization of the ELT is equal to the minimization of the TotTx subjected to

s.t

Ptx=a

PacketSize=b

DataRate=c

where a, b and c are constants for a particular type of mote and application.

Proposed ELT objective function balances the overburdened nodes by considering the data traffic of the nodes. Hence, fair energy consumption among the nodes is expected using this metric. We have modified the DIO message packet; the new data traffic metric is introduced in RPL. This data traffic metric is updated by the ELT approximation for each node and advertised to the other nodes in DIO packet. Generally, the root sends the first DIO message and the nodes upon receiving this DIO select the root as a preferred parent. Based on the objective function the nodes start calculating their own rank and update DIO message for advertisement to all neighbors. The root ignores the DIOs that come from its child. At this stage, we aim to allow each parent to calculate its total transmissions. The transmissions include data packet sending attempts, the number of control messages and the number of collisions.

In Fig. 2, the child nodes send DIS messages to the parent nodes, if they failed to receive the DIOs from the parent nodes. The parent node updates the information which is used as the decision criteria for the child nodes at the time of parent selection. This information can be nodes rank, channel quality, delay, ETX based on the application demand or network requirement. After updating information, these parameters are broadcasted to the child nodes using DIOs. The parent node receives the DAO message if it is selected by the child node.



Fig. 2. Flow Chart of Parent Nodes.

In Fig. 3, the preferred parent nodes send DIO messages to the child nodes to update the DODAG after each DIO interval. The child node upon receiving DIOs from the parent node extracts the information that is updated based on the metrics of the objective functions of RPL. The selection of the parent node is decided by the child node after comparing the information of all the parent nodes. The decision criterion is governed by the objective function. After the parent selection, the child node sends DAO message to the respective parent node. If any child does not succeed in getting DIO message from the parent nodes, it can send DIS messages to all the parent nodes to ask them to send DIOs.



Fig. 3. Flow Chart of Child Nodes.

## V. EVALUATION PLATFORM AND PERFORMANCE METRICS

Software: Contiki platform and the COOJA Simulator are used for the simulation study [29, 30]. Contiki is the most widely used operating systems for the real-time implementation and simulation of IoT wireless networks. Its development language is C and comes with a lot of features like open source, memory efficient, full IP (IPv4/IPv6) network stacks with standard protocols (UDP, TCP, HTTP, 6lowpan, RPL, COAP), power awareness, radio duty cycling, hardware portable, proto-threads & multi-threads and small footprint (10k RAM & 30k ROM) (Contiki O.S and COOJA simulator http://www.contiki-os.org/) COOJA is a Java-based hybrid approach in terms of the cross-level emulation and simulation tool based on Contiki operating system [30]. It also provides the facility to implement the sensors' software in C language (Contiki O.S and COOJA simulator http://www.contiki-os.org/).

### A. Performance Metrics

The RPL routing protocol is subjected to different standard performance metrics for its evaluation. These are as follows:

*1) Control traffic overhead:* The number of control messages that are DIO, DIS and DAO messages of each node constructs control traffic overhead. Majorly, DIO messages are responsible for the control traffic overhead.

*2) Packet latency:* The difference in the sending time of a packet from a node and the receiving time of the similar packet at the root node is termed as packet latency. It is offered by the path to the packet to reach a root node from a node.

*3) Packet delivery ratio:* The ratio of the number of packets received at the root node to the number of packets sent by a node is called a node's packet delivery ratio.

*4) Energy consumption:* It is the energy utilization of a node during radio on time.

*5) Parent transitions:* Parent transitions are the number of the parent changes by a node on the basis of the information in the objective function.

*6) Total transmissions:* It is the sum of all control messages, payload and re-transmissions of packets.

## VI. EXPERIMENTAL RESULTS

### A. Network Setup

To stimulate the actual lossy environment in COOJA, UDGM model is used. UDGM model has two different ranges as shown in Fig. 4. First one is the transmission range and the other range is interference range. The Tx ratio is set to 100% in the simulations whereas two Rx ratio values: 100% and 80% are used. A network is deployed, having 20 nodes with one node as a root/server node and 19 nodes are used as client nodes. The client nodes sense data and transmit it to the root node. The locations of the network nodes are recorded and the distances of the nodes from the root node are calculated. The locations are shown in Fig. 5. The sensor node periodically transmits UDP packets. The sensor nodes use RDC (radio duty cycling) to minimize energy consumption. Simulations run

time is one hour. The other simulation parameters and their values are given in Table I.

### B. Simulations and Performance Analysis

The network is simulated for all the objective functions one by one. To obtain the results, the network is simulated multiple times for a single performance parameter and the average of all the simulation results is obtained for HOP, ETX and ELT one by one. Log files are used for extracting the data. The simulation instance is shown in Fig. 6. The network is simulated for one hour and results are obtained by averaging the results of simulations.



Fig. 4. UDGM model with two different ranges



Fig. 5. Network of nodes with their locations

TABLE I. SIMULATION PARAMETERS

| Simulation Parameters | Value |
|---|---|
| OF | Hop, ELT, ETX |
| Rx Ratio | 80%, 100% |
| Tx Ratio | 100% |
| Tx Range | 60m |
| Interference Range | 100m |
| Simulation Time | 1 hr. |
| No. of Nodes | 1 server node, 19 client nodes |
| Packet Sending Interval | 600ms approx. |

Fig. 6.    COOJA Simulator Environment.

*1) PDR vs distance:* The PDR is plotted versus distances of the nodes from the root node as shown in Fig. 7. The PDR for HOP, ELT and ETX is quite similar for nodes closer to the root node in Fig. 7. However, the PDR for ETX is better as compared to HOP for the nodes far from the root node. Since, ETX keeps into account the channel quality and thus offers better probability of packets being successfully transmitted. In case of ELT and HOP the performance of the nodes far from the root node is unfortunately lower in ETX metric based objective functions because ELT and HOP metrics don't consider channel quality or path lengths while selecting parents.

*2) Packet latency:* The packet latency versus distances of the nodes from the root node are plotted as shown in Fig. 8. The packet delay increases when the distance from the root node is increased. HOP offers lower packet delays as compared to ELT and ETX metrics. The control traffic overhead in ETX and ELT is the reason for the higher delays than HOP. However, the HOP metric based objective function offers lower delays due to the shorter path length selection.

*3) Parent transitions:* The number of parent changes by the nodes is directly proportional to the network stability. Therefore, it is of vital importance that the parent changes governed by the routing metric of the objective function should be evaluated. From Fig. 9, the ELT is unstable with more parent changes as compared to HOP and ETX because parent selections are based on the number of transmissions of the nodes. Moreover, the middle nodes have more possible preferred parents than the leaf nodes and the nodes nearer to the root node that is the reason the rate of parent selection is higher for the nodes located in the middle. However, the nodes closer to the root have only a small set of possible preferred parents, thus the parent changes are lesser.

*4) Control traffic:* In control packets, the most frequently sent packet is the DIO packet. In Fig. 10, the control traffic sent by the nodes is shown with respect to the distances of nodes from the root node. The HOP performs well. The middle nodes have a large set of preferred parents so these nodes send more DIO packets as compared to the nodes closer to the root node and the leaf nodes.



Fig. 7.    PDR of HOP, ETX and ELT Plotted Against the Distances of Nodes from the Root Node.



Fig. 8.    Network Latency of HOP, ETX and ELT Plotted Against the Distances of Nodes from the Root Node.



Fig. 9.    Parent Transitions of HOP, ETX and ELT Plotted Against the Distances of Nodes from the Root Node.

Fig. 10. Control Traffic Overhead of HOP, ETX and ELT Plotted Against the Distances of Nodes from the Root Node.

*5) Energy consumption:* The energy consumption of the nodes can be estimated as the number of transmissions of the node which includes the control overhead and payload. As the number of the transmissions increases, the energy consumption by the node is also increased. The energy consumption pattern shown in Fig. 11 depicts that the nodes closer to the root node consume more energy than the leaf nodes. This is due to the fact that a smaller number of nodes are present near the root node to forward the traffic to the root node received from all their respective child nodes. The energy consumption of the nodes in case of ETX is lower, but almost similar to the ELT. However, HOP has low energy consumption due to low PDR and less parent changes.

*6) Total* transmissions*:* The number of transmissions of the nodes can be estimated as the number of successful and unsuccessful packet's sending attempts of the node. The nodes closer to the root node transmit more packets to the root node. However, the leaf nodes face poor packet delivery ratio and highest number of the parent changes so this lead in the increased number of transmissions. The behavior noted for these three metrics are quite similar to the energy consumption performance as shown in Fig. 12.



Fig. 11. Energy Consumption of HOP, ETX and ELT Plotted Against the Distances of Nodes from the Root Node.



Fig. 12. Total Transmissions of HOP, ETX and ELT Plotted Against the Distances of Nodes from the Root Node.

## VII. CONCLUSION

In this paper, a comparative performance analysis of the lightweight routing protocol RPL is presented due to the highly flexible nature of the RPL. A new method is used to implement the ELT objective function and the performance is compared with the other two objective functions (HOP and ETX) with respect to the distances of the nodes from the root node. The literature survey suggests that RPL has been widely accepted for IoT applications and looked up as an energy efficient routing protocol. In this paper, performance analysis of RPL based on its energy efficiency and stability is done using some simulation findings in Contiki OS provided COOJA network simulator. Variations in performance are done by considering the different objective functions and the factors are discussed which are responsible for these variations in the performance metrics of the network.

The performance of network metrics; Network Latency, PDR, Energy Consumption, Control Packet Overhead can be analyzed based on network parameters such as Packet Reception Ratio, Dio Interval Minimum, Dio Interval Doubling, Send Interval. The evaluation of the RPL can be extended to the mobile nodes which will be an enhancement of RPL.

REFERENCES

[1] Atzori L, Iera A, Morabito G, "The internet of things: A survey", Computer networks. 2010 Oct 28;54(15):2787-805.

[2] Zorzi M, Gluhak A, Lange S, Bassi A, "From today's intranet of things to a future internet of things: a wireless-and mobility related view", IEEE Wireless Communications, 2010 Dec;17(6)

[3] Winter T, Thubert P, Brandt A, "RPL: IPv6 routing protocol for low power and lossy networks", draft-ietf-roll-rpl-19, Internet Draft, 2011

[4] Qasem M, Al-Dubai A, Romdhani I, Ghaleb B, Gharibi W, "A new efficient objective function for routing in internet of things paradigm", in Standards for Communications and Networking (CSCN), 2016 IEEE Conference on 2016 Oct 31 (pp. 1-6). IEEE.

[5] Thubert P, "Objective function zero for the routing protocol for low-power and lossy networks (RPL)".

[6] Gnawali O, "The minimum rank with hysteresis objective function".

[7] Accettura N, Grieco LA, Boggia G, Camarda P, "Performance analysis of the RPL routing protocol", in Mechatronics (ICM), 2011 IEEE International Conference on 2011 Apr 13 (pp. 767-772). IEEE.

[8] Qasem M, Altawssi H, Yassien MB, Al-Dubai A, "Performance evaluation of RPL objective functions", in Computer and Information Technology; Ubiquitous Computing and Communications; Dependable, Autonomic and Secure Computing; Pervasive Intelligence and Computing (CIT/IUCC/DASC/PICOM), 2015 IEEE International Conference on 2015 Oct 26 (pp. 1606-1613). IEEE.

[9] Levis P, Clausen T, Hui J, Gnawali O, Ko J, "The trickle algorithm", Internet Engineering Task Force.2011.

[10] Tripathi J, de Oliveira JC, Vasseur JP, "A performance evaluation study of RPL: Routing protocol for low power and lossy networks", in Information Sciences and Systems (CISS), 2010 44th Annual Conference on 2010 Mar 17 (pp. 1-6). IEEE.

[11] Accettura N, Grieco LA, Boggia G, Camarda P, "Performance analysis of the RPL routing protocol", in Mechatronics (ICM), 2011 IEEE International Conference on 2011 Apr 13 (pp. 767-772). IEEE.

[12] Liao CY, Chang LH, Lee TH, Chen SJ, "An Energy-Effciency-Oriented routing algorithm over RPL". 2013

[13] Landsiedel O, Ghadimi E, Duquennoy S, Johansson M, "Low power, low delay: opportunistic routing meets duty cycling", in Proceedings of the 11th international conference on Information Processing in Sensor Networks 2012 Apr 16 (pp. 185-196). ACM

[14] Barbato A, Barrano M, Capone A, Figiani N, "Resource oriented and energy efficient routing protocol for IPv6 wireless sensor networks", in Online Conference on Green Communications (GreenCom), 2013 IEEE 2013 Oct 29 (pp. 163-168). IEEE.

[15] Rezaei E, "Energy efficient RPL routing protocol in smart buildings".2014

[16] Korte KD, Sehgal A, Sch•onw•alder J, "A study of the RPL repair process using ContikiRPL", in IFIP International Conference on Autonomous Infrastructure, Management and Security 2012 Jun 4 (pp. 50-61). Springer, Berlin, Heidelberg

[17] Yang X, Guo J, Orlik P, Parsons K, Ishibashi K, "Stability metric based routing protocol for low-power and lossy networks", in Communications (ICC), 2014 IEEE International Conference on 2014 Jun 10 (pp. 3688-3693). IEEE

[18] Alvi SA, ul Hassan F, Mian AN, "On the energy efficiency and stability of RPL routing protocol", in Wireless Communications and Mobile Computing Conference (IWCMC), 2017 13th International 2017 Jun 26 (pp. 1927-1932). IEEE.

[19] Lamaazi H, Benamar N, "RPL enhancement using a new objective function based on combined metrics", in Wireless Communications and Mobile Computing Conference (IWCMC), 2017 13th International 2017 Jun 26 (pp. 1459-1464). IEEE.

[20] Gaddour O, Koubaa A, Chaudhry S, Tezeghdanti M, Chaari R, Abid M, "Simulation and performance evaluation of DAG construction with RPL", in Communications and Networking (ComNet), 2012 Third International Conference on 2012 Mar 29 (pp. 1-8). IEEE

[21] Liu S, Wu M, Chen C, Lv B, Li S, "A high-throughput routing metric for multi-hop Ad hoc networks based on real time test bed", in TENCON 2013-2013 IEEE Region 10 Conference (31194) 2013 Oct 22 (pp. 1-4). IEEE

[22] Gonizzi P, Monica R, Ferrari G, "Design and evaluation of a delay-efficient RPL routing metric", in Wireless Communications and Mobile Computing Conference (IWCMC), 2013 9th International 2013 July (pp. 1573-1577). IEEE

[23] Le Q, Ngo-Quynh T, Magedanz T, "RPL-based multipath routing protocols for internet of things on wireless sensor networks", in Advanced Technologies for Communications (ATC), 2014 International Conference on 2014 Oct 15 (pp. 424-429). IEEE.

[24] Aljarrah E, Yassein MB, Aljawarneh S, "Routing protocol of low-power and lossy network: Survey and open issues", in Engineering MIS (ICEMIS), International Conference on 2016 Sep 22 (pp. 1-6). IEEE.

[25] Cotrim JR, Kleinschmidt JH, "Performance evaluation of RPL on a mobile scenario with different ContikiMAC radio duty cycles", in High Performance Switching and Routing (HPSR), 2017 IEEE 18th International Conference on 2017 Jun 18 (pp. 1-6). IEEE.

[26] Banh M, Mac H, Nguyen N, Phung KH, Thanh NH, Steenhaut K, "Performance evaluation of multiple RPL routing tree instances for internet of things applications", in Advanced Technologies for Communications (ATC), 2015 International Conference on 2015 Oct 14 (pp. 206-211). IEEE.

[27] Iova O, Theoleyre F, Noel T, "Using multi-parent routing in RPL to increase the stability and the lifetime of the network", Ad Hoc Networks. 2015 Jun 1;29:45-62.

[28] Iova O, Theoleyre F, Noel T, "Stability and efficiency of RPL under realistic conditions in wireless sensor networks", in Personal Indoor and Mobile Radio Communications (PIMRC), 2013 IEEE 24th International Symposium on 2013 Sep 8 (pp. 2098-2102). IEEE

[29] Dunkels A, Gronvall B, Voigt T, "Contiki-a lightweight and flexible operating system for tiny networked sensors", in Local Computer Networks, 2004, 29th Annual IEEE International Conference on 2004 Nov 16 (pp. 455-462). IEEE

[30] Qasem M, Altawssi H, Yassein MB, Al-Dubai A, "Performance evaluation of RPL objective functions", in Computer and Information Technology; Ubiquitous Computing and Communications; Dependable, Autonomic and Secure Computing; Pervasive Intelligence and Computing (CIT/IUCC/DASC/PICOM), 2015 IEEE International Conference on 2015 Oct 26 (pp. 1606-1613). IEEE.

# Towards Effective Service Discovery using Feature Selection and Supervised Learning Algorithms

Heyam H. Al-Baity[1], Norah I. AlShowiman[2]

IT department, College of Computer and Information Sciences
King Saud University, Riyadh, Saudi Arabia

*Abstract*—With the rapid development of web service technologies, the number and variety of web services available on the internet are rapidly increasing. Currently, service registries support human classification, which has been observed to have certain limitations, such as poor query results with low precision and recall rates. With the huge amount of available web services, efficient web service discovery has become a challenging issue. Therefore, to support the effective application of web services, automatic web service classification is required. In recent years, many researchers have approached web service classification problems by applying machine learning methods to automatically classify web services. The ultimate goal of our work is to construct a classifier model that can accurately classify previously unseen web services into the proper categories. This paper presents an intensive investigation on the impact of incorporating feature selection methods (filter and wrapper) on the performance of four state-of-the-art machine learning classifiers. The purpose of employing feature selection is to find a subset of features that maximizes classification accuracy and improves the speed of traditional machine learning classifiers. The effectiveness of the proposed classification method has been evaluated through comprehensive experiments on real-world web service datasets. The results demonstrated that our approach outperforms other state-of-the-art methods.

*Keywords*—*Web service discovery; Web Service Description Language (WSDL); supervised machine learning; classification; feature selection*

## I. INTRODUCTION

The number of the available services that are published on the internet is increasing rapidly. These services are provided by different domains (e.g., education, finance, and health) and are available anywhere and anytime. Therefore, finding suitable web services for users quickly and efficiently has become a challenging issue and key problem to be solved. Web services are client and server applications that communicate over the World Wide Web [1, 2]. Basically, a web service works as a request-response form, where a client requests a service from a service provider through a request message. Upon receiving a request, a service provider will respond with a response message. The Web Service Description Language (WSDL), which uses an XML format, can be used to describe web services [3]. WSDL documents are stored in a centralized web service repository called Universal Description, Discovery, and Integration (UDDI). UDDI allows service providers to register their services and clients to discover these services.

Classifying web services into different categories according to their functionality is an efficient method of web service classification. This classification process is typically performed manually by domain experts. However, with the rapid growth of web services on the internet, it has become impractical to organize, classify, and manage web services manually as this requires intensive human effort. Additionally, it is an error-prone task due to the large number of categories in web service registries [4, 5]. Therefore, combining machine learning (ML) techniques to classify and manage web services automatically is an important task.

Based on the popularity of web services and the potential benefits that can be obtained from automatic web service classification, research in this field has recently gained significant attention. Several ML methods have been proposed to automatically classify web services [6]. However, it is still an open problem and further improvements can be achieved.

In this paper, we propose an enhanced method for the automatic classification for web services. Our approach is essentially based on the combination of feature selection methods and supervised ML classifiers. Specifically, feature selection methods have been used prior to performing the classification tasks. The main goal of incorporating feature selection is to select only the significant attributes from the dataset for the classifier. This reduces the size of the dataset and significantly improves the efficiency and accuracy of the classifier. To the best of our knowledge, no extensive work has been performed in this area.

The remainder of this paper is organized as follows. Section II describes the key concepts that are required to implement the proposed approach. In Section III, an overview of previous works related to solving the problem of classification of web services is presented. Section IV describes our methodology. Section V provides a detailed description of our approach and experimental results for real-world service description data. Finally, Section VI concludes this paper and discusses potential future work.

## II. BACKGROUND

In this section, some important concepts necessary for the rest of the paper will be presented. Firstly, we will focus on the main employed classifiers. Secondly, the feature selection strategy will be explained.

## A. Overview of Employed Classifiers

*1) Support Vector Machine (SVM):* SVM is one of the most popular ML algorithms for classification and regression analysis. It operates based on the concept of finding a hyperplane that maximizes the margin between two classes [7]. SVM learns from a training dataset, where each data sample is associated with a class label. It is effective for problems with large dimensionality, such as image and text categorization, because each data sample in the training set is represented as a point in an n-dimensional space, where n is the number of features. SVM then maps new data samples to the closest classes [8, 9]. In general, data samples will be represented in different categories when a large gap divides them. This gap is called a hyperplane. Many hyperplanes can separate a group of data, but the most optimal hyperplane is the hyperplane that creates the largest separation or maximum gap space between classes.

*2) Decision Tree (DT):* The DT is a recursive predictive classifier that predicts outcomes based on input data. It uses a tree structure to visualize a dataset, and represents sequences and consequences [10]. It can be represented using "IF, THEN" rules to be easily understood. The core concept of a DT is to repeatedly split a dataset into smaller datasets according to descriptive features until a sufficiently small dataset is obtained that contains data points with a single label. A DT has three types of nodes. First, the root node is the topmost node in the tree and has two or more branches. It is used to store predictors. Second, internal nodes or non-leaf nodes represent attributes of the root node. Finally, leaf nodes have no further branches and represent the outcomes of all prior decisions. A branch represents the outcome of a test, which is labeled in leaf nodes. The depth of a node is the minimum number of steps required to reach the node when starting from the root node. The path from the root node to a leaf node contains series decisions that can be converted into a decision rule. There are many types of DT algorithms, with the most popular being C4.5.

*3) Naive Bayes (NB):* NB is a probabilistic prediction classifier. It uses Bayes' theorem of probability to predict the probability of a tuple belonging to a specific class [11]. NB is useful for very large datasets and is easy to implement because it does not require prior knowledge of data and assumes that attributes are independent. NB begins by learning the conditional probabilities of any input attributes representing categorical data and outputs a class label with a corresponding probability score.

*4) Neural Network (NN):* NN algorithms have become very popular over the past few years. NNs are inspired by the architectural depth of the human brain. They can be used for regression or classification problems. The main purpose of an NN is to convert inputs into significant outputs by allowing the computer to behave like a human brain to solve problems [12]. Generally, NNs are organized in layers. An artificial NN refers to any structure of interconnected neurons. NNs have achieved excellent results for speech recognition, visual object recognition, object detection, natural language processing, etc.

## B. Feature Selection Strategy

Feature selection is the process of automatically selecting a subset of the most relevant features for a problem from an original feature set for use in a model. In this manner, irrelevant or redundant features can be removed without losing any important information. The goal of this technique is to increase the ability of a model by minimizing redundancy and maximizing relevant data. Furthermore, it decreases the required storage space and time for processing [13]. The basic steps of feature selection methods are subset generation, subset evaluation, stopping criterion, and result validation. Typically, each feature has a binary weight of one if it is selected and zero if it is not selected.

There are three general classes of feature selection algorithms: filter methods, wrapper methods and embedded methods.

*1) Filter methods:* Filter methods measure the general characteristics of training data, such as distance, consistency, dependency, information, and correlation, without using any specific classifiers. Filter methods apply a statistical measure to assign a score to each feature [14]. Such methods consist of two steps. In the first step, features are arranged based on two types of schemes: univariate or multivariate schemes. A univariate scheme rates each feature independently, whereas a multivariate scheme ranks features together in one step. Multivariate schemes are well suited to redundant features. In the second step, the features with the highest scores are used in the classification process. Examples of filter methods are relief, Fisher score, and information gain filters.

*2) Wrapper methods:* Wrapper methods use the predictive accuracy of a classifier itself to determine the quality of selected features. The general steps for wrapper models are to (1) select a subset of features, and (2) use the target classifier to evaluate the selected features [15]. Steps one and two are repeated until an acceptable performance level is reached. Examples of this wrapper method are sequential forward selection, sequential backward selection, and genetic algorithm wrappers.

The main difference between filter and wrapper methods is that filter methods act as a preprocessing step and work independently of the learning algorithm when selecting the features, whereas wrapper methods operate in the context of the learning algorithm. Typically, wrapper methods provide better predictive accuracy than filter methods. However, wrapper methods can be very slow because of the repeated calls to the learning algorithm.

*3) Embedded methods:* Embedded methods attempt to narrow the gap between filter and wrapper methods by combining the advantages of each type of method. Such methods interact with learning algorithms and perform feature selection during model construction without splitting the data

into training and testing sets. Examples of embedded methods are DT, random forest, NB, and SVM methods.

## III. RELATED WORK

Classifying web services automatically is considered to be one of the most important issues for web services. Many researchers have focused on this problem in their research. We will highlight some recent studies on web service classification that are based on ML approaches. The reviewed research papers will be categorized according to the type of ML algorithms used (supervised, unsupervised, and hybrid (supervised and unsupervised)).

### A. Supervised Approaches

Laachemi et al. [16] proposed an approach that uses an SVM and enhances classification accuracy based on a stochastic local search (SLS). As a preprocessing step, they scaled or resized the values of quality web service (QWS) dataset attributes [17] to use them in the SVM classifier to finding optimal solutions. The accuracy of this SLS-SVM approach was 84.86%, which is better than the accuracy of an SVM alone or other similar classifiers (NB, metaAdaBoostMl, lazyL Wang-Landau, lazyKstar and treesDS).

Liu et al. [18] combined an SVM classifier with latent Dirichlet allocation (LDA)-based topic models to classify large-scale services and scale down the cost of manually labeling services for training a classifier. The LDA algorithm was used to extract high-level topics from services. The SVM was used as the base classifier for this approach because it performs well on text classification. This is important for classifying web services based on their descriptions. Additionally, they introduced a pool-based active learning strategy to decrease the cost of manual labeling of services, which is required for building a training set. The dataset used for testing was the distributed reliability assessment mechanism for web services (WS-DREAM). In the preprocessing stage, they manually labeled services with informative descriptions based on their functionality. Next, they looked for categories with high numbers of services and selected ten categories. Therefore, as the number of the service increases, the LDA-SVM will provide more accurate results than an SVM alone. Experimental results clearly demonstrated the effectiveness of this active learning service classification framework.

Mustafa et al. [19] proposed a novel classification model using a multi-layer perceptron neural network (MLPNN) optimized via Tabu search (TS). The MLPNN is a model inspired by neuroscience that is used to predict events. They used MLPs with back propagation (BPP) and the Levenberg-Marquardt (LM) algorithm to train an MLP classifier and TS to optimize the classifier. They used a QWS dataset similar to the dataset used in [16]. Experimental results demonstrated that the MLP-TS model achieved superior accuracy, precision, recall, and root mean squared error (RMSE) compared to the un-optimized MLP-LM and MLP-BPP models.

Raj et al. [4] proposed a method to improve web service selection using the K-nearest neighbors (KNN) algorithm based on quality of service (QoS) parameters. They implemented this method on a large dataset containing 5,825 web services. The proposed classification approach begins by utilizing the KNN algorithm as a feature selection method to select a smaller, but more discriminative set of features based on QoS parameters. The classification process is then applied using the selected features. The results indicate that this approach can speed up the selection process compared to manually selected results.

Liu et al. [20] used a semantic web service classification method based on NB to enhance web service classification accuracy. They trained a classifier starting from the service interface level and then used an NB model to classify web services. They used the OWLS-TC dataset, which consists of 1,000 web services. To identify optimal classification features, they used the information gain metric, which led to improved service classification efficiency. The proposed model involves three stages: data preparation, classifier training, and application. The final accuracy of the classification was over 90% and the information gain values differed based on differences in the service attributes.

### B. Unsupervised Approaches

Zhao et al. [21], proposed a novel clustering method called the multiple attribute object NetClus (MAO-NetClus) for web service classification to improve the accuracy of service recommendation. This method is based on a heterogeneous information network that focuses on geographic information (i.e., the relationships between web services and their locations). This approach was tested on the WS-DREAM web service dataset. They evaluated the performance of the MAO-NetClus algorithm in different scenarios with different data sparseness, cluster quantity, and iteration numbers. Their approach overcame the limitations of typical web service clustering methods, which are generally based on web service description information, and achieved better performance than the original NetClus algorithm. Additionally, it improved the accuracy of service recommendation.

Tian et al. [22] proposed a web service clustering approach based on clustering web services with short textual descriptions. They used a tag-aided dual author topical model (TD-ATM) to enhance short text clustering by using tags to find long texts on Wikipedia. They used the programmable web dataset, which contains 11,339 web services. Feature extraction was incorporated as a final step during preprocessing to produce the inputs of for the TD-ATM. To evaluate the performance of the TD-ATM, they crawled 4,402 web services and 1,065 Wikipedia pages with a focus on the tags of the web services. They evaluated their approach using the common metrics of entropy, purity, and F-measure. The TD-ATM achieved superior results when compared to previous web service clustering approaches, such as the K-means and agglomerative models, and ATM Long (ATM-L) and ATM Short (ATMS).

### C. Hybrid Approaches

Rupasingha et al. [23] improved web service clustering by using ontology learning and a SVM. They used a SVM to calculate the semantic similarity in a generated ontology of web services instead of an edge-count-based method. This approach calculates similarity based on the summary of hybrid term similarity (HTS) and summary of context-aware similarity

(CAS) methods. To cluster web services, they used WSDL files. The proposed hybrid approach includes five phases. Phase one is a feature extraction process that describes the characteristics of each web service. Creating an ontology for each extracted feature is performed in phase two. Phase three calculates the web service similarity values of the ontologies using the SVM. The integration of five different features to calculate a final similarity value is performed in phase four. In the final phase, an agglomerative clustering algorithm (i.e., HTS) is used to cluster the web services. The results indicate decreased purity and increased entropy compared to other approaches with an increasing number of web services. The efficiency and accuracy of this hybrid approach combining HTS and CAS are better than those of each method individually.

Helmy et al. [24] proposed a novel approach that enhances web service clustering by using supervised machine learning techniques. They used DT, NB, and deep learning (DL) classification methods. The WSDL service retrieval test collection dataset was used in their study. This dataset contains 1,088 WSDL services classified into nine domains. The process for web service clustering uses different steps based on the RapidMiner software package. In the first step, which focuses on the preprocessing of data, tokenization based on non-letters and English stop word removal, are used to filter the WSDL files. Feature extraction is the second step of the classification process. In this step, information is extracted from each service's WSDL file. The third step is trimming the extracted features and building a feature-vector-space model, which is applied in the clustering step to produce outputs. To evaluate the effectiveness of this algorithm, the authors computed its precision, recall, accuracy, and F1-score. For the sake of comparison, they tested three techniques, namely the DT, NB, and DL methods, to determine which techniques provide the best results in terms of accuracy and efficiency. The DL technique resulted in the highest accuracy compared to the other approaches, but required more processing time compared to the NB and DT methods.

### D. Discussion of Related Work

The aim of the review of literature is to investigate and classify recent research on web service classification using ML approaches. Table I contains a summary of the main features of previous research papers, including information regarding the publication year, where or not feature extraction or selection was used, types of ML approaches, applied algorithms, datasets, evaluation metrics, and results of experiments. During our literature review, the following trends have been noticed:

- The field of web service classification is becoming increasingly important based on the increasing number of web services available on the internet.

- Recently, researchers have largely focused on ML algorithms, which can provide efficient web service discovery and overcome manual classification problems.

- The average number of web services in the test datasets was 3,738.

- Feature extraction mechanisms are used to optimize data based on both supervised and unsupervised approaches.

- The most popular classifiers in the literature are the SVM, DT, and NB model.

- The most commonly used evaluation metrics are accuracy, entropy, purity, precision, recall, and F-measure.

Based on previous work on supervised approaches, the algorithms that achieved the best accuracy for classifying the web services can be summarized as follows:

- An NB model achieved 90% accuracy in [20].

- Approaches that combined SVM classifier with LDA [3] or SLS [16] achieved better accuracy compared to SVM classifiers alone.

- Deploying metaheuristic optimization techniques (e.g., TS) with MLPNN classifier enhanced the accuracy of model results [19].

Therefore, based on the above research efforts, the aim of this work is to improve the classification accuracy for web services using supervised ML algorithms combined with feature selection methods. Such combination has the potential to achieve superior accuracy and has not yet been investigated in the literature. In the proposed model, two feature selection methods (filter and wrapper) are combined with four ML classifiers that achieved the highest accuracy for classifying web services in literature. The four classifiers are: SVM, NB model, DT (C4.5), and NN. Comprehensive experiments were conducted to test the performance of the classifiers with and without feature selection methods in order to show the effect of incorporating feature selection methods to the classification process and to find the best approach for web service classification.

TABLE I.        COMPARISON OF LITERATURE REVIEW PAPERS

| PaperRef. No. | Year | Feature Select./ Extract. | Machine Learning | | Evaluation | | |
|---|---|---|---|---|---|---|---|
| | | | *Approach* | *Algorithms* | *Dataset* | *Metrics* | *Results* |
| *[23]* | 2015 | Extraction | Hybrid | Ontology learning and SVM | Not specified | Purity, entropy, precision, recall, and F-measure | Purity decreased and entropy increased with an increasing number of web services. Average precision: 24.59%,4.69%, 9.16% Average recall: 29.04%, 2.04%, and 1.59% Average F-measure: 28.31%, 3.59%, and 5.42% |
| *[19]* | 2015 | ------ | Supervised | MMLP-TS | 364 Web services | Accuracy, precision, recall, RMSE | Accuracy: 97% |
| *[4]* | 2015 | KNN for Feature Selection | Supervised | KNN | 5,825 Web services | Accuracy | The approach speed-up the selection process compared to the manually selected results. |
| *[16]* | 2016 | ------ | Supervised | SVM and SLS | 364 Web services | Accuracy | Accuracy: 84.86% |
| *[3]* | 2016 | ------ | Supervised | SVM classifier with LDA | 3,738 Web services | Accuracy | LDA-SVM active learning yields more accurate results than SVM |
| *[20]* | 2016 | ------ | Supervised | NB | 1,000 Web services | Accuracy | Accuracy: >90% |
| *[21]* | 2016 | ------ | Unsupervised | NNovel web service clustering method: MAO-NetClus. | 3,738 Web services | Accuracy | MAO-NetClus yields better performance than original NetClus |
| *[22]* | 2016 | Extraction | Unsupervised | Tags | 11,339 Web service | Entropy, Purity, F-measure | Entropy: 19.6% , purity: 27.9% , and F-measure: 26%, |
| *[24]* | 2017 | Extraction | Hybrid | DT, NB, and DL | 1,088 WSDL services | Precision (p), recall (R), accuracy (A), F1-score (F) | DT: (p): 86.76, (R): 86.78, (A): 85.63, (F): 86.8 NB: (p): 90.5, (R): 90.4, (A): 90.34, (F): 90.4 DL: (p): 91.24, (R): 93.79, (A): 90.80, (F): 91.8 |

## IV. METHODOLOGY

In this study, Anaconda Python is used as our experimental platform. To implement the classification algorithms, the Spyder (3.2.3) compiler is used, which is an advanced interactive development environment for Python language. Spyder has built-in integration with a number of popular scientific packages that are used for different programming purposes, such as Matplotlib, NumPy, and scikit-learn.

The implementation process begins by importing the required packages. First, the NumPy library is imported, which is used to initialize the classification models. Second, the scikit-learn library is imported, which is the most widely used library for implementing ML algorithms.

For SVM classifier, the sklearn.svm.SVC class of scikit-learn is imported. Multiclass support was handled based on a one-versus-one scheme. The kernel parameter was set to a nonlinear kernel radial basis function. The remaining parameters were set empirically as follows: the regularization parameter C was set to 1 and the maximum number of iterations was set to 1000. The function gridSearchCV() is used to optimize the parameters.

For DT classifier, the DecisionTreeClassifier class has been imported to perform multi-class classification on the used dataset with a fully grown tree.

The DT employs a greedy algorithm that divides inputs via recursive binary splitting until reaching the leaf nodes. The maximum depth of the tree was set to 3 and the minimum number of samples required to create a leaf node was set to 5 for the initial values.

For NN classifier, MLP is used by importing the MLPClassifier class from the sklearn.neural_network library (sklearn.neural_network.MLPClassifier). This class trains a

model iteratively to prevent overfitting. The number of hidden layers was set to 10 and each layer contains 8 neurons. The maximum number of iterations was set to 200. For the NB classifier, the built-in NB algorithm (MultinomialNB) is used, which is suitable for classifying discrete features.

Following the split percentage validation technique, the used datasets have been divided into training set (60%) for constructing the classification models and the rest of the percentage data is used to test the developed models.

## V. EXPERIMENTS AND RESULTS

### A. Experiment Design

Each classification process for the proposed models consists of four general steps: processing of dataset, applying feature selection to the dataset, classifying web services based on the subset of features obtained from feature selection methods, and finally selecting the appropriate web service from the classification results.

### B. Datasets

In our experiments, a real-world web service dataset called QWS dataset [17] is used as our test collection for QoS prediction. This dataset contains many different web services and each web service is defined by the nine quality attributes presented in Table II. QWS attributes indicate the performance of the web service and determine which services satisfy a given set of user requirements. Web services in the QWS dataset are pre-classified into four categories: 1) platinum (high quality), 2) gold, 3) silver, and 4) bronze (low quality). The prediction of QoS is based on the quality rating provided by the web service relevancy function (WSRF), which ranks web service quality using the nine attributes. Based on its WSRF rank, each web service will have a service classification number in the range of 1–4. This number is useful for classifying web services into service categories. The QWS dataset covers many domains and uses a web service crawler engine to collect services.

The QWS[1] dataset is publicly available and has two free versions. The first version (Version 1.0) consists of 364 web services and was created in 2007. The updated QWS dataset (Version 2.0) has a set of 2,507 web services and QWS measurements taken in 2008 using a web service broker framework. In this study, both versions of the dataset are used. The first version is used for the sake of comparison to previous methods from [16] and [19], and the second version is used to test the performance of the proposed approach on a larger dataset. In our experiments, services were divided into two parts. One part was used as training data, and the other was used as testing data.

### C. Evaluation Metrics

To assess the performance of the four proposed classification models, the following commonly used performance metrics were calculated: accuracy, precision, sensitivity, specificity, error rate, and execution time. A brief description of these evaluation techniques is provided below [25].

- Accuracy, sometimes referred to as classification rate, is the total number of correct predictions over the total number of samples in the dataset. The maximum accuracy rate is 100%.

$$\text{Accuracy} = \frac{TP + TN}{P + N},$$

where TP=true positive, TN=true negative, P=positive, and N=negative.

- Error rate is the total number of incorrect predictions over the total number of samples in the dataset. The best error rate is zero and the worst is one.

$$\text{Error Rate} = \frac{FP + FN}{P + N}$$

- Sensitivity, sometimes referred to as recall or the true positive rate, is the total number of correct positive predictions over the total number of positive samples. The best sensitivity is one and the worst is zero.

$$\text{Sensitivity} = \frac{TP}{P}$$

- Specificity, sometimes referred to as the true negative rate, is the total number of correct negative predictions over the total number of negative samples. The best specificity is one and the worst is zero.

$$\text{Specificity} = \frac{TN}{N}$$

- Precision, sometimes referred to as the positive predictive value, is the total number of correct positive predictions over the total number of positive predictions. The best precision is one and the worst is zero.

$$\text{Precision} = \frac{TP}{TP + FP}$$

TABLE II.    QWS DATASET QUALITY ATTRIBUTES

| # | Attribute Name | Description | Units |
|---|---|---|---|
| 1 | Response Time | Time required sending a request and receiving a response. | ms |
| 2 | Availability | Number of successful invocations/total invocations. | % |
| 3 | Throughput | Total number of invocations for a given time period. | invocations/second |
| 4 | Successability | Number of responses/number of request messages. | % |
| 5 | Reliability | Ratio of the number of error messages to total messages. | % |
| 6 | Compliance | The extent to which a WSDL document follows WSDL specifications. | % |
| 7 | Best Practices | The extent to which a web service follows the WS-I basic profile. | % |
| 8 | Latency | Time required for the server to process a given request. | ms |
| 9 | Documentation | Measure of documentation (i.e., description tags) in WSDL document. | % |

---

[1]http://www.uoguelph.ca/~qmahmoud/qws/#Service_Classification_,

### D. Results Analysis and Comparisons

The proposed approach was validated in three phases. First, we compared the enhanced web service classification approach to two previous methods from [16] and [19] that used the QWS dataset (Version 1.0). Second, we tested the enhanced classification approach on the QWS dataset (Version 1.0) with and without using feature selection strategies. Third, we tested the enhanced classification approach on the QWS dataset (Version 2.0) with and without using feature selection strategies. The objective of evaluating the proposed classification models on different sizes of datasets was to study their scalability.

*1) QWS dataset (Version 1.0) Experiment 1:* For the sake of comparison and to clarify the impact of incorporating feature selection methods to the classification process, the proposed approach is compared to the method from [16], which used an NB classifier and SVM classifier combined with SLS. In addition, the proposed approach is compared to the method from [19], which used an NN classifier combined with TS.

As shown in Fig. 1 and Table III, the SVM classifier in our study (without incorporating feature selection) achieved a higher accuracy value (93%) compared to the accuracy value (84.86%) of the SVM classifier with SLS from [16]. Additionally, the SVM classifier with both feature selection methods (filter and wrapper) achieved better accuracy compared to the SVM classifier with SLS from [16].

The NB classifier in our study (without using feature selection) achieved the same accuracy value as the NB classifier from [16], which was 81%. However, when incorporating the feature selection methods (filter and wrapper) into NB classifier, the accuracy of the classifier slightly decreased compared to the classification accuracy of the NB classifier from [16].

The NN classifier in [19] was implemented using the MLP-TS. The accuracy of the MLP-TS model was 97% and that of our NN classifier (without using feature selection) was 92%. This value decreased to 84% when using feature selection (filter and wrapper), as shown in Table III. This indicates that the NN classifier performs better when using stochastic TS algorithms on small datasets. Fig. 2 presents the accuracy results of these two classifiers.



Fig. 2. Accuracy Comparison between our NN Classifier and Method from [19].

*2) Experiment 2*: The accuracy values of the four employed classifiers using QWS dataset (Version 0.1) are presented in Fig. 3. A detailed comparison between the four classifiers (with and without feature selection strategies) is provided in Table III. It is clear that the SVM classifier achieves the best accuracy results when using the wrapper feature selection method or no feature selection method. The NB classifier with the filter feature selection method achieved the worst accuracy value. NB is considered to be the fastest classifier and required only 0.01 second to complete the classification task with or without using any feature selection methods. The classifier with the lowest error rate (0.04) was the NN classifier with the wrapper feature selection method, meaning that the NN classifier predicted nearly all samples correctly.

Generally, it has been observed that deploying the wrapper method for the feature selection process enhanced the accuracy and error rate more than the filter method. However, the wrapper method has longer execution time.

It is important to note that the impact of using feature selection methods during the classification process did not clearly appear in the results of the above two experiments due to the small size of the QWS dataset used.



Fig. 1. Accuracy Comparison between our SVM and NB Classifiers (with and without Feature Selection) and the Method from [16].



Fig. 3. Accuracy Comparison between the Employed Classifiers (SVM, DT, NB, and NN) with and without Feature Selection Strategies on QWS Dataset (Version 1.0).

TABLE III.     EVALUATION RESULTS OF THE EMPLOYED CLASSIFIERS ON QWS DATASET (VERSION 1.0)

| Dataset | | Dataset Version 1.0 | | |
|---|---|---|---|---|
| *Classifier* | *Evaluation Metrics* | *Without Feature Selection* | *Filter Feature Selection* | *Wrapper Feature Selection* |
| *SVM* | Execution time in sec | 49.41 | 74.39 | 115.31 |
| | Accuracy | 0.93 | 0.85 | 0.93 |
| | Specificity | 0.95 | 0.89 | 0.95 |
| | Precision | 0.88 | 0.72 | 0.87 |
| | Error Rate | 0.065 | 0.15 | 0.06 |
| | Sensitivity | 0.87 | 0.72 | 0.88 |
| *DT* | Execution time in sec | 0.01 | 0.01 | 0.01 |
| | Accuracy | 0.88 | 0.83 | 0.86 |
| | Specificity | 0.91 | 0.88 | 0.91 |
| | Precision | 0.75 | 0.65 | 0.71 |
| | Error Rate | 0.13 | 0.16 | 0.13 |
| | Sensitivity | 0.73 | 0.72 | 0.72 |
| *NB* | Execution time in sec | 0.01 | 0.01 | 0.01 |
| | Accuracy | 0.81 | 0.73 | 0.80 |
| | Specificity | 0.86 | 0.82 | 0.86 |
| | Precision | 0.73 | 0.54 | 0.66 |
| | Error Rate | 0.18 | 0.26 | 0.19 |
| | Sensitivity | 0.63 | 0.53 | 0.65 |
| *NN* | Execution time in sec | 0.12 | 0.12 | 0.16 |
| | Accuracy | 0.92 | 0.84 | 0.84 |
| | Specificity | 0.94 | 0.89 | 0.89 |
| | Precision | 0.83 | 0.69 | 0.74 |
| | Error Rate | 0.07 | 0.15 | 0.15 |
| | Sensitivity | 0.86 | 0.68 | 0.70 |

*3) QWS dataset (Version 2.0) Experiment 3:* In this experiment, performance comparisons were performed between the employed four classifiers (SVM, NB, DT, and NN) on a larger version of QWS dataset (Version 2.0) with and without incorporating feature selection strategies. This new version of QWS was publicly available after conducting Experiments 1 and 2.

It is clear from Table IV that the SVM classifier with the wrapper method achieved the best accuracy value, but had a longest execution time compared to SVM without feature selection and the remaining classifiers. This indicates that the SVM classifier takes full advantage of the wrapper feature selection method to enhance its accuracy, but incurs significant computational costs. However, the error rates of SVM without feature selection and SVM with the wrapper method were the same.

The DT classifier (with and without feature selection methods) has faster execution times compared to SVM, but achieved lower classification accuracy and a higher error rate. The execution time of NB classifier was faster than that of the DT classifier when incorporating the two feature selection methods. However, it achieved worse accuracy and error rate values.

The NN classifier had the fastest execution time among all classifiers. The accuracy values of the NN classifier with the wrapper method and without feature selection were nearly the same. However, the accuracy deteriorated when using the filter method. In contrast, the NN classifier without feature selection showed a higher error rate compared to the wrapper method. This implies that unlike the filter method, the use of the wrapper method with the NN classifier does not adversely affect the accuracy and error rate of the classifier.

Although the SVM classifier achieved superior performance compared to the other classifiers in this experiment, the NN classifier achieved competitive results with a significant decrease in execution time compared to SVM.

As shown in Fig. 4, the classifier with the highest accuracy result was SVM with the wrapper method, followed by NN classifier. The classifiers with the filter method generally achieved worse results compared to their performance with the wrapper method.

In conclusion, after evaluating the classifiers, it has been found that the fastest classifiers when using the QWS dataset (Version 1.0) were the NB and DT classifiers. In contrast, for the QWS dataset (Version 2.0), the NN was the fastest. The classifiers that achieved the best classification accuracy and lowest error rates was the SVM followed by the NN on both versions of the QWS dataset.

The results revealed that the NN classifier achieved slightly better results on the larger QWS dataset (Version 2.0). Therefore, it can conclude that the NN classifier achieves better performance when considering larger datasets. However, incorporating feature selection with the NN classifier did not significantly enhance the accuracy of the model due to the nature of the NN classifier, which provides automatic prediction of hidden features.



Fig. 4.    Accuracy Comparison between the Employed Classifiers (SVM, DT, NB, and NN) with and without Feature Selection Strategies on QWS Dataset (Version 2.0).

TABLE IV. EVALUATION RESULTS FOR THE PROPOSED CLASSIFIERS ON THE QWS DATASET (VERSION 2.0)

| Dataset | | Dataset version 2.0 | | |
|---|---|---|---|---|
| *Classifier* | *Evaluation Metrics* | *Without Feature Selection* | *Filter Feature Selection* | *Wrapper Feature Selection* |
| *SVM* | Execution time in sec | 200.41 | 243.39 | 287.31 |
| | Accuracy | 0.91 | 0.85 | 0.93 |
| | Specificity | 0.94 | 0.89 | 0.95 |
| | Precision | 0.86 | 0.72 | 0.87 |
| | Error Rate | 0.064 | 0.15 | 0.06 |
| | Sensitivity | 0.85 | 0.72 | 0.88 |
| *DT* | Execution time in sec | 0.56 | 0.76 | 0.88 |
| | Accuracy | 0.89 | 0.85 | 0.88 |
| | Specificity | 0.98 | 0.90 | 0.91 |
| | Precision | 0.72 | 0.79 | 0.71 |
| | Error Rate | 0.11 | 0.19 | 0.13 |
| | Sensitivity | 0.70 | 0.72 | 0.72 |
| *NB* | Execution time in sec | 0.56 | 0.64 | 0.77 |
| | Accuracy | 0.81 | 0.70 | 0.82 |
| | Specificity | 0.86 | 0.77 | 0.88 |
| | Precision | 0.73 | 0.50 | 0.69 |
| | Error Rate | 0.18 | 0.27 | 0.22 |
| | Sensitivity | 0.63 | 0.50 | 0.69 |
| *NN* | Execution time in sec | 0.43 | 0.43 | 0.55 |
| | Accuracy | 0.92 | 0.86 | 0.91 |
| | Specificity | 0.94 | 0.58 | 0.93 |
| | Precision | 0.83 | 0.73 | 0.82 |
| | Error Rate | 0.07 | 0.11 | 0.04 |
| | Sensitivity | 0.86 | 0.68 | 0.80 |

Furthermore, when combining the wrapper method with the four classifiers, the overall accuracy of the classifiers enhanced. Generally speaking, the wrapper method yields better accuracy than the filter method because it uses a preselected learning algorithm to evaluate and select an optimal subset of the features, which results in the best classifier performance. However, wrapper methods suffer from being computationally expensive. The SVM classifier achieved the highest accuracy compared to the other classifiers when using the wrapper method on both versions of the QWS dataset. The NN classifier came in second place when using the wrapper method on the QWS dataset (Version 2.0).

Moreover, it has been noticed that deploying the filter method with all four classifiers on both versions of the QWS dataset reduced their performance. The reason of this lower-than-expected performance is the size and dimensionality of the dataset. Filter methods work better on large, high-dimensional datasets

Finally, to enhance classifier performance by using feature selection strategies (filter and wrapper methods), several important factors must be considered, including dataset size, dimensionality of the dataset, the nature of the classifier, and nature of the classification/prediction problem. For example, filter methods are better suited to problems that must be solved online or in batches because they are faster compared to wrapper methods.

## VI. CONCLUSIONS AND FUTURE WORK

Effective web service classification is a crucial issue for web services. In this study, we focused on the problem of supervised classification of web services. A novel automated classification method that combines state-of-the-art ML classifiers (SVM, DT, NB, and NN) with feature selection methods (filter and wrapper) is proposed. The purpose of employing feature selection methods in the classification process is to find the effective feature subset prior to training and testing the classifier. We expected the classification accuracy of the ML classifiers to improve when using these methods. Intensive experiments to evaluate the proposed approach were conducted on a publicly available real-world dataset for web services called the QWS dataset and comparisons to related methods were made. Preliminary results revealed a slight improvement in classification accuracy. Additionally, the proposed approach outperformed other algorithms discussed in the literature.

Our future work will proceed in two different directions. First, we will extend our experimental study to a larger dataset with higher dimensionality to investigate the impact of these factors on the performance of the proposed approach. Second, we will study the effects of employing DL algorithms for web service classification problems.

## REFERENCES

[1] A. S. Mustafa and Y. S. Kumaraswamy, "Data mining algorithms for Web-services classification," in International Conference on Contemporary Computing and Informatics (IC3I), 2014, pp. 951-956. IEEE.

[2] P. Kaur, "Web Content Classification: A Survey," International Journal of Computer Trends and Technology (IJCTT), vol. 10, no. 2, 2014.

[3] X. Liu, S. Agarwal, C. Ding, and Q. Yu, "An LDA-SVM Active Learning Framework for Web Service Classification," IEEE International Conference on Web Services (ICWS), San Francisco, CA, 2016, pp. 49–56.

[4] M. Raj and S. Pragasam, "QoS based classification using K-Nearest Neighbor algorithm for effective Web service selection," IEEE International Conference on Electrical, Computer and Communication Technologies (ICECCT), Coimbatore, 2015, pp. 1–4.

[5] F. Deng, "Web Service Matching based on Semantic Classification," M.S. thesis, School of Health and Society, Department of Computer Science, pp. 9–20, 2012.

[6] M. Cracknell, "Machine learning for geological mapping: algorithms and applications," Ph.D. dissertation. University of Tasmania, Australia, 2015.

[7] M. Praveena and V. Jaiganesh, "A literature review on supervised machine learning algorithms and boosting process," International Journal of Computer Applications, vol.169, no. 8, pp. 32-35, 2017.

[8] S. B. Kotsiantis, I. Zaharakis, and P. Pintelas, "Supervised machine learning: A review of classification techniques," Emerging artificial intelligence applications in computer engineering, pp. 3–24, 2007.

[9] J. Yang and X. Zhou, "Semi-automatic Web Service Classification Using Machine Learning," International Journal of u-and e-Service, Science and Technology, vol. 8, no. 4, pp. 339-348, 2015.

[10] J. Fürnkranz, D. Gamberger, and N. Lavrač, "Foundations of rule learning," Springer Science & Business Media, 2012.

[11] T. Raj, T. F. Michael, K. Ravichandran, and K. Rajesh, "Domain specific Web service composition by parameter classification using naïve Bayes algorithm," World Applied Sciences Journal, vol. 29, pp. 99–105, 2014.

[12] J. Schmidhuber, "Deep learning in neural networks: An overview," Neural networks, vol. 61, pp. 85–117, 2015.

[13] J. Gui, Z. Sun, S. Ji, D. Tao, and T. Tan, "Feature selection based on structured sparsity: A comprehensive study," IEEE transactions on neural networks and learning systems, vol. 28, no. 7, pp. 1490-507, 2016.

[14] M. Mwadulo, "A review on feature selection methods for classification tasks," International Journal of Computer Applications Technology and Research, vol. 5, no. 6, pp. 395-402, 2016.

[15] Y. Dhote, S. Agrawal, and A. J. Deen, "A Survey on Feature Selection Techniques for Internet Traffic Classification," in International Conference on Computational Intelligence and Communication Networks (CICN), Jabalpur, 2015, pp. 1375–1380.

[16] A. Laachemi and D. Boughaci, "A stochastic local search combined with support vector machine for Web services classification,"

[17] E. Al-Masri and Q. H. Mahmoud, "Discovering the best Web service: A neural network-based solution," in IEEE International Conference on Systems, Man and Cybernetics, 2009, pp. 4250–4255.

[18] X. Liu, S. Agarwal, C. Ding, and Q. Yu, "An LDA-SVM Active Learning Framework for Web Service Classification," in IEEE International Conference on Web Services (ICWS), San Francisco, CA, 2016, pp. 49–56.

[19] A. Syed Mustafa and Y. Kumara Swamy, "Web Service classification using Multi-Layer Perceptron optimized with Tabu search," in IEEE International Advance Computing Conference (IACC), Banglore, 2015, pp. 290–294.

[20] J. Liu, Z. Tian, P. Liu, J. Jiang, and Z. Li, "An Approach of Semantic Web Service Classification Based on Naive Bayes," in IEEE International Conference on Services Computing (SCC), San Francisco, CA, 2016, pp. 356–362.

[21] H. Zhao, J. Wen, J. Zhao, and F. Luo, "A new model-based Web service clustering algorithm," in IEEE Region 10 Conference (TENCON), Singapore, 2016, pp. 3468–3472.

[22] G. Tian, J. Wang, K. He, and C. Sun, "Leveraging Auxiliary Knowledge for Web Service Clustering," Chinese Journal of Electronics, vol. 25, no. 5, pp. 858–865, 2016.

[23] M. Rupasingha, I. Paik, and B. T. G. S. Kumara, "Calculating Web service similarity using ontology learning with machine learning," in IEEE International Conference on Computational Intelligence and Computing Research (ICCIC), Madurai, 2015, pp. 1–8.

[24] A. Helmy and M. Geith, "An Enhanced Approach for Web Services Clustering using Supervised Machine Learning Techniques," International Journal of Scientific & Engineering Research, vol. 8, no. 1, pp. 158–170, 2017.

[25] M. Sokolova and G. Lapalme, "A systematic analysis of performance measures for classification tasks," Information Processing & Management, vol. 45, no. 4, pp. 427-437, 2009

in International Conference on Advanced Aspects of Software Engineering (ICAASE), Constantine, 2016, pp. 9–16.

# An Advanced Emergency Warning Message Scheme based on Vehicles Speed and Traffic Densities

Mustafa Banikhalaf[1], Saleh Ali Alomari[2], Mowafaq Salem Alzboon[3]

Faculty of Information Technology and Computer science, Yarmouk University, Jordan, Irbid[1]
Faculty of Science and Information Technology, Jadara University, Jordan, Irbid[2],[3]

*Abstract*—**In intelligent transportation systems, broadcasting Warning Messages (WMs) by Vehicular Ad hoc Networks (VANETs) communication is a significant task. Designing efficient dissemination schemes for fast and reliable delivery of WMs is still an open research question. In this paper, we propose a novel messaging scheme, Advanced Speed and Density Warning Message (ASDWM). ASDWM is a broadcast-based scheme that meets design objectives and achieves high saved rebroadcast and reachability, as well as low end-to-end latency of WM delivery. The ASDWM uses vehicle speeds and vehicles density degrees to help emergency vehicles to send WM according to a road condition, adaptively. Simulation results demonstrate the effectiveness and superiority of the ASDWM over its counterparts.**

*Keywords—Warning message; the broadcast storm problem; emergency vehicles*

## I. INTRODUCTION

Inside a city or on the highway, unexpected events such as traffic accidents and medical emergencies occur every day. Hence, it is very critical for emergency vehicles to reach an accident spot as soon as possible. In case of such unexpected events, an emergency vehicle should inform other vehicles and traffic lights ahead to clear the way for it. Basically, all emergency vehicles are equipped with a wireless card to detect the event and utilize the underlying VANET architecture to issue WMs. One of the most important challenges in this scenario is to design a warning message dissemination scheme to transfer the WM in a reliable and low-latency broadcasting. It is crucial that all vehicles and traffic lights in front of emergency vehicles receive the WM with high probability, and with the minimum possible delay to take a proper action very quickly.

A simple solution to the above problem is handled using blind flooding [1]. It is a primitive scheme that allows each vehicle to rebroadcast the WM when it receives it for the first time to all surrounding vehicles exactly once. Blind flooding trends to be an optimal solution in low sparse networks. However, using it in cases of high vehicular traffic densities cause several communication channel problems. Network performance drops down due to a considerable number of duplicated generated messages. These duplicated messages can cause message collisions, increase occupancy or contention on available channel capacity, which leads to high latency in delivering WM. This problem is widely known as the broadcast storm problem [1], and several solutions have been proposed in the literature to mitigate its effects [2][3][4]. The

main idea of these solutions has been to limit rebroadcasting the message to candidate vehicles and to guarantee that all other vehicles received it. Each proposed solution uses a different method to choose group of vehicles that are responsible to forward WM with minimum overhead, latency and high reachability.

In VANETs, the most effective and reliable broadcast approach is to privilege a broadcast operation to any vehicle that located on the transmission range of the source vehicle [5][6]. However, this approach cannot be operated without Global Positioning System (GPS) availability. It is commonly known that the strength of a GPS signal is often influenced by external environment conditions, which imposes several issues when using this approach. Other proposed approaches use a WM prioritization technique [7] to make GPS based schemes more efficient. It is assumed that the WM should be assigned a high priority at MAC layer and broadcasted before other types of messages. Although this technique helps some safety applications to guarantee that the WM is delivered as fast as possible, but time delay to broadcast two consecutive WMs depends on a fixed slot time. This increases number of unnecessary WMs, which reduces network communication efficiency. The main contribution in this work is to: (i) develop (ASDWM) scheme, (ii) and to adjust the required time delay before broadcasting the next WM dynamically, (iii) which eventually mitigates side effect of the broadcast storm problem.

The rest of the paper is organized as follows. Section II reviews the related work on the probabilistic broadcasting schemes. Section III presents ASDWM. Section IV shows the simulation environment used to validate the ASDWM, presents and discusses the obtained results. Finally, Section V concludes this paper.

## II. RELATED WORKS

Over the last decades, probabilistic broadcasting schemes have been implemented using several techniques such as probabilistic density-based schemes [3]-[8], probabilistic counter-based schemes [9][10] and [11], probabilistic GPS-based schemes [12][13][14], and probabilistic algebraic-based scheme [15]. Several studies show that probabilistic schemes that depends on a GPS device are the appropriate solutions for several communication scenarios in VANETs. However, GPS availability and reliability are not always guaranteed, which makes these studies not always applicable. In this paper, we shed the light on some important probabilistic schemes which require GPS availability, followed by the most effective and relevant studies that depend on the vehicles speed as a

parameter to calculate the rebroadcast probability [16] and [17].

### A. *Probabilistic GPS based-Schemes*

Irresponsible Forwarding (IF) [12], is a broadcasting scheme presented for VANETs which combines advantages of both distance and density schemes to calculate a dynamic forwarding probability value. The inter-vehicle spacing distribution in the network, and the distance between two vehicles are used to calculate the forwarding probability using the following expression:

$$p = (1 - F_x(z - d))^{1/k} \qquad (1)$$

where X denotes the distance between consecutive vehicles, d is a distance between a transmitter and a receiver, and z is the transmission range, $F_x(z - d))$ is the Cumulative Distributed Function (CDF), and k is a shaping parameter used to adjust the rebroadcast probability. Authors in [18] have proposed CAREFOR scheme based on IF. It assumes that is not practical for each vehicle in the network has the same transmission range. Hence, CAREFOR includes ratio accounts for difference in the vehicle's transmission range.

Weighted P-Persistence (WP-P) [19] is another well-known scheme that uses the distance value between vehicles as an input to (2), to determine the forwarding probability, where the rebroadcast operation is privileged to farthest vehicles always:

$$P_{ij} = \frac{Dist_{ij}}{T_r} \qquad (2)$$

where $Dist_{ij}$ is the distance between vehicle i and vehicle j, and $T_r$ is the transmission radio range. Several similar GPS based schemes have been proposed in [20] and [21], which use vehicles coordination positions to calculate the rebroadcast value carefully. Interested readers can refer to these studies to gain deep knowledge and understanding.

### B. *Probabilistic Speed-based Schemes*

Speed Adaptive Broadcast (SAB) is proposed in [16] to adaptively adjust the forwarding probability based on the speed of traveling vehicle. A ratio between the speed value of current vehicle i at time t, and the limit of speed allowed $V_x$ for road vehicles, are used by the following equation to calculate the forwarding probability:

$$P(i, t) = \frac{V(i,t)}{V_x} \qquad (3)$$

In [17], the authors proposed Speed Adaptive Probabilistic Flooding (SAPF) to estimate traffic densities in VANETs based on the vehicle's speed, and to adjust the forwarding probability according to the following equation:

$$P = 0.055v - 0.033 \qquad (4)$$

where v is the vehicle's speed. SAPF defines two types of speed thresholds are vl and vh which represents vehicles that move with low speed, and vehicles that move with high speed respectively. Estimating the vehicle's density when v > vh is impossible. Hence, the direct blind flooding is used instead to guarantee high message reachability. On the other hand, if v < vl, a fixed value of the forwarding probability is used as the maximum network capacity has been reached.

### III. ASDWM DESCRIPTION

The significant feature of the ASDWM is to regulate the rebroadcast probability according to traffic densities and the vehicle speed. The main aim is to guarantee that WM is reached all vehicles with low latency for a period of time. Vehicles move with low speeds inside a city indicate high vehicle density due to traffic jam, accidents or other potential hazards. In this case, it is sufficient to use a low forwarding probability value to disseminate the WM to all vehicles with minimum cost. In the other hand, when vehicle speeds inside a city reach the maximum speeds this implies low vehicular densities, and high rebroadcast probability values are required to achieve a high percentage of the WM delivery. Therefore, the ASDWM categorizes traffic densities based on the vehicle speed into three regimes; low-density, medium-density and high-density. The low-density regime means that the speed of vehicle reaches to the maximum speed and the rebroadcast probability is set to be equal to 1 (i.e., flooding). We use the flooding technique in this regime to guarantee message delivery to all vehicles. Normally, inside the medium-density regime a vehicle travels at a speed above 10km/h and below the maximum speed. Hence, the rebroadcast probability is set to be equal 0.7 to balance between high message reachability and low retransmission [3]. In the high-density regime, when the vehicle moves with a speed less than 10km/h, it is not practical to set a fixed value for the rebroadcast probability. Therefore, the following strategy is used to estimate the accurate density level and adjust the rebroadcast probability accordingly. Three density levels (density level 1, density level 2 and density level 3) are calculated based on neighboring information to calculate the best of value the rebroadcast probability.

**Density Level 1 (DL1)** represents a density measurement of a set of 1-hop neighborhood information that can be covered via a set of 2-hop neighborhood information. Assume that $S_v^1$ is a set of 1-hop neighbors of a vehicle v, and $S_v^2$ is a set of 2-hop neighbors of v. N denotes the number of elements in each set. DL1 is calculated as follows:

$$DL1 = \frac{N(S_v^1)}{N(S_v^1) + N(S_v^2)} \qquad (5)$$

**Density Level 2 (DL2)** represents a density measurement of a set of 2-hop neighborhood information within 2-hop neighborhood information. DL2 is calculated as follows:

$$DL = \frac{N(S_v^2)}{N(S_v^1) + N(S_v^2)} \qquad (6)$$

**Density Level 3 (DL3)** represents a density measurement of a set of 2-hop neighborhood information that can only be reached via a 1-hop neighbor. Assume that $S_{v,x_k}^3$ is the set of 2-hop neighborhood information that can only be reached via 1-hop neighborhood information $x_k$ from a given vehicle v (for k = 1, 2, 3,….,n). DL3 is calculated as follows:

$$DL3 = \frac{1}{N(S_v^1)} \sum_{k=1}^{N(S_v^1)} N(S_{v,x_k}^3) \qquad (7)$$

Based on the collected information from DL1, DL2 and DL3, the rebroadcast probability value that to be used in the high-density regime is calculated as follows:

$$P = \frac{DL1+DL2+DL3}{3} \qquad (8)$$

The WM must be queued with the high priority at MAC layer. Enhanced Distributed Channel Access (EDCA) is the fundamental channel access mechanism of Wireless Access for Vehicular Environment (WAVE). It offers four Access Classes (ACs) to priorities messages based on their types and applications. AC0 indicates that a broadcast message has the lowest priority, and AC3 it has the highest priority [22]. The WM in this work is assigned AC3 to make sure it is delivered before any other types of messages.

The most related studies [7] and [23] use fixed delay time between broadcasting two consecutive WMs. In fact, this implementation is not practical in real life scenarios, and makes the broadcast storm problem worse. the ASDWM adopts the same density principle to adjust time delay when broadcasting WMs. The following equation is used:

$$\text{Time\_Delay} = \left(\frac{DL1+DL2+DL3}{3}\right) \times \Delta(t) \, \text{max} \qquad (9)$$

where $\Delta(t)max$ is a small random waiting time. Using the delay time (9) can regulate the period between consecutive broadcasting message adaptively. For instance, when a vehicle moves inside the high-density regime, delay time is preferable to be long as one broadcast is enough to reach all vehicles. While in the low-density regime delay time should be too short to keep sending the WM until arrives a destination area. The ASDWM steps are organized in Fig. 1.

---

**Sender: Emergency Vehicle (EV)**
**STEP1: WHEN** (a vehicle sends WM) **Do**
**STEP2:** Set WM = AC3
**STEP3:** Set( $Time\_Delay_i$ )
**STEP4:** Broadcast (WM)
**STEP5:** Set next ( $Time\_Delay_{i+1}$ )
**STEP6: WHILE** ( $Time\_Delay_i$ ) is not expired **THEN**
**STEP7:** {Broadcast (WM)}
**Receiver: Normal Vehicle (NV)**
**STEP1:** WHEN (a vehicle receives WM for the first time) **DO** {
**STEP2:**     **IF** (vehicle speed = MAX\_SPEED) **THEN**
**STEP3:**        {Broadcast WM with P = 1}
**STEP4: ELSE\_IF** (vehicle speed>10km/h and<MAX\_SPEED) **THEN**
**STEP5:**        {Broadcast WM with P = 0.7}
**STEP6:**        **ELSE** **IF** (vehicle speed<10 km/h) **THEN**
**STEP7:**        {Broadcast WM with P =(DL1+DL2+DL3)/3}

---

Fig 1.     The ASDWM Scheme Logical Steps.

*C. Example on ASDWM*

In Fig. 2 scenario (a), the emergency vehicle travels with 60km/h, which is the maximum speed limit on the road. Normally, when the driver notices that the road is empty, he accelerates the vehicle speed until it reaches the maximum speed limit. For this reason, the emergency vehicle keeps sending the WM every short regular period based on the density input to (9). As DL1, DL2 and DL3 is equal to 0. In Fig. 2 scenario (a), Time\_Delay is set to $10^{-3}$ms. Once the vehicle A receives the WM from the emergency vehicle, it retransmits it with probability p = 1, as no vehicle in the neighborhood rebroadcasts to prevent the WM's dying out.



Fig 2.     Adjusting the rebroadcast probability value according to the traffic densities.

In the medium-density regime as in Fig. 2 scenario (b), small probability value may lead to poor message reachability, and high probability value leads to the broadcast storm problem.

Hence, the rebroadcasting probability should be set to at least 0.7 to balance between message reachability and low retransmission [3]. Time\_delay between consecutives WMs is also calculated based on (9).

The serious side effect of the broadcast storm problem often appears in the high-density regime as shown in Fig. 2 scenario (c). Hence, the rebroadcasting probability must be chosen carefully to mitigate its side effect as much as possible. Usually, vehicles are considered inside the high-density regime if it travels with speed less than 10km/h. In this case, the rebroadcasting probability is set dynamically with respect to different density levels that are calculated in (8). For instance, when the vehicle A receives the WM from the emergency vehicle, it calculates DL1, DL2, DL3, P and Time\_Delay as follows:

$$DL1 = \frac{N(S_v^1)}{N(S_v^1) + N(S_v^2)} = \frac{3}{3+2} = 0.6$$

$$DL2 = \frac{N(S_v^2)}{N(S_v^1) + N(S_v^2)} = \frac{2}{3+2} = \frac{2}{5} = 0.4$$

$$DL3 = \frac{1}{N(S_v^1)} \sum_{k=1}^{N(S_v^1)} N(S_{v,x_k}^3) = \frac{1}{3} \sum_{k=1}^{3} 1 + 1 + 0 = \frac{2}{3} = 0.66$$

$$P = \frac{DL1 + DL2 + DL3}{3} = \frac{0.6 + 0.4 + 0.66}{3} = 0.55$$

$$Time\_Delay = \left(\frac{DL1 + DL2 + DL3}{3}\right) \times \Delta(t)max$$
$$= \left(\frac{0.6 + 0.4 + 0.66}{3}\right) \times 10^{-3}$$
$$= 0.55 \times 10^{-3}$$

where $N(S_v^1)$ and $N(S_v^2)$ are set of 1-hop neighbors [ B, E, F] = [3], and set of 2-hop neighbors [D, C] = [2] of the vehicle A, respectively. Set of 2-hops neighbors that can only be reached via 1-hop neighbor $N(S_{v,x_k}^3)$ for the vehicle A is collected as follows: $S_{A,B}^3 = [C] = [1]$, $S_{A,E}^3 = [D] = [1]$ and $S_{A,F}^3 = [0]$.

## IV. SIMULATION SETUP AND RESULT ANALYSIS

The performance of ASDWM is evaluated through simulation using NS-2 simulation environment [24], and it is compared against the most related scheme SAPF [17]. Traffic flows are generated using SUMO [25]. The physical layer frequency is adjusted to 5.9 GHz according to DSRC standard, and Bandwidth is set to 10MHz [26]. The transmission range communicating value is adjust to 250 meters. Between 100 and 500 vehicles are uniformly distributed on a road consisting of two lanes 5 Km length that is similar to Fig. 2. Vehicles travel at different speed between minimum 0km/h up to max 60km/h throughout the lanes. The following metrices are used to observe the network performance:

- **Reachability:** It is measured by the percentage of vehicles receiving the WM, divided by the total number of vehicles that are reachable.

- **Saved Rebroadcast:** It shows the ratio between the numbers of vehicles receiving the WM and the number of vehicles rebroadcasting it.

- **Latency:** It is the time between sending the WM from the source, until the time it reached the destination.

Fig. 3 illustrates the reachability achieved by Flooding, SAPF and ASDWM when the number of vehicles is varied. Basically, reachability always increases with increased number of vehicles. This is because when the density of vehicle increases, the distance between neighbors decreased and the number of vehicles covering a road segment increase. For instance, reachability achieved by Flooding increases from 70% for 50 vehicles to 100% for 200 vehicles, while the reachability achieved by SAPF and ASDWM increases from 30% to 100% for 30 and 200 vehicles, respectively. As expected, Flooding achieves the best reachability performance

compared to the other schemes, as it allows all vehicle to rebroadcast the WM. On the other hand, the reachability performance results in the dense network for SAPF and ASDWM are similar and comparable to Flooding.

Fig. 4 shows that the latency incurred by Flooding, SAPF and ASDWM increases with increased network density. All the schemes have similar latency when number of vehicles is equal 100. After this point, Flooding incurs the highest latency compared to other schemes. It is noticed from the Fig. 4 that ASDWM incurred the lowest latency compared to the Flooding and SAPF in the sparse and the dense network. This is because the ASDWM adjusts the Time_Delay between WMs according to traffic densities, which reduces the number of retransmissions, prevents contentions and message collisions in the vehicular networks. Fig. 5 shows the saved rebroadcast reported by Flooding, SAPF and ASDWM as network density increased. Normally, the vehicular network becomes denser as the number of vehicles increased, which makes all the schemes to rebroadcast a larger number of unnecessary WMs. it shows that the saved rebroadcast for Flooding scheme is equal to zero as all vehicles are eligible to retransmit the received WM. The figure also shows that as the network densities increased, the ASDWM achieves significant percentage of the saved rebroadcast compared to the SAPF. For instance, compared with the SAPF, the saved rebroadcast in the ASDWM can increase further by approximately 20% when the number of vehicles is relatively large (e.g., 500 vehicles).



Fig 3. Reachability (%) vs Number of Vehicles.



Fig 4. Latency (ms) vs Number of Vehicles.

Fig 5.    Saved Rebroadcasts vs Number of Vehicles.

## V.    CONCLUSION AND FUTURE WORKS

In this paper, we presented ASDWM, which is designed to disseminate WMs over VANETs in urban environments. The WMs can be quickly and reliably delivered by the emergency vehicles to the target destination with lowest latency and highest saved rebroadcast, without scarifying reachability. Simulation results show that in comparison with Flooding and SAPF protocols, ASDWM achieves the minimal message delivery latency, and keeps reachability is equivalent to Flooding in high densities area. Moreover, ASDWM can maintain more than 30% and 100% saved rebroadcast compared to SAPF and Flooding, respectively. For the future works, ASDWM can be extended further to include vehicle's direction when broadcasting WMs. This makes it more applicable with real life scenarios and VANETs applications.

### REFERENCES

[1]  Y.-C. Tseng, S.-Y. Ni, Y.-S. Chen and J.-P. Sheu," The broadcast storm problem in a mobile ad hoc network," Wireless Networks, vol.8, no.2-3, pp.153-167, Mar.2002.

[2]  M. B.Yassein, M. B. Khalaf and A. Y. Al-Dubai," A new probabilistic broadcasting scheme for mobile ad hoc on-demand distance vector ( AODV ) routed" The Journal of Supercomputing,vol.53, no.1,pp.196–211, 2010.

[3]  Z. J. Haas, J. Y. Halpern and L. Li, "Gossip-based ad hoc routing," IEEE/ACM Transactions on Networking, vol. 14, no. 3, pp. 479-491, June.2006.

[4]  J. Cartigny and D. Simplot," Border node retransmission based probabilistic broadcast protocols in ad-Hoc networks," Telecommunication Systems, vol. 22, pp.189–204 , 2003.

[5]  D. G. Reina, S. L. Toral, P. Johnson and F. Barrero "Improving discovery phase of reactive ad hoc routing protocols using Jaccard distance," The Journal of Supercomputing, vol.67, pp 131–152, January 2014.

[6]  G. Korkmaz, E. Ekici and F. Ozguner, "An efficient fully ad-hoc multi-hop broadcast protocol for inter-vehicular communication systems," IEEE International Conference on Communications (ICC 2006) , Istanbul, Turkey, pp. 423-428.

[7]  F. J. Martinez, C.-K. Toh, J.-C. Cano, C. T. Calafate, P.Manzoni, "A Street Broadcast Reduction Scheme (SBR) to Mitigate the Broadcast Storm Problem in VANETs," Wireless Personal Communications, vol.56, pp. 559–572, Feb. 2011.

[8]  A. Y. Al-Dubai, M. B. Khalaf, W. Gharibi and J. Ouenniche, "A new adaptive probabilistic broadcast protocol for vehicular networks," IEEE 81st Vehicular Technology Conference (VTC 2015), Glasgow, pp. 1-5.

[9]  M. Chekhar, K. Zine-Dine, M. Bakhouya and A. Aaroud, "A dynamic threshold-based probabilistic scheme for broadcasting in ad hoc networks," 15th International Conference on Intelligent Systems Design and Applications (ISDA 2015), Marrakech, pp. 511-516.

[10] J. Sospeter, D. Wu, S. Hussain, and T. Tesfa, "An Effective and Efficient Adaptive Probability Data Dissemination Protocol in VANET," Data, vol. 4, no. 1, p. 1, 2018.

[11] A. Mohammed, M. Ould-Khaoua, L. M. Mackenzie and J. Abdulai, "Dynamic probabilistic counter-based broadcasting in mobile ad hoc networks," 2nd International Conference on Adaptive Science & Technology (ICAST 2009), Accra, pp. 120-127.

[12] S. Panichpapiboon and L. Cheng, "Irresponsible forwarding under real intervehicle spacing distributions," IEEE Transactions on Vehicular Technology, vol. 62, no. 5, pp. 2264-2272, Jun.2013.

[13] M. B. Khalaf , A. Y. Al-Dubaia  and G. Minb, "New efficient velocity-aware probabilistic route discovery schemes for high mobility Ad hoc networks", Journal of Computer and System Sciences,vol.81, pp.97-109, Feb.2015.

[14] W. Wang, T. Luo, H. Kang, and M. Elhoseny, "A local information sensing-based broadcast scheme for disseminating emergency safety messages in IoV," Mob. Inf. Syst., vol. 2019, 2019.

[15] G. Koufoudakis, K. Oikonomou, K. Giannakis, and S. Aïssa, "Probabilistic flooding coverage analysis for efficient information dissemination in wireless networks," Comput. Networks, vol. 140, pp. 51–61, 2018.

[16] M. Chaqfeh and A. Lakas, "Speed adaptive probabilistic broadcast for scalable data dissemination in Vehicular Ad Hoc Networks," International Wireless Communications and Mobile Computing Conference (IWCMC 2014), Nicosia, pp. 207-212.

[17] Y. Mylonas, M. Lestas, A. Pitsillides and P. Ioannou, "Speed Adaptive Probabilistic Flooding for vehicular ad-hoc networks," IEEE 22nd International Symposium on Personal, Indoor and Mobile Radio Communications (PIMRC.2011), Toronto, 2011, pp. 719-723.

[18] A. Mostafa, A. M. Vegni, and D. P. Agrawal, "A probabilistic routing by using multi-hop retransmission forecast with packet collision-aware constraints in vehicular networks," Ad Hoc Networks, vol.14, pp.118-1129, 2014.

[19] N. Wisitpongphan, et al., "Broadcast storm mitigation techniques in vehicular ad hoc networks," IEEE Wireless Communications, vol. 14, no. 6, pp. 84-94, Dec. 2007.

[20] B. Yuan, A. jie and Z. Huibing,"Location aided probabilistic broadcast algorithm for mobile Ad-hoc network routing," The Journal of China Universities of Posts and Telecommunications, vol.24, pp. 66-71, 2017.

[21] I. A. Khan, A. Javaid  and H. L. Qian, "Distance-based dynamically adjusted probabilistic forwarding for wireless mobile Ad Hoc Networks," 5th IFIP International Conference on Wireless and Optical Communications Networks (WOCN 2008), Surabaya, pp. 1-6.

[22] S. Eichler, "Performance evaluation of the IEEE 802.11p WAVE communication standard," 66th IEEE Vehicular Technology Conference (VTC 2007), Baltimore, MD, pp.2199-2203.

[23] F. J. Martinez, et al., "Evaluating the impact of a novel warning message dissemination scheme for VANETs using real city maps", 9th International IFIP- TC6 Networking Conference, Chennai, India, pp.265-276, May. 2010.

[24] McCanne, S., Floyd, S.: The Network Simulator - ns-2.

[25] D. Krajzewicz, G. Hertkorn, C. Rossel, and P. Wagner, "SUMO (Simulation of Urban MObility)—An open-source traffic simulation," Proceedings of the 4th Middle East Symposium on Simulation and Modelling (MESM 2002), Sharjah, UAE, pp.183–187.

[26] Standard Specification for Telecommunications and Information Exchange between Roadside and Vehicle Systems, _5GHz Band Dedicated Short Range Comm. (DSRC) Medium Access Control (MAC) and Physical Layer (PHY) Specifications, ASTM E2213-03, Sept. 2003.

# Domain and Schema Independent Semantic Model Verbalization: A Conceptual Overview

Kaneeka Vidanage[1], Noor Maizura Mohamad Noor[2], Rosmayati Mohemad[3], Zuriana Abu Bakar[4]

School of Informatics and Applied Mathematics, University Malaysia Terengganu (UMT)

Kuala Nerus, Terengganu, Malaysia

*Abstract*—**Semantic Web-based technologies have become extremely popular and its a success that has spread across many domains, additional to the computer science domain. Nevertheless, the reusability aspects associated with the created and available semantic knowledge models are very low. The main bottleneck associated with this issue is, the difficulty associated in understanding the complex schema of a knowledge model created and barriers associated with querying the knowledge models using SPARQL or SQWRL query formulations. This research emphasizes on proposing a verbalizer which can go beyond existing Controlled Natural Language (CNL) type verbalizers and to verbalizer knowledge stored in a knowledge model file written in either RDF or OWL format, despite its domain and schematics.**

*Keywords*—*Ontology; OWL; RDF; Verbalize; Schema*

## I. INTRODUCTION

Ontologies are domain rich conceptualizations [1]. That is the definition given for ontologies by Spasic et.al in [1]. Resource Description Framework (RDF) and the Ontology Web Language (OWL) are the most prominent and World Wide Web Consortium (W3C) accredited standards for creating ontologies [2]. The initial idea of ontologies was elicited from the concept of Semantic Web by Tim Berners-Lee for the first time [3]. However, though the idea initially emerged in 2001, by 2013, almost more than 4 million web domains have incorporated semantic web technologies to their web sites [4]. This clearly depicts the massive growth of semantic web across the entire globe proving its remarkable success. Further to that, as claimed by Feigenbaum in [5] and Kashyap in [12] the noteworthy feature related to semantic web`s utmost success is the potential of both human and machine readability of semantic web`s knowledge representations.

The concept of ontologies that emerged from the initial idea of semantic web can be recognized as a very effective technology among contemporary computer researchers and enthusiasts. For means of justification, it can be easily pointed out, that, already thousands of ontologies developed for a variety of purposes are available in online repositories, almost with free accessibility. To name a few repositories where these predefined ontologies are available will be Vocab.Org [14], Swoogle [15], LOV [16], Protégé Wiki [17], AberOWL [18] and BioPortal [19]. Among them also, it's significant to emphasize both [18] and [19] are ontology repositories solely with human bioscience and diseases related aspects. This is critical evidence to point out that the ontologies as domain rich

conceptualizations [1] are doing extremely well in other domains as well, without limiting itself to the computer science domain.

The taxonomic structure, positioning the individuals, assertions between individuals, object and data properties can be effectively visualized through a tool like Protégé [6] or Top Braid composer [7]. Even though, graphical visualization will also not be adequate in most cases, as it`s not only the computer scientists or ontologists who are expecting to seek and gain the benefits of ontologies. Even for computer scientists or ontologists as well, it would be a really challenging task to understand the schema of an ontology developed by another set of researchers [13]. On the other hand, it has been stated that the development of an ontology from scratch is not an easy task as so far, no 100% automated mechanisms are available on ontology construction. Human intervention is essential [20]. Therefore, as claimed by [15][16], methods and mechanisms need to be sought after to enhance the reusability and overcome technological barriers associated, with an understanding of already created ontologies. The effective outcome of this would be, knowledge dissemination associated with existing ontologies will be further improved, enhancing the reusability aspects as well. Additionally, it will prevent a precious piece of information resource being stagnated on the internet after serving, only the one specific purpose it had been created for, which is recognized as an utmost cognitive waste as well by [21].

The ontology-based applications are not only limited to the computing domain. Medical sciences, pharmaceutical sciences [8][9], library sciences [5], law [10], criminology sciences [11] and ample other industries also comprehensively utilize the benefits of ontologies. This brings out the argument, potential capabilities of ontologies are not only sought after within the computer domain. As already conversed, several other disciplines are also very much keen on integrating the capabilities of ontologies to fulfil their discipline-specific requirements as well.

This setting will clearly open up the atmosphere to point out the greatest two bottlenecks associated with semantic web based ontologies, which will be leading to the research question discussed in this paper.

Firstly, understanding of the schema of an ontology written in RDF or OWL is a greatly challenging task even for computer scientists or ontologists as well. Therefore, for non-technical consultants like medical professionals, lawyers, criminologists, it would be a great obstacle [10] [11] [13].

Because without properly understanding the schema of the ontology, queries cannot be written to fulfil appropriate knowledge requirements. Secondly, writing of SPARQL or SQWRL queries to prorogate knowledge retrieval could be mostly an infeasible and unfair task to be expected from a non-technical specialist [22] [23] [24].

Therefore, as already conversed, these two issues will act as critical bottlenecks hindering the effective usage of semantic technologies within and outside of the computing domain. One of the potential solutions to overcome this technical barrier is to introduce ontology verbalizers. Verbalizers are capable of extracting knowledge represented in an OWL or RDF knowledge model and presenting it in a human understandable natural language [24].

But there are several problems associated with existing ontology verbalizers as well, hence most of them are domain and schema dependent, which means, they can work only with one domain as they have been tightly glued to one specific onology`s schema only [25][26]. The other issue is most of the verbalizers cannot work with both RDF and OWL formats and they work with either one of these and not with both, which again increases complexity in finding a suitable verbalizer for a required task [26] [27]. Eventually, most of the existing verbalizers can verbalize the knowledge in an ontology to a Controlled Natural Language (CNL) format only. CNL is a primitive English representation of triple sequences stored in an ontology model, which is not a conversational and readable English output which could be understood by anyone [28].

All these bottlenecks form the pathway to the research question to be discussed in this research which is to be "How to effectively verbalize both OWL or RDF based ontology, despite its domain and schema?"

The remaining section of the paper will discuss, about related works, methodology, results and discussion, evaluation and conclusion, respectively.

## II. RELATED WORKS

It`s already conversed in the introduction section as well, there are two critical bottlenecks recognized to be hindering the reusability of existing knowledge models as well as adversely affecting the use of ontologies in other domains as well. To quickly revise, firstly the complexity associated with comprehension of the schematics, as without properly knowing the schematics of the ontology, writing appropriate queries for the knowledge retrieval would be infeasible [10] [11] [13]. Secondly, even after the hurdle of comprehending schema is achieved as the next step, writing of accurate SPARQL or SQWRL queries to achieve the knowledge retrieval demands [22] [23] [24] will become a critical challenge. Users should have a sound knowledge about triple concepts of the ontologies and the relevant syntaxes as well as RDF and OWL axiom related concepts to properly write a query to fulfil knowledge requirements.

It was already stated in the introduction section, it`s not always computer specialists or ontologists only who will be seeking the usage of ontologies [13]. Hence, these challenges would hinder the spreading of benefits of the ontologies to wider audiences within and out of the computer science arena.

One possible potential to overcome this barrier is to use ontology verbalizers. Nevertheless, there are ample issues associated with ontology verbalizers as well, as already conferred in the introduction. The researchers will investigate that aspect more deeply through the assessment of existing verbalizers.

Two such pieces of evidence for verbalizers are [29] and [30], which are acting as domain and schema dependent because both of these verbalizers are statically mapped to DBpedia and accommodation ontology. Hence, those verbalizers are not open to verbalize any other ontology file fed into it. MIKAT [31] is another verbalizer which is specifically defined for the breast cancer domain. This verbalizer acts by providing necessary assistance and guidance to the clinical investigations made by the consultants on their patients, related to breast cancer ailments and diagnosis [31]. In the same way [32] points outs another verbalizer which is specifically defined for the colonoscopy domain. This verbalizer has the capability of annotating video footages of an ongoing colonoscopy. Therefore, none of the verbalizers discussed above is capable of functioning as a generalistic verbalizer.

On the other hand, Noy et.al [33] pointed out, the most popular two formats associated with ontological knowledge representation are RDF and OWL. Even though most of the verbalizers available currently cannot work with both, but only with either OWL or RDF which is another bottleneck to be sorted out, when finding a verbalizer for a verbalization task. For instance, in [34] there is one such verbalizer which can work only with OWL and not in RDF format.

The other issue associated with the verbalizers are, they have still not reached the level of verbalizing the knowledge in the form of conversational English which can be read and understood by everybody. Most verbalizers extract and present the knowledge in CNL formats as already discussed in the introduction section as well. Attempt-to-Control–English (ACE) is one such popular form of CNL output [30] [35]. ACE is again a primitive English representation extracting out the triple arrangement in the knowledge file and it`s not enhanced as conversational English which could be read and understood by everybody.

Therefore, it's very apparent; there is a research gap to be addressed, on verbalizing an ontology despite its RDF or OWL as well as regardless of its domain and schema as well. In other words, the requirement of a generalized verbalizer, which can verbalize knowledge with a more mature level than ACE, is the research gap to be addressed.

## III. METHODOLOGY

The main emphasis of this research is to come up with a generalized verbalizer, which can extract knowledge from ontologies despite their domain and schema as well as regardless of whether they are written in RDF or OWL formats. In fulfilling these goals, as the initial step, the comprehensive literature analysis was conducted via seeking for latest research and journal articles from credible repositories such as Springer, ACM Science Direct, IEEE etc. Keywords such as "ontology verbalizers, generalistic verbalizers, domain and schema independent verbalizers etc.."

are used to streamline the search results received from the research repositories mentioned above.

Even though lots of valid pieces of information has been collected, however, a proper solution addressing the research gap of generalized verbalizer which can work with both RDF and OWL, despite the domain and schema is not located. Further to that, it is also found, there is a deficiency issue in terms of verbalizing the knowledge as most of the existing verbalizers have not reached to a level beyond ACE as already conferred in detail, in the related works section.

Consequently, all these facts collected created a solid platform to further brainstorm and to continue the research. After completion of multiple brainstorming sessions with field experts and consultants, eventually, the following execution flow is derived as the initial step of extracting the required raw facts from the RDF / OWL knowledge models to commence up with the verbalization process. The proposed flow for the required facts extraction is denoted below, as in Fig. 1.

As illustrated in Fig. 1, the initial step is to check for the format of the knowledge model type. Because depending on its RDF or OWL appropriate extraction procedures needs to be triggered. Unless the user has to be notified with a suitable error message, claiming uploaded file type is not supported etc. Once the format verification stage is passed, the information extraction phase can be continued.

It`s decided to extract information, in sequential order of the individuals one at a time and one after another, rather than extracting information from here and there of the knowledge model, as it would adversely affect to the coherence related with further processing of the extracted information.



Fig. 1. Information Extraction Process from RDF / OWL Ontology File.

The process would be to select the first individual located in the knowledge model file. Here the individual means the entity of the knowledge instance located in the knowledge file. Then a sequential scan can be conducted throughout the entire document to extract all, subject, object, data properties, object properties, axioms, schematic information etc of the considered individual. Then using separately defined decision analysis and information extraction methodologies, all those information extracted can be again verified and carefully stored in a series of database tables designed according to the schematic structure specified below in Fig. 2. The individual element extraction processors will not be discussed in this paper as it is out of the scope of this paper, whereas the main emphasis of this paper will be on the verbalization process.

Here as denoted in Fig. 2, for each of the individual captured, an autogenerated fact Id can be introduced as a primary key as a mechanism of preserving information consistency associated with the respective individual. Then, using inheritance, properties are derived as data properties and object properties. These properties could not be overlapping, in the sense; a data property cannot be an object property and vice-versa. That`s why the "OR" disjoint constraint is used. Apart from that, another schematics relation is also introduced in the schema to take a track of special RDFS or OWL axioms linked with the individual's expressions. This measurement is taken to overcome the possible information losses associated with RDFS or OWL constraints defined in the knowledge model file when describing the individual's capabilities, domains, ranges etc.

Additionally, another table is defined to keep a track of individual's contexts. Because when it comes to verbalizing, the knowledge stored in an ontology, tracking the context would be very important in expressing the proper meaning out to the end user. It is intended to use discourse representation theory proposed by [36] [37] for the purpose of capturing the individuals' contexts. Fig. 3 mentioned below depicts the overall process flow associated with the discourse representation theory applied for this research. But the entire concept of discourse representation theory is not fully elaborated here, as it is outside the bounds of this paper. Therefore, the interested reader is encouraged to read the article mentioned in [37].

After completion of the process associated with extracting facts from the knowledge model file (depicted in Fig. 1) all extracted facts will be stored in the database schema proposed in Fig. 2. Then, that information can be again accessed in an individual-specific manner, assuring the triple sequence order to perform context assessment as per the discourse representation theory, as illustrated in Fig. 3 above. All extracted pieces of information, associated with the individual can be fed to a hash map and can perform, part of the speech tagging assessment and lexical analysis to trace potential changes causing on contextual differentiation. Then accordingly, context alert flag needs to be updated and it has to be supplied back to the verbalization module, along with the suggested pronouns to be used, which is technically referred as the discourse referent.

Fig. 2.   Database Schema.



Fig. 3.   Discourse Representation Theory Associated Process Flow.

Now all steps are in line to commence the verbalization process. The flow associated with the verbalization mechanism will be discussed in detailed under, results and discussions section, which is the next part of this research article.

## IV. RESULTS AND DISCUSSION

Up to now, the process associated with extracting important pieces of information from the RDF or OWL knowledge files, storing them in the database schema, application of discourse representation theory to ensure context sensing is, already discussed and illustrated in Fig. 1, 2 and 3, respectively.

The next important aspect is to discuss the verbalization process in detail. Upon the completion of the information extraction phase illustrated in Fig. 1, all individual-specific information is stored in the database schema presented in Fig. 2. Further, as in Fig. 3, discourse representation theory is applied to ensure context sensing information is recognized and proper discourse referents are introduced and context flags are updated as required. Then all these processed information, stored back in the database can be retrieved again to a hash map. It`s very important to keep on track, only the specific individual associated information is required to be extracted from the database. This will be feasible, as the information stored in the database is also governed and linked via entity and referential integrity constraints defined in the database schema.

The subsequent step would be to apply the Rapid Automatic Keyword Extraction (RAKE) algorithm [38] to all individual-specific lexical information derived as of one pool. RAKE algorithm will compute correlations amidst all the other individual-specific lexicons, within the pool and will derive a context relevancy value.    Then all these lexicons' context relevancy values, lexicons, additionally added identification index values, need to be carefully stored in a temporary table which will be used later as a look-up grid for appropriate lexicon extraction.   This process is graphically visualized in Fig. 4.

Once all these individual specific, lexicon`s context relevancy values, are extracted to the temporary table, the next step, intended is to apply, K-Means clustering algorithm [39]. K-Means is an unsupervised machine learning clustering algorithm, which could be used to cluster heterogeneous pieces of information into mostly feasible homogenous sets of clusters. Hence the intended goal is on generic verbalization, it has been decided to have 04 defined clusters representing, "introduction related—C1", "elaboration related–C2", "analysis related–C3" and eventually "conclusions related—C4" as those would be the ideal coverage which could be expected in terms of generic verbalization. Therefore, in terms of the K-Means clustering algorithm, it has been decided to specify K value as K=4, representing the four clusters from C1 to C4 derived above.

The ultimate expectation would be, once the K-Means clustering algorithm is applied to the consolidated individual, lexicon`s context relevancy values, need to be segmented into a mostly optimal 04 clusters, where within the cluster, information has to be mostly homogeneous. K-Means clustering algorithm works on the underlying concept of Euclidian distance. Therefore, the RAKE algorithms` context relevancy values would be much useful to segment the lexicons, considering their context relevancy values and grouping them into 04 homogeneous clusters.    The entire process associated with the application of the K-Means algorithm is graphically illustrated in Fig. 5.



Fig. 4.   Applying RAKE Algorithm to Get Context Relevancy Values.



Fig. 5.   Applying the K-Means Algorithm.

Afterwards, context-sensitive values residing in each cluster can be converted up back to the individual`s lexicons via referring them back with the look-up grid maintained earlier. This will result in four clusters of homogenous lexical groups belonging to the same individual.

The next step is to define four specified templates to cover up the verbalization scope of "introduction", "elaboration", "analysis" and "conclusion" which is synchronized up with the four clusters derived from K-Means clustering. But the issue is hence the K-Means clustering is an unsupervised clustering algorithm, and it`s not practical to directly mention which cluster can be mapped as "introduction" or "elaboration" or "analysis" or "conclusion" or vice-versa.

To overcome that issue, first –order-logic based Prologue rules can be introduced inside each of the four templates proposed, denoting their phased specific execution level, as "introduction" or "elaboration" or "analysis" or "conclusion". Then at the time of inferencing, these rules will execute and fill the templates with appropriate data elements derived from the knowledge model file. Then a thematic mapping algorithm like Latent Semantic Indexing (LSI) [40] can be used to determine, which template matches with which cluster. LSI is a very intelligent algorithm, which does a lot more than simple keyword comparison. The Single Value Decomposition (SVD) mechanism used in LSI functions as a dimensionality reduction mechanism via integrating all related dimensions to one specific theme. This SVD strategy is a key contributor resulting in the refined intelligent behaviour of the LSI algorithm [41]. Here the template does contain only a gist of information and using it as the triggering point, with the help of LSI algorithm, the appropriate cluster can be identified. The process of clusters and template matching is clearly noted in Fig. 6 mentioned below.

Eventually, via referring to the context-sensitive values of each of the lexicons belonging to a specified cluster, all lexicons within the cluster can be sorted from max to min of its context sensitive values. This will allow locating the nucleus and the satellites as per the Rhetoric Structure Theory (RST) which plays a vital role in micro-planning of sentences [42]. This will improve the readability of the text generated from the verbalizer. According to RST, the nucleus will postulate the most important fact associated within the context and the satellites will become the associative facts which will be elaborating the nucleus.

Ultimately, prologue governed phase specific templates can instruct the SimpleNLG framework to carry out the domain and schema independent verbalization process. The final step associated with domain and schema independent verbalization is as mentioned in Fig. 7.



Fig. 6. Cluster and Template Matching.



Fig. 7. Final Steps of the Verbalization Process.

## V. EVALUATION

For the evaluation purposes of the domain and schema independent verbalizer, the process depicted in Fig. 8 below is utilized. The crime domain is selected as the application domain to test the verbalizer. The reason for that is, the research team already works with the government police and crime officers for several other crime analysis types of research, ongoing.

Initially, an existing knowledge model on the domain of crime knowledge is sought after. Then, later on, it's decided, in order to more accurately perform the evaluation process, to create a simple knowledge model associated with the sub-discipline of evidence handling related to a crime scene. Few of the crime officers (i.e. -6) are interviewed and a potential knowledge model on the domain of evidence handling was created via the use of Protégé in RDF. Then through Protégé Integrated Development Environment (IDE), same RDF knowledge model is converted to it`s OWL counterpart.

Subsequently, the RDF version of the knowledge model is uploaded to the verbalizer and the contents are verbalized in English. Then, the generated output is provided to a few of the crime experts in the department and they are asked to verify the effectiveness in the facets of understandability, information loss and reliability aspects. Henceforth, the same process is followed for the OWL version of the knowledge model as well. Interviewed feedbacks obtained from the crime specialists are thematically and statistically assessed via thematic analysis evaluation methodology.

As specified in the literature, thematic analysis is a very effective mechanism, which can be applied to any discipline in evaluating qualitative data [43]. Use of thematic analysis in computer science discipline is also very prominent. For instance, in [44] Porter et al. have stated the effectiveness in the use of thematic analysis for designing a proper UI/UX for e-personas leading into an e-government identification project. Likewise as suggested by multiple pieces of evidence [44-45], in developing information systems which make close interactions with humans, critical emphasis should be given to the interaction experiences [45] Without utilizing proper subjectivist approaches like thematic analysis, unforeseen negative interaction experiences cannot be extracted, which would ultimately lead to a failed deployment of the information system [43-45].

The suggested verbalizer in this research is also going to be a system, which closely interacts with the end user. Because the content verbalized by the verbalizer should be understood by the end user with almost no ambiguity. In order to check that dimension, users' personal interaction experience with the

system is very important. Therefore, these arrangements clearly point out the suitability of the thematic analysis evaluation technique to be used for this research as well. In the process of thematic analysis, the first step would be, detailed and repetitive reading of the interview feedbacks, documented through the interview process.

Then, rationales emerged out are carefully analyzed. Ideologies which have a close coherence are aggregated into groups and organized as facets. When interviewing more and more end users/subject specialists (i.e. crime officers), coherent feedbacks obtained are caused to accumulate the facet counters defined for each facet group, representing a numerical/statistical overview on the insights collected. The concept of Evaluation Onion proposed by [46] is utilized in mapping the facts emerged out from repetitive interview feedback reading, with facet criteria presented as facet groups.

Evaluation onion, which is also referred to as CCP framework [46] is recommended by Eslami et al. (2017) in [47] and many other researchers as well, in determining evaluation criteria's be used in the assessment of Information Systems which closely interacts with the ends users. Fig. 9 denoted below will present the overview of evaluation onion concept which is also referred to as CCP framework.

Table I will clearly illustrate how the recognized facts from interview feedbacks can be mapped within the CCP framework. Bold "Wh " question criteria's will clearly demonstrate how the CCP framework aspects are interlinked with interview comments facets.



Fig. 8. Complete Evaluation Process.



Fig. 9. CCP Framework by Farby and others (Farbey, Land, & Targett, 1993).

TABLE I. FACETS MAPPING INTO CCP FRAMEWORK

| Facets | Indexed Code | CCP framework mapping question |
|---|---|---|
| Verbalizer Accuracy. | **PAC** | **What** is the level of accuracy experienced in this system? |
| Verbalizer Applications. | **PAP** | **How** this system could be useful for end users? |
| Verbalizer Assistance. | **PAS** | **Who** would be benefitted by the use of this system? |
| Verbalizer Importance | **PAI** | **Why** this research is important to end-users / experts? |

Fig. 10 depicts the distribution of frequencies associated with each facet mentioned above.



Fig. 10. End users Response Distribution Against Defined Facets.

As the second phase of evaluation, a statistical assessment is also conducted with the use of evaluation metrics. True Positives (TP), True Negatives (TN), False Positives (FP), False Negatives (FN), associated with the verbalization process is verified. The crime officers involved in the interview process, for the output evaluation, also involved in the creation of the knowledge model. Hence, they have the proper ideology in verifying the accuracy aspects associated with the verbalization, to recognize, are there any information losses, misinterpretations etc.

Subsequently, typical test measurements such as recall, precision and F-measures are derived, on the True Positive (TP), True Negative (TN), False Positive (FP) and False Negative (FN) confusion matrix element values derived. Table II depicts the statistics derived in one specific verbalization instance associated with 315 text expressions. Fig. 11 illustrates the calculation process initiated from the confusion matrices concept.

Here TP denotes accurate verbalization of expressions extracted from the knowledge model and FPs denotes incorrect verbalizations of the existing expressions in the knowledge model. Likewise, FN denotes misinterpretations resulted in verbalization process and eventually, TN denotes exclusion of less important axioms which occurred during the verbalization process. This will mainly occur via the functionality of the RAKE algorithm discussed above.

Fig. 11. Formulation of Test Statistics for Verbalization.

Quantitative test statistics derived for the verbalization took place for a specified instance, is logged in Table II.

TABLE II.     TEST STATISTICS FOR A VERBALIZATION INSTANCE

| Measurement | Accomplishment |
|---|---|
| Sensitivity | 0.78 |
| Precision | 0.90 |
| Accuracy | 0.86 |
| F-Measure | 0.8 |

The verbalizer proposed in this research is qualitatively (via thematic analysis and CCP framework) and quantitatively assessed (via test statistics) as conversed above.

## VI. CONCLUSION

The research gap attempted to address in this research article is to derive a potential resolution on the issue of domain and schema independent ontology verbalization, despite RDF or OWL formats. As already conversed in the paper, even there are few existing verbalizers located; most of them have multiples of issues, which are already discussed in related works section etc. Therefore, as a means of overcoming those deficiencies, this new conceptual arrangement of the verbalizer and its internal algorithmic functionalities are reviewed and evaluated in this paper. It is assumed, these findings will further contribute in making use of semantic technologies more applicable and addressable across a vast range of domains, despite the technical bottlenecks. However, the verbalizer proposed needs to be tested on several multiple domains, to further enhance and stably justify its accuracy.

## VII. LIMITATIONS AND FUTURE WORK.

The most challenging aspect of the domain and schema independent verbalization is inability to use a large dataset as the training corpus to train the verbalizer to effectively perform the verbalization process, on any given domain. Because, at the stance, used a specified training dataset to train the verbalizer it will not further be a domain-independent verbalizer. Therefore, to rationally handle this requirement, a combination of algorithms and techniques such as RAKE algorithm on key phrase extraction, K-Means unsupervised learning algorithm, Prologue enabled phased specific templates and Rhetoric Structure Theory and SimpleNLG framework has been utilized as a pipeline of technologies (see Fig. 12).



Fig. 12. Pipeline of Technologies.

At the moment, the verbalizer is only evaluated on the crime domain and it has to be further tested on other domains, such as medicine, law, management etc. Then test statistics such as recall, precision, F-measures can be derived for the verbalizer assessing the functionality on several other domains and, it will yield to derive a more normalized and stable outcome on the verbalizer performance.

REFERENCES

[1]  Spasic, I., Ananiadou, S., McNaught, J., & Kumar, A. (2005). Text mining and ontologies in biomedicine: Making sense of raw text.Briefings in Bioinformatics, 6(3), 239-251.doi:10.1093/bib/6.3.239

[2]  Caldarola, E. G., & Rinaldi, A. M. (2016). An Approach to Ontology Integration for Ontology Reuse. 2016 IEEE 17th International Conference on Information Reuse and Integration (IRI). doi:10.1109/iri.2016.58

[3]  Berners-Lee, Tim (May 17, 2001). "The Semantic Web" (PDF). Scientific American.

[4]  Ramanathan V. Guha (2013). "Light at the End of the Tunnel". International Semantic Web Conference 2013 Keynote.

[5]  Lee Feigenbaum (May 1, 2007). "The Semantic Web in Action". Scientific American

[6]  Musen, M. A., &amp; The Protégé Team. (2013). Protégé Ontology Editor. Encyclopedia of Systems Biology, 1763-1765. doi:10.1007/978-1-4419-9863-7_1104

[7]  Topbraid Enterprise Data Governance (2019) Retrieved March 6, 2019, from, https://www.topquadrant.com/products/topbraid-enterprise-data-governance/

[8]  Bontcheva, K., & Wilks, Y. (2004). Automatic Report Generation from Ontologies: The MIAKT Approach. Natural Language Processing and Information Systems, 324-335. doi:10.1007/978-3540-27779-8_28

[9]  Bao, J., Cao, Y., Tavanapong, W., & Honavar, V. (2004). Integration of Domain-Specific and DomainIndependent Ontologies for Colonoscopy Video Database Annotation. Artificial Intelligence Research Laboratory-Iowa State University.

[10] Ku, C., & Leroy, G. (2014). A decision support system: Automated crime report analysis and classification for e-government. Government Information Quarterly, 31(4), 534-544. doi:10.1016/j.giq.2014.08.003

[11] Pinheiro, V., Furtado, V., Pequeno, T., & Nogueira, D. (2010). Natural Language Processing based on Semantic inferentialism for extracting crime information from text. 2010 IEEE International Conference on Intelligence and Security Informatics.

[12] Kashyap V (2008) Ontologies and Schemas. The Semantic Web, 79-135. doi:10.1007/978-3-540 764526_5

[13] Rusu, D., Dali, L., Fortuna, B., Grobelnik, M., & Mladnnic, D. Triple Extraction from sentences. Paper presented at Technical University of Cluj-Napoca, Romania.

[14] Davis, I. (2014). vocab.org - A URI space for vocabularies. Retrieved February 16, 2019, from http://vocab.org/

[15] Yu, L. (2007). Swoogle. Introduction to the Semantic Web and Semantic Web Services,145-157. doi:10.1201/9781584889342.pt3

[16] Vandenbussche, P., Atemezing, G. A., Poveda-Villalón, M., &Vatant, B. (2016). Linked Open Vocabularies (LOV): A gateway to reusable semantic vocabularies on the Web. Semantic Web, 8(3), 437-452. doi:10.3233/sw-160213

[17] Protege. (2018). Protege Ontology Library - Protege Wiki. Retrieved February 16, 2019, fromhttps://protegewiki.stanford.edu/wiki/Protege_Ontology_Library

[18] Slater, L., Gkoutos, G. V., Schofield, P. N., &Hoehndorf, R. (2016). Using AberOWL for fast and scalable reasoning over BioPortal ontologies. Journal of Biomedical Semantics, 7(1). doi:10.1186/s13326-016-0090-0

[19] Faria, D., Jiménez-Ruiz, E., Pesquita, C., Santos, E., &Couto, F. M. (2014). Towards Annotating Potential Incoherences in BioPortal Mappings. The Semantic Web – ISWC 2014, 17-32. doi:10.1007/978-3-319-11915-1_2

[20] Trokanas, N., & Cecelja, F. (2016). Ontology evaluation for reuse in the domain of Process Systems Engineering. Computers & Chemical Engineering, 85, 177-187. doi:10.1016/j.compchemeng.2015.12.003

[21] Zenuni, X., Raufi, B., Ismaili, F., & Ajdari, J. (2015). State of the Art of Semantic Web for Healthcare. Procedia - Social and Behavioral Sciences, 195, 1990-1998. doi:10.1016/j.sbspro.2015.06.213

[22] Chergui, W., Zidat, S., & Marir, F. (2018). An approach to the acquisition of tacit knowledge based on an ontological model. Journal of King Saud University - Computer and Information Sciences. doi:10.1016/j.jksuci.2018.09.012 9.

[23] Alavi, M., Leidner, D.E., 2001. Knowledge Management and Knowledge Management Systems: Conceptual Foundations and Research Issues. Manag. Inf. Syst. Q. 25, 107–136. https://doi.org/10.2307/3250961. Anderson, J.R., 1983. The Architecture of Cognition. Harvard University Press, Cambridge, MA.

[24] Gutierrez-Basulto, V., Ibanez-Garcia, Y., Kontchakov, R., & Kostylev, E. V. (2015). Queries with Negation and Inequalities over Lightweight Ontologies. SSRN Electronic Journal. doi:10.2139/ssrn.3199213

[25] Williams, S., Third, A., & Power, R. (2011). Levels of organisation in ontology verbalization. ENLG. Retrieved from https://www.semanticscholar.org/paper/Levels-of-organisation-in-ontologyverbalisation-Williams-Third/08c6a058f5f78cf497 01d2534bf9 c6a f3683f9e9

[26] Habernal, I., & Konopík, M. (2013). SWSNL: Semantic Web Search Using Natural Language. Expert Systems with Applications, 40(9), 3649-3664. doi:10.1016/j.eswa.2012.12.070

[27] Poulovassilis, A., Selmer, P., & Wood, P. T. (2016). Approximation and Relaxation of Semantic Web Path Queries. SSRN ElectronicJournal. doi:10.2139/ssrn.3199265

[28] Kaarel Kaljurand and Norbert E. Fuchs. 2007. Verbalizing owl in attempt to controlled English. In Proceedings of Third International Workshop on OWL: Experiences and Directions, Innsbruck, Austria (6th–7th June 2007), volume 258

[29] Poulovassilis, A., Selmer, P., & Wood, P. T. (2016). Approximation and Relaxation of Semantic Web Path Queries. SSRN ElectronicJournal. doi:10.2139/ssrn.3199265

[30] Kaarel Kaljurand and Norbert E. Fuchs. 2007. Verbalizing owl in attemp to controlled English. In Proceedings of Third International Workshop on OWL: Experiences and Directions, Innsbruck, Austria (6th–7th June 2007), volume 258

[31] Bontcheva, K., & Wilks, Y. (2004). Automatic Report Generation from Ontologies: The MIAKT Approach. Natural Language Processing and Information Systems, 324-335. doi:10.1007/978-3-540-27779-8_28

[32] Bao, J., Cao, Y., Tavanapong, W., & Honavar, V. (2004). Integration of Domain-Specific and DomainIndependent Ontologies for Colonoscopy Video Database Annotation. Artificial Intelligence Research Laboratory-Iowa State University.

[33] Noy, N., & McGuiness, D. (2001). Ontology Development 101: A Guide to Creating Your First Ontology. Stanford University, Stanford

[34] Bontcheva, K., & Wilks, Y. (2004). Automatic Report Generation from Ontologies: The MIAKT Approach. Natural Language Processing and Information Systems, 324-335. doi:10.1007/978-3-540-27779-8_28

[35] Bojars, U., Liepins, R., Gruzitis, N., Cerans, K., & Celms, E. (2016). Extending OWL Ontology Visualizations with Interactive Contextual Verbalization. VOILA@ISWC.

[36] Lascarides, A., & Asher, N. (2008). Segmented Discourse Representation Theory: Dynamic Semantics With Discourse Structure. Computing Meaning, 87-124. doi:10.1007/978-1-4020-5958-2_5 .

[37] Van Eijck, J. (2006). Discourse Representation Theory. Encyclopedia of Language & Linguistics, 660668. doi:10.1016/b0-08-044854-2/01090-7

[38] Gupta, S., Mittal, N., & Kumar, A. (2016). Rake-Pmi Automated Keyphrase Extraction. Proceedings of the International Conference on Informatics and Analytics - ICIA-16. doi:10.1145/2980258.2980463

[39] Wei, S., Yonglin, O., Qingcai, Z., Jiaqiang, H., & Yaying, S. (2018). Unsupervised Machine Learning: Kmeans Clustering Velocity Semblance Auto-Picking. 80th EAGE Conference and Exhibition 2018. doi:10.3997/2214-4609.201800919

[40] Al-Anzi, F. S., & AbuZeina, D. (2018). Enhanced Search for Arabic Language Using Latent Semantic Indexing (LSI). 2018 International Conference on Intelligent and Innovative Computing Applications (ICONIC). doi:10.1109/iconic.2018.8601096

[41] Amini, B., Ibrahim, R., Othman, M. S., & Nematbakhsh, M. A. (2015). A reference ontology for profiling scholar's background knowledge in recommender systems. Expert Systems with Applications, 42(2), 913-928. doi:10.1016/j.eswa.2014.08.031

[42] MANN, W. C., & THOMPSON, S. A. (1988). Rhetorical Structure Theory: Toward a functional theory of text organization. Text - Interdisciplinary Journal for the Study of Discourse, 8(3). doi:10.1515/text.1.1988.8.3.243

[43] Lapadat, J. (2010). Encyclopedia of Case Study Research. In Encyclopedia of Case Study Research. SAGE Publications

[44] Porter, C., Sasse, M., & Letier, E. (2013). Giving a voice to personas in the design of egovernment identity processes. Research to Design: Challenges of Qualitative Data Representation and Interpretation in HCI-in BCS HCI.

[45] Adams, A., Lunt, P., & Cairns, P. (2008.). A qualitative approach to HCI research. Research Methods for Human–Computer Interaction, 138-157. doi:10.1017/cbo9780511814570.008

[46] Farbey, B. Land and Targett, D (1993). How to assess your IT investment. A Study of Methods and Practice. Butterworth Heinemann, Oxford

[47] Eslami Andargoli, A., Scheepers, H., Rajendran, D., & Sohal, A. (2017). Health information systems evaluation frameworks: A systematic review. International Journal of Medical Informatics, 97, 195-209. doi:10.1016/j.ijmedinf.2016.10.008

# A Modified Adaptive Thresholding Method using Cuckoo Search Algorithm for Detecting Surface Defects

Yasir Aslam[1]

Research Scholar, Department of Electronics and Communication Engineering, Noorul Islam Centre for Higher Education, Kumaracoil, India

Santhi N[2]

Associate Professor, Department of Electronics and Communication Engineering, Noorul Islam Centre for Higher Education, Kumaracoil, India

Ramasamy N[3]

Associate Professor, Department of Mechanical Engineering Noorul Islam Centre for Higher Education Kumaracoil, India

K. Ramar[4]

Principal &Professor, Department of Computer Science and Engineering, Einstein College of Engineering, Tirunelveli, India

*Abstract*—There are various mathematical optimization problems that can be effectively solved by meta-heuristic algorithms. The improvement of these algorithms is that they carry out iterative search processes which resourcefully act upon exploration and exploitation in spatial domain containing global and local optima. An innovative robust Cuckoo Optimization Algorithm (COA) with adaptive thresholding is proposed to solve the problem of detection and estimation of surface defects on metal coating surfaces. The proposed method is developed through implementing changes to COA and improved the performance. For improving capability of local search as well to keep the global search effect, the enhanced methods such as level set is associated with the proposed method. Also, the method adapts dynamic step size, adaptively changing with the search process for improving the rate of convergence and the ability of local search. The algorithm performance is scrutinized from the experimental analysis and results. Also, the segmentation effectiveness is further enhanced by adapting suitable methods for preprocessing and post processing. The comparison and analysis of the results accomplished with the proposed method and results of earlier methods shows superior performance of the proposed method.

*Keywords*—*Thresholding; surface defect; optimization; image processing; coated surface*

## I. Introduction

The quality control is a significant feature of today's highly competitive industry. Each manufacturing process output inspection is an imperative way to enrich the end product quality. The manual inspection of end products retards the whole procedure as it gets expensive, time consuming as well as impact the efficiency of human because of the unsafe atmosphere in industries. The inspection process has been automated and the methodology ought to be considered as resourceful layout of human intellect and comprehend together with rapidity of machine [1]. Automated visual inspection is an exceptionally imperative non-contact method in industries.

It detects diminutive defects which turn out as local anomalies relative to the adjacent background in the acquired image. A robust automated visual inspection method for identifying restrained defects in the pattern surface and the nondestructive testing technique become commonly used method for defect detection and classification. The microwave nondestructive testing technique [2] employed to detect defects in non-ceramic insulators. Currently, machine vision system has turn out to be the main stream nondestructive approach to resolve this kind of problem, reviewing its distinctiveness of quick response and non-contact. Accordingly, various algorithms have been developed to recognize and categorize the surface defect [3].

The method of active microwave thermography identifies the defects on steel surface and carbon fiber reinforced polymer based materials. The detection approach based on the computer vision technology make possible the compilation of surface quality insights and locates appropriate disruption all through the production process. The methods of defect detection could be organized into supervised and unsupervised manner. The supervised methods primarily depend on sufficient and standard set of training data. But, in exceptional instances, there exist an absolute defect dimensions set in an existent production environment. The unsupervised methods [4] on the basis of irregularity detection might be of immense practical significance, which is capable to find the anomalous regions via evaluation and comparison of local patches within the image. In general, existing methods for anomaly detection could be organized into three types such as, spectral, model based and statistical methods. Previously, statistical methods utilized to evaluate the texture through calculating spatial distribution of pixel intensities. The defects were detected with the first-order statistics in those approaches and the mean, variance and histogram based computations together with the second-order statistics depending upon the co-occurrence matrix.

Automated surface inspection is the field of utilizing the computer vision algorithms for surface inspection. The most surface defects are local variances or anomalies in homogeneous surface. The commonly used automated surface inspection algorithms [5] construct local features for surface defect detection. The statistical, structural, model based and filter based methods are the most important types of automated surface inspection methods. A statistical method determines the neighboring distribution of pixel standards and the structural methods form recurring patterns through dislocation and texture primitives.

The machine vision based method utilized for surface inspection as a non-contact detection technology in this decade which is capable to track the dynamic details of the object surface. The method makes use of the images confined by the camera for computing the target dimensions. The grey level of every pixel in the image characterizes the potency of the light reflected against the surface of the measured object. The pixels of an image correspond one-to-one with points on the surface of object and be utilized for computing the defects positions. For inspecting the surface defects of tunnels [6] this method has been applied. The new-fangled unidentified defects require restructuring the automated existing algorithms [7] so as to distinguish and categorize new defects, causing extended development cycles, impediments and have need of human endeavor toward sustaining the design and improvement of the system constantly.

In image analysis, an important process is thresholding, which turns down the color or grayscale image into binary. The simplest technique is thresholding to segment an image into regions exhibiting generic properties. The method constitutes black and white binary images from color or grayscale images with the intensity variation of all pixels to distinct values of zero or else one. Thresholding decreases the intensity of pixels lower than a definite value to zero (black), while the pixels higher than the selected value are specified as one (white). In image segmentation, for rapid evaluation the thresholding method is useful because of its ease and quick processing rate. The images might have different lighting conditions in different areas but utilizing a global threshold value may not be good choice, in such conditions the adaptive thresholding [8] is used.

In this paper we propose Modified Adaptive Thresholding Method using Cuckoo Search [9] Algorithm for Detecting Surface Defects, which is flexible to accord with these problems and further improves the segmentation performance. The thresholding technique will separate the image segments in to defected and non-defected regions. The proposed method is then compared with existing thresholding methods.

## II. IMAGE THRESHOLDING AND BINARIZATION

Initially for defect detection, the color image is changed to grayscale and the thresholding method is used to get the appropriate black and white image. The defect is highlighted in a distinctive color as white and the rest of the coating as black.

### A. Otsu's Method for Thresholding

Otsu's method is considered as a basic method to adapt a threshold value of an image. Through maximizing the between-class variance, the Otsu's method finds the optimum threshold $d^*$ as [10].

$$d^* = \underset{1 \leq d \leq K}{\operatorname{argmax}} \rho_R^2 \tag{1}$$

where, $\rho_R^2$ the between-class variance given by

$$\rho_R^2 = \beta_0(\omega_0 - \omega_d)^2 + \beta_1(\omega_1 - \omega_d)^2 \tag{2}$$

Where, the occurrence of probabilities of the two classes are $\beta_0$ and $\beta_1$, the average of the two classes are $\omega_0$ and $\omega_1$, the total average is $\omega_D$, the foreground and the background pixels are separated using the threshold value.

### B. Median-Based Extension of Otsu's Method

This method has been shown to be more vigorous in the resolving the threshold, in particular while in case of skewed and heavy-tailed distribution. The median-based Otsu's method is functional using median values rather than the mean values of the gray scale level distribution. Xue and Titterington find the mean absolute deviation (MAD) [11] from the median on the way to estimate the distribution rather than the variance. The threshold is selected by

$$d^* = \underset{1 \leq d \leq K}{\operatorname{argmin}} \{\beta_0 MAD_1(d) + \beta_1 MAD_2(d)\} \tag{3}$$

where, $MAD_1$ is the median for foreground class $F_0$ and $MAD_2$ is the median for background classes $F_1$ which can be as follows,

$$MAD_1(d) = \sum_{m=1}^{d} \frac{p(m)}{\beta_0(d)} |m - s_1(d)| \tag{4}$$

$$MAD_2(d) = \sum_{m=d+1}^{Z} \frac{p(m)}{\beta_1(d)} |m - s_2(d)| \tag{5}$$

where, $s_1$ and $s_2$ are the sample median for $F_0$ and $F_1$, respectively.

### C. Adaptive Thresholding with PSO

The adaptive thresholding in general takes a grayscale or color image as input and outputs a binary image describing the segmentation. A threshold has to be calculated for each pixel in the image and for this the binary image pixels are used to calculate the threshold value. If the pixel value is lower to the threshold value then the image is set to background; or else it set to a foreground. This can be viewed as an optimization problem to obtain an optimal threshold value for the adaptive threshold method. Therefore, Particle swarm optimization (PSO) is used to get an optimal adaptive threshold value [12]. PSO is surrounded to predict the behavior of birds in search for nourishment on a cornfield or fish school. The technique can capably determine best or close to best solutions in excessive search spaces. The swarm molecule movement in the search space is demonstrated in the following equations:

$$V_i^n = V_i^n + c_1 \cdot ma_1 \cdot (pb_i^n - x_i^n) + c_2 \cdot ma_2 \cdot (gb^n - x_i^n) \tag{6}$$

$$x_i^n = x_i^n + \delta V_i^n \tag{7}$$

From the equations (6) and (7), $c_1, c_2$ are the coefficients with the range of 2.0, $ma_1, ma_2$ are the independent random values developed in the limit between 0 and 1, $V_i^n$ is the velocity of i$^{th}$ particle. $x_i^n$, $pb_i^n$ represents the current position $i$ and the optimal fitness value of the molecule at the present iteration, $gb^n$ is the optimum global values in the swarm.

## III. PROPOSED METHOD

### A. Modified Adaptive Thresholding Method using Cuckoo Search Algorithm

A highly efficient defect detection system is proposed for identifying surface defects on metal coatings. The different steps involved in the proposed method are shown in Fig. 1.The high resolution images are used for the processing. Adaptive thresholding [13] utilized by which the gray scale image is taken as input. A threshold value is automatically calculated at every pixel of image and when the value of each pixel falls under the threshold, it is set with the background value or else foreground value is set. The threshold value is optimized by the proposed Cuckoo Optimization Algorithm (COA) method, which quickly finds the optimal threshold value for segmentation. The steps involved in defect detection are depicted in the block diagram. Initially the images are subjected to preprocessing with Contrast Stretching method. The Cuckoo optimization method is used by which the each group fitness function is evaluated and an optimal threshold value is calculated. Then with the level set method, region segmentation is performed. The final step in processing is morphological operation which is used to refine the segmentation.

### B. Preprocessing

It is an effortless method applied for image enrichment that takes effort by stretching the series of intensity values [14] incorporates to get better contrast in an image and extent preferred choice of values. The upper and lower pixel value limits must be specified before stretching can be performed over which the image is to be normalized. The lower and upper limits for 8-bit gray level images might be 0 and 255.

Consider the upper and lower limits $v$ and $u$ respectively. The highest and lowest pixel values [15] now present in the image as y and x respectively. The following function is used for scaling each of the pixels P in an image, where $P_{out}$ and $P_{in}$ are the output and input pixel intensities respectively at each individual pixel present in the image.



Fig. 1. Proposed Block Diagram.

$$P_{out} = (P_{in} - x)\left(\frac{v-u}{y-x}\right) + u \qquad (8)$$

### C. Cuckoo Optimization Algorithm (COA)

The COA local search involves, every cuckoo lays eggs surrounded by an explicit radius and that confides on the quantity of eggs, immigration of population performs global search. It is clear from the experiments exhibited, that COA [16] needs less number of iterations toward reaching the global optimum. Cuckoo Search Algorithm has some rules: a) Each of the cuckoo merely lay single egg by each instance while depot within arbitrarily nest preferred. b) Nests consisting of eggs with superior quality are considered as finest nest where about accepted above to subsequent formation. c) There is a fixed amount of host nest available. The probability in which the host bird detect the unfamiliar egg in its nest is taken as $p_a \in [0, 1]$ since the quantity of host nest obtainable is preset. The host bird moreover extinguishes the egg otherwise nest is disregarded whereas built a fresh nest.

Cuckoo Optimization Algorithm

*a)* Initialize the habitat of cuckoo with some random points on the profit function $c_p$ at habitat $x_1, x_2, \ldots, x_{N_{var}}$

$$\text{Profit} = c_p(\text{habitat}) = c_p(x_1, x_2, \ldots, x_{N_{var}}) \qquad (9)$$

where, $1 \text{x } N_{var}$ is an array of current living position or habitat of cuckoo

*b)* Contribute some eggs to each cuckoo

*c)* Define the Egg Laying Radius (ELR) for each cuckoo:

$$\text{ELR} = \propto \times \frac{\text{current number of cuckoo's eggs}}{\text{total number of eggs}}$$

$$\times \text{var}_{high} \text{var}_{low} \qquad (10)$$

*d)* Let cuckoos lay eggs inside their consequent ELR

*e)* Eradicate those eggs that are renowned by host birds

*f)* Assess the habitat of every newly grown cuckoo

*g)* Limit cuckoos maximum number in the location and eradicate those who live in worst habitat Group cuckoos and locate best group and choose goal habitat:

$$J = \sum_{j=1}^{MaxIter} \sum_{i=1}^{Cuckoopop} \left\| cuckoo_i^j - c_j \right\|^2 \qquad (11)$$

*h)* Let new cuckoo population immigrate in the direction of goal habitat

*i)* Stop if population size exceeds maximum iteration, if not, then go to step (b).

In Cuckoo Optimization Algorithm, the cuckoo with its egg prompt a optimization problem through candidate solution and consists the range $d_1 = [k_{l,1}, k_{l,2}, \ldots, k_{l,(j-1)}, k_{l,j}, k_{l,(j+1)}, \ldots, k_{l,n}]$ whereas $1 \leq l \leq t$ with its population size, $d_1$ the cuckoo on the l'th position, $k_{l,j}$ the j describes the dimension of it. To evaluate each habitat a fitness function is used by the COA [17] and targets to normalize it. The cuckoos are produced randomly at the initial step [18].

$$k_{l,j} = \text{rand x } M_j + z_{min,j} \qquad (12)$$

Here $Z_{min,j}$ and $M_j$ are j'th dimensions of vector, M and $Z_{min}$ it indicates range and minimum space of input correspondingly,

$$Z_{min} = [\ z_{min,1}\ , z_{min,2}\ ,...., z_{min,n}\ ] \qquad (13)$$

$$M = [M_1, M_2, ......, M_n\ ] \qquad (14)$$

Each one cuckoo lays numeral eggs in other birds nest contained by the explicit radius, KM. The value for cuckoo in the l'th position is

$$KM_l = \omega \text{xMxE}_l\ /\ TE \qquad (15)$$

where, $\omega$ is controller value of $KM$. $E_l$ is the amount of eggs within the cuckoo of l'th position whereas TE indicates entire eggs in population. The radius of egg at position p, $KM_{l,p}$ can be represented as

$$KM_{l,p} = rand * KM_l;\ t = 1, 2, ...., E_l \qquad (16)$$

The percent of low fitness cuckoo eggs are renowned which the host bird eliminates. If the size of population exceeds a maximum size provided, the low fitness weaker cuckoos are discarded. The existent cuckoo groups are partite into numerous groups. The average fitness is calculated for each group. The best group [19][20] containing the finest point is considered as target point. The movement of cuckoos with reference to this point in view of the distance interim among their present or current positions, terminal point and diversion angle. The algorithm is terminated if the population size exceeds maximum number of iterations.

### D. Level Set Segmentation for Region Segmentation

The Level-set methods are a conceptual framework for segmenting specific area using level sets as a tool for statistical analysis of surfaces and shapes. Initially the area is defined, then it evolved itself by values of pixels [21], it moves towards the object if the pixel value is same else it is not. So it is well suitable for analyzing the defected area. This method handles topological variation of surfacing boundary. The implicit representation is employed for contour representing zero level set $R$, that is, $\emptyset(R)=0$ of the level set function $\emptyset$. As the $\emptyset$ (a, b, t) the level set function is moved, it inflates, plunge, widen, etc. The level set functions is tend to be moved at first defining a velocity field **S** which illustrates by what means the contour points move within time. Next construct a primary value for the level set function, $\emptyset$(a, b, t=0), stationed on the initial contour position and alter $\emptyset$ with time, the current contour is described as $\emptyset$(a(t), b(t), t) = 0. In two dimensions the level set method amounts to describing a closed curve $\Gamma$ using a subordinate function$\emptyset$. $\Gamma$ is denoted as the zero level set that corresponds to the actual position of the curve at a given frame of $\emptyset$ by $\Gamma$ = {(a, b) |$\emptyset$(a, b) = 0}, moreover through the function $\emptyset$, employs $\Gamma$ implicitly through utilizing level set function. The positive values inside the region are taken by the function $\emptyset$ [22] delimited with the curve $\Gamma$ and outside the region it takes negative values. The evolution equation for level set is obtained by manipulating $\emptyset$ to indirectly move R as $\emptyset$(R) = 0

$$\frac{d\emptyset(R)}{dt} = \frac{\partial R}{\partial t}.\nabla\emptyset + \frac{\partial\emptyset}{\partial t} = 0 \qquad (17)$$

$$\frac{\partial\emptyset}{\partial t} = -S|\nabla\emptyset| \qquad (18)$$

where, S is the speed function normal to the curve. It is said that the distance function inside the curve is negative and outside it is positive. A suitable speed function S is chosen and may segment an object in an image. The function is accelerated using standard level set segmentation [23][24] and can be represented as:

$$S = 1 - \varepsilon M + \beta(\nabla\emptyset.\nabla|\nabla I|) \qquad (19)$$

The inflation inside the object is caused due to term 1 in the contour. The curvature of the contour is reduced by $-\varepsilon M$ (viscosity) term. The final term is edge attraction which pulls the contour to the edges. The equation can be updated as follows:

$$\emptyset_i^{n+1} = \emptyset_i^n - \Delta t\left(\sqrt{A_i^{-x}} + \sqrt{A_i^{+x}}\right) \qquad (20)$$

where, $\sqrt{A_i^{-x}}$ and $\sqrt{A_i^{+x}}$ are forward and backward differences, $\emptyset_i^{n+1}$- current (target) level set function point of curve, $\emptyset_i^n$ – initial level set point, $\Delta t$ - difference in time.

### E. Post Processing by Morphological Operation

The final step for defect detection is post processing by morphological operation. Here, the morphological operation [25] termed opening is used on the binary or grayscale image by way of the structuring element. There must be a distinct structuring element object, while divergent to group of objects. There upon for both the cases the similar structuring element and the morphological open operation are used followed by dilation. The morphological operator is used to refine the segmentation, it allows to add or delete the unwanted information in border of defects and smoothing the area of the defect, it provide the best refining result.



(a)       (b)

(c)       (d)

Fig. 2. Represents the Output Segmented Images of Sample 1using, (a) Adaptive thresholding with Cuckoo, (b) Adaptive thresholding with PSO, (c) Median based Otsu, (d) Otsu's Method.

Fig. 3. Represents the Output Segmented Images of Sample 2 using, (e) Adaptive thresholding with Cuckoo, (f) Adaptive thresholding with PSO, (g) Median based Otsu, (h) Otsu's Method.

The above figures denote the processed output images of the two samples using the adaptive thresholding with cuckoo optimization, adaptive thresholding with PSO, Median based Otsu and Otsu's Method. Fig. 2 and Fig. 3 represent the output segmented images of sample 1 and sample 2, respectively. The white portion of the image shows the defected area or uncoated regions and black portion represents coated area. The proposed method shows better detection compared to existing methods.

## IV. RESULT AND DISCUSSION

This section comprises the proposed technique evaluation based on experimental results. The achievement of the proposed Adaptive thresholding with Cuckoo optimization is compared with the results obtained from adaptive thresholding with PSO, median based Otsu and Otsu's method. In this regard, we used four performance measures such as sensitivity, specificity, accuracy and precision.

### A. Parameter Calculation

*a) Sensitivity:* It measures the proportion of actual positives that are suitably recognized. It is also referred as true positive rate or probability of detection.

$$Sensitivity = \frac{TP}{TP+FN} \qquad (21)$$

*b) Specificity:* It measures the proportion of actual negatives that are suitably recognized. It is also referred as true negative rate.

$$Specificity = \frac{TN}{TN+FP} \qquad (22)$$

*c) Accuracy:* Accuracy is defined as closeness of a measured value to a generally known or standard value.

$$Accuracy = \frac{TN+TP}{(TN+TP+FN+FP)} \qquad (23)$$

The following tables and its corresponding graphical representation shows the proposed method is suitable for examining the coated and uncoated regions of metal surfaces. The results in Table I clearly showed that proposed method greatly outperformed other algorithms in terms of the sensitivity measure. Table II shows the performance evaluation of proposed method for the specificity measures and Table III shows the performance evaluation of accuracy measures. The precision measures are tabulated in Table IV. The table values shows that adaptive thresholding based cuckoo method provided more robust results than other algorithms on the majority of the considered test input images. Fig. 4 and Fig. 5 show the graphical representation of proposed method with other algorithms, obtained from the experimental result using sensitivity measures and accuracy measures respectively. It is clear from the table and graph, the results obtained from the proposed method gives better result compared to other three techniques such as adaptive thresholding with PSO, median based Otsu and Otsu's method for defect detection and segmentation. The proposed method yields much better results when considered the evaluation parameters.



Fig. 4. Graphical Representation of Proposed method with other Algorithms Obtained from the Experimental Result using Sensitivity Measures.

TABLE I. COMPARISON TABLE OF ADAPTIVE THRESHOLDING WITH CUCKOO METHOD, ADAPTIVE THRESHOLDING WITH PSO, MEDIAN BASED OTSU AND OTSU'S METHOD IN TERMS OF SENSITIVITY

| Image No | Adaptive thresholding with Cuckoo | Adaptive thresholding with PSO | Median based Otsu | Otsu's Method |
|---|---|---|---|---|
| 1 | 0.80156 | 0.77951 | 0.7376423 | 0.70417 |
| 2 | 0.79483 | 0.77461 | 0.7392594 | 0.71506 |
| 3 | 0.80269 | 0.7394 | 0.7269316 | 0.70475 |
| 4 | 0.79461 | 0.77319 | 0.7372715 | 0.72497 |
| 5 | 0.82909 | 0.80573 | 0.7554124 | 0.71522 |

TABLE II.    COMPARISON TABLE OF ADAPTIVE THRESHOLDING WITH CUCKOO METHOD, ADAPTIVE THRESHOLDING WITH PSO, MEDIAN BASED OTSU AND OTSU'S METHOD IN TERMS OF SPECIFICITY

| Image No | Adaptive thresholding with Cuckoo | Adaptive thresholding with PSO | Median based Otsu | Otsu's Method |
|---|---|---|---|---|
| 1 | 0.84616 | 0.77326 | 0.71958 | 0.70082 |
| 2 | 0.79163 | 0.76796 | 0.70047 | 0.70239 |
| 3 | 0.81188 | 0.78531 | 0.70016 | 0.70013 |
| 4 | 0.82231 | 0.76642 | 0.70038 | 0.70046 |
| 5 | 0.81617 | 0.80177 | 0.71063 | 0.73762 |

TABLE III.    PERFORMANCE EVALUATION OF ADAPTIVE THRESHOLDING WITH CUCKOO METHOD, ADAPTIVE THRESHOLDING WITH PSO, MEDIAN BASED OTSU AND OTSU'S METHOD IN TERMS OF ACCURACY

| Image No | Adaptive thresholding with Cuckoo | Adaptive thresholding with PSO | Median based Otsu | Otsu's Method |
|---|---|---|---|---|
| 1 | 0.84955 | 0.80179 | 0.72619 | 0.72699 |
| 2 | 0.81606 | 0.7984 | 0.71677 | 0.71799 |
| 3 | 0.80856 | 0.741 | 0.71657 | 0.71655 |
| 4 | 0.82895 | 0.79742 | 0.72671 | 0.72676 |
| 5 | 0.84133 | 0.81969 | 0.73324 | 0.7501 |

TABLE IV.    PERFORMANCE EVALUATION OF ADAPTIVE THRESHOLDING BASED CUCKOO METHOD WITH ADAPTIVE THRESHOLDING WITH PSO, MEDIAN BASED OTSU AND OTSU'S METHOD IN TERMS OF PRECISION

| Image No | Adaptive thresholding with Cuckoo | Adaptive thresholding with PSO | Median based Otsu | Otsu's Method |
|---|---|---|---|---|
| 1 | 0.90722 | 0.8604 | 0.75221 | 0.75254 |
| 2 | 0.89441 | 0.8289 | 0.75244 | 0.75295 |
| 3 | 0.85472 | 0.83667 | 0.75236 | 0.75235 |
| 4 | 0.89941 | 0.84846 | 0.75242 | 0.75244 |
| 5 | 0.90032 | 0.83824 | 0.75508 | 0.76181 |



Fig. 5.    Graphical Representation of Proposed Method with other Algorithms Obtained from the Experimental Result using Accuracy Measures.

## V.    CONCLUSION

This paper has presented comparison of the three existing techniques with the proposed method through calibrating its performance by means of the evaluation parameters. The Cuckoo optimization based adaptive thresholding method proposed is verified to be more appropriate for solving the problems, compared with conventional methods. Furthermore, the experimental results convey that the Cuckoo optimization is considered to be improved throughout the time taken or instance for generating the optimal solution. Hence, the proposed method exhibits superior performance for detecting the surface defects while comparing with the other three techniques.

In future, we plan to use other powerful optimization method like firefly optimization to the proposed method and then apply it to different types of metallic surface images.

REFERENCES

[1]    Hao Shen, Shuxiao Li, Duoyu Gu and Hongxing Chang, "Bearing defect inspection based on machine vision", Measurement, vol.45, pp. 719–733, Jan 2012.

[2]    Hui Zhang , Xiating Jin, Q. M. Jonathan Wu, Yaonan Wang, Zhendong He and Yimin Yang, "Automatic Visual Detection System of Railway Surface Defects With Curvature Filter and Improved Gaussian Mixture Model", IEEE Transactions on Instrumentation and Measurement, vol.67, Issue.7, pp.1593–1608, 2018.

[3]    Eissa Negahdarzadeh, Effat Yahaghi, Behrouz Rokrok, Amir Movafeghi and Abolfazl Keshavarz Khani, "Diagnosis of design and defects in radiography of ceramic antiqueobjects using the wavelet-domain hidden Markov models", Journal of Cultural Heritage, vol. 35, pp. 56-63, Feb 2018.

[4]    Kun Liu, Heying Wang, Haiyong Chen, ErqingQu, Ying Tian and Hexu Sun, "Steel Surface Defect Detection Using a New Haar–Weibull-Variance Model in Unsupervised Manner", IEEE Transactions on Instrumentation and Measurement, vol.66, Issue.10, pp.2585–2596, July 2017.

[5]    RuoxuRen, Terence Hung and Kay Chen Tan, "Automatic Microstructure Defect Detection of Ti-6Al-4V Titanium Alloy by Regions-based Graph", IEEE Transactions on Emerging Topics in Computational Intelligence, vol.1, Issue: 2, pp.87–96, April 2017.

[6]    Hongwei Huang, Yan Sun, Yadong Xueand Fei Wang, "Inspection equipment study for subway tunnel defects by grey-scale image processing", Advanced Engineering Informatics, vol.32, pp. 188–201, March 2017.

[7]    Bashar M. Haddad, Sen Yang, Lina J. Karam, Jieping Ye, Nital S. Patel and Martin W. Braun, "Multifeature, Sparse-Based Approach for Defects Detection and Classification in Semiconductor Units", IEEE Transactions on Automation Science and Engineering, vol.15, pp. 145–159, Jan. 2018.

[8]    Payel Roy, Goutami Dey, "Adaptive Thresholding: A comparative study", IEEE International Conference on Control, Instrumentation, Communication and Computational Technologies, 2014.

[9]    Wenhua Han, Jun Xu, Mengchu Zhou, GuiyunTian, Xiaohui Shen and Sui-Hoi Edwin Hou, "Cuckoo-search and Particle-filter-based Inversing Approach to Estimating Defects via Magnetic Flux Leakage Signals", IEEE Transactions on Magnetics, vol.52, Issue.4, April 2016.

[10]    N. Otsu, "A threshold selection method from gray-level histograms," Automatica, vol. 11, no. 285–296, pp. 23–27, 1975.

[11]    J. H. Xue and D. M. Titterington, "Median-based image thresholding," Image Vis. Comput., vol. 29, no. 9, pp. 631–637, 2011.

[12]    Yasir Aslam , Santhi N , Ramasamy N , K. Ramar, "An Effective Surface Defect Detection Method Using Adaptive Thresholding Fused With PSO Algorithm", International Journal of. Simulation. Systems, Science & Technology,Volume 19, Number 6, Page 1, December 2018.

[13] Kai Hu, Xieping Gao and Fei Li, "Detection of Suspicious Lesions by Adaptive Thresholding Based on Multi resolution Analysis in Mammograms", IEEE Transactions on Instrumentation and Measurement, VOL. 60, NO. 2, Feb 2011.

[14] Andrey Kuznetsov and Vladislav Myasnikov, "A new copy-move forgery detection algorithm using image preprocessing procedure", Procedia Engineering, 3rd International Conference "Information Technology and Nanotechnology , vol.201, pp.436–444, 2017.

[15] EnginAkar,Sadık Kara, Hidayet Akdemirand AdemKırıs, "Fractal analysis of MR images in patients with chiari malformation: the importance of preprocessing", Biomedical Signal Processing and Control, vol.31, pp. 63–70 , 2017.

[16] Shadi Mahmoudi and Shahriar Lotfi, "Modified cuckoo optimization algorithm (MCOA) to solve graph coloring problem", Applied Soft Computing, vol.33, pp.48-64, Aug 2015.

[17] MahyaAmeryan and Mohammad Reza Akbarzadeh, "Clustering Based on Cuckoo OptimizationAlgorithm",Iranian Conference on Intelligent Systems (ICIS), 2014

[18] Sankalap Aroraand Satvir Singh, "A Conceptual Comparison of Firefly Algorithm ,Bat Algorithm and Cuckoo Search", International Conference on Control, Computing, Communication and Materials (ICCCCM) , 2013

[19] Seyed Alireza Moezi, Ehsan Zakeri and Amin Zare, "Structural single and multiple crack detection in cantilever beams using a hybrid Cuckoo-Nelder-Mead optimization Method", Mechanical Systems and Signal Processing, vol. 99, pp. 805–831, 2018.

[20] Seyed Alireza Moezi, Ehsan Zakeri and Amin Zare, "A generally modified cuckoo optimization algorithm for crack detection in cantilever Euler-Bernoulli beams", Precision Engineering, vol.52, pp. 227-241, Dec 2017.

[21] Cheng Liu, Weibin Liu and Weiwei Xing, "An improved edge-based level set method combining local regional fitting information for noisy image segmentation", Signal Processing, vol. 130, pp. 12–21, June 2016.

[22] Qingqi Hong andBeizhan Wang, "Segmentation of Vessel Images using a Localized Hybrid Levelset Method", 6th International Congress on Image and Signal Processing(CISP 2013), vol.2, pp. 631–635, Dec2013.

[23] Alizadeh. M, Soltanian Zadeh. H and Haji Maghsoudi. O, "Segmentation of Small Bowel Tumors in Wireless Capsule Endoscopy Using Level Set Method", IEEE 27th International Symposium on Computer-Based Medical Systems, 2014.

[24] Xin-Jiang, Renjie-Zhang and Shengdong-Nie, "Image Segmentation Based on Level Set Method", International Conference on Medical Physics and Biomedical Engineering, Physics Procedia, vol.33 pp.840–845, 2012.

[25] Priyadarsan Parida and Nilamani Bhoi, "2-D Gabor filter based transition region extraction and morphological operation for image segmentation", Computers & Electrical Engineering, vol.62, pp. 119-134, Aug 2017.

# Secure Data Provenance in Internet of Things based Networks by Outsourcing Attribute based Signatures and using Bloom Filters

Muhammad Shoaib Siddiqui[1], Atiqur Rahman[2], Adnan Nadeem[3], Ali M. Alzahrani[4]

Faculty of Computer and Information Systems
Islamic University of Madinah
Kingdom of Saudi Arabia

*Abstract*—With the dawn of autonomous organization and network and service management, the integration of existing networks with Internet of Things (IoT) based networks is becoming a reality. With minimal human interaction, the security of IoT data moving through the network becomes prone to attacks. IoT networks require a secure provenance mechanism, which is efficient and lightweight because of the scarce computing and storage resources at the IoT nodes. In this paper, we have proposed a secure mechanism to sign and authenticate provenance messages using Ciphertext-Policy Attribute Based Encryption (CP-ABE) based signatures. The proposed technique uses Bloom filters to reduce storage requirements and an outsourced ABE mechanism to use lessen the computational requirements at the IoT devices. The proposed technique helps in reducing the storage requirements and computation time in IoT devices. The performance of the proposed mechanism is evaluated and the results show that the proposed solution is best suited for resourced constrained IoT network.

*Keywords—Data provenance; bloom filter; ciphertext policy attribute based encryption; IoT*

## I. INTRODUCTION

The reason behind widespread adoption of ubiquitous computing is the use of smart devices and the Internet of Things (IoT). There are now more than 20 billion smartphones and IoT devices [1]. With the existing Wi-Fi and 3G/4G connections and 5G in the near future, Internet is available as on-the-go gateway to the rest of the world, making IoT devices an important factor in most solutions. IoT applications, such as, telemedicine, healthcare, sensing and diagnostic reporting, remote doctor consultation, security, surveillance, industrial automation and monitoring, telemetry, asset tracking, etc., with IoT-based sensor networks providing remote monitoring, are making significant progress. According to Statista 2018 [1], there would be more than 75 billion IoT devices connected to the Internet by 2025 (see Fig. 1).

One of the major characteristics behind the success of IoT is the least human involvement. With the penetration of smart-devices and intelligent home-appliances, IoT has opened the doors of new prospects. However, to maintain the accuracy of this technology, reliability and integrity of data must be ensured. The data received by these devices is vulnerable to

attacks. The most critical attacks in the field of monitoring and surveillance are false data injection and the corruption of data [3]. For ensuring the data integrity and sender's authentication, researchers have proposed data provenance mechanisms, which can identify a change in the data and provide a list of sender and forwarder of the data [6].

However, the data provenance mechanisms are also prone to attacks in which the attacker manipulates the provenance data to spoof its identity or hides the change in the data [4]. By attacking the provenance mechanism, the adversary can hide the loss of data integrity. Therefore, it is important to develop secure data provenance mechanism, which can provide data integrity and ensure that the mechanism itself is safe from attacks. However, the secure data provenance schemes have high cost of storage and computation for resource constrained IoT devices.

As the solution to the above-mentioned problem, we propose a secure provenance mechanism based on the attribute-based encryption [14]. We used ciphertext-policy attribute-based encryption (CP-ABE) [14] instead of Public key Cryptography [17] for designing our secure provenance mechanism. In [2], the authors have discussed how the algorithm of CP-ABE can be distributed. We used the algorithm, with some modifications, to convert our secure provenance mechanism into a distributed mechanism by outsourcing most of the computational load on fog/edge nodes. Due to the scarce resources at IoT nodes, we needed to devise a solution which required low storage for keeping the record of each data packet. For this, we have utilized the Bloom filters.



Fig. 1. The Number of IoT Devices Connected to the Internet.

The rest of the paper is articulated as follows. Section II presents the literature review by discussing the related work in the field of secure provenance. Section III presents the proposed mechanism. Section IV discusses the details of experimentation and results, while Section V concludes the paper.

## II. RELATED WORK

Data provenance provides the source of information and the historical path it has been forwarded from. Scientist is interested in the origin and historical transformation of the data as it contains critical information about the reliability of the data [5]. From it, one can determine the quality of the data based on data origin and past derivations, trace the sources of errors, automated re-creation of the derivations to update the data and provide the attribution of the data sources. Provenance is also essential for the commercial domain, where it can be used to deepen the data source in a data warehouse, trace the creation of intellectual property and provide an audit trail for regulatory purposes.

The use of a data source in a distributed system has been used to track records using a data source, in identifying the data flow from a subset of the original inputs, and in mending the data flow. To do so, one needs to keep track of the set of inputs for each operator used to derive each output. Although there are many ways for data provenance, such as, copy-provenance and how-provenance [4] [6], the information needed is a simple form of why-provenance, or lineage, as defined by Cui et al. [7].

In IoT based networks, reliability of the data received is a critical issue as a lot of attacks are made to corrupt data or inject false data in the networks. Data provenance is a mechanism to ensure that the data is forwarded by reliable nodes such that data integrity is maintained; however, the data provenance mechanisms are also prone to attacks. By attacking the provenance mechanism, the attacker can hide the loss of data integrity. Therefore, it is important to develop secure data provenance mechanism, which can ensure data integrity and is safe from attacks.

In [8], authors have argued that the one of requirement for IoT networks is to allow a user to trust the data regarding its origin and location. They provide a secure provenance mechanism based on secret key cryptography; however, they didn't mention the key distribution and revocation mechanisms.

In [3], the authors have identified the requirements and challenges of implementing a secure provenance system for IoT based networks. They claim that if data traces of IoT devices are recorded then provenance can play a vital role as it solves many issues related to data trustworthiness, decision-making, data reconciliation and data replication. The challenges of implementing a secure provenance mechanism in IoT, WSN, and RFID based networks are identified as data storage and processing, as these devices lack computational power.

Sultan et al. have presented a secure provenance scheme for wireless sensor network (WSN) in [9]. They argued that data provenance in sensor networks introduce some difficult requirements due to their scarce energy resources, and provenance mechanism's high bandwidth consumption. Therefore, it needs efficient storage and secure transport. Their scheme is based on Bloom filters, like [10] and is a light-weight solution; however, the security of the mechanism has a few question marks. Furthermore, their solution may not give the complete provenance path.

Kamal and Tariq provide a solution to provenance in IoT based network, which is also light-weight [11]. They have generated link fingerprints using the Received Signal Strength Indicator (RSSI) value of the IoT devices. The path to source node is identified at server, using the correlated coefficient based on the matching of the link fingerprints. However, they have relied on the physical parameters of RSSI for ensuring security also, without utilizing an authorization mechanism. The attacker can easily spoof the RSSI values to hide its location.

Liang et al. present a block-chain based provenance mechanism in cloud environment, in which they have proposed a framework for gathering and authenticate data provenance in the cloud, by using blockchain of provenance data [12]. It operates mainly in three phases: (1) provenance data collection, (2) provenance data storage, and (3) provenance data validation.

In our proposed schemes, we aim to develop a secure and light-weight provenance mechanism, specifically design for IoT devices, which are a part of a Machine to Machine Network. Our scheme uses Bloom filter to store the provenance information to reduce the storage requirement significantly at the resource constrained IoT devices. We have utilized the solution provided by [13] to outsource the signature generation using attribute-based encryption to lessen the load of signing and verification from the IoT nodes also. Following section discuss the proposed mechanism in detail.

## III. PROPOSED SYSTEM

### A. Communication Model

The communication model is a Cloud environment with IoT devices connected to the cloud. For the users of the cloud, some fog/edge nodes would be assigned that would be closed to the IoT devices as compared to the main servers of the Cloud. These dedicated edge nodes are computationally powerful and close to the IoT devices; therefore, IoT devices can delegate their signature generation and verification responsibilities to these nodes. Fig. 2 shows the scenario of the proposed system.

### B. Threat Model

The threat model for the proposed system is as follows:

- There is no physical protection for the IoT nodes. An attacker can access the IoT device and initiate a physical attack. The enemy's purpose is to access the memory of the device and compromise the data integrity.

- An attacker can mimic as an authorized IoT node. In this way attacker can insert wrong or fallacious data

into the system for compromising the accuracy of the system.

- An attacker can snoop, alter, reiterate, and infuse invalid data and messages.

- An attacker can modify data sent from an IoT node to the sink and compromise the provenance mechanism.

### C. Secure Provenance based on CP-ABE

For securing the provenance mechanism, Ciphertext-Policy Attribute based Encryption (CP-ABE) is used [14] [15]. Fig. 3 shows the conceptual model of implementing the CP-ABE in IoT based Cloud environment. The distribution of CP-ABE algorithm is adopted from [13], where the authors have implemented the load sharing of cryptographical computation

by outsourcing the signing and signature verification to the edge nodes.

By using the distributed CP-ABE algorithm, the partial signature is created by the Fog/Edge node, which takes up most of the computational overhead. The partial signature is sent to the IoT node which creates the complete signature.

Whenever an IoT node wants to send a message to another node, it creates the signature for the message and attaches it with the message. By using a signature, it is ensured that the message is from the authorized user and the data integrity is maintained. Using our proposed mechanism, we ensure the correctness of the data, which includes the identity of the sender, the identity of the forwarding node, and the original data of the IoT node.



Fig. 2. A Conceptual Scenario of the Cloud based IoT Network with Edge/Fog Node Supporting the Computational Load.



Fig. 3. Conceptual Model – IoT based Cloud/Fog Computing Scenario.

Fig. 4 shows the signing process used to send the provenance data from the IoT node. The signing process is mainly divided into five steps. These five steps are; Setup, Key Generation, Signing-outsourced and Signing. Signing-outsourced has most of the computational load that is why it is performed by the fog/edge node, while the step with significantly lower overhead, i.e. Signing is performed by the IoT node. The details of each step are as follows:

*1) Setup:* this step is the first step which is performed at the attribute authority, IoT node, or the edge node. The inputs to this step are the security parameter δ, a universal set of attributes μ, and an auxiliary information α. The setup process provides the master key Ψ and a public key U.

*2) Key generation:* This is the second step, which is to be performed by the attribute authority. The IoT node would initiate the process of signing by requesting the attribute authority for secret key for the IoT node with attribute set $A$, and the corresponding edge/fog node. The key generation phase is provided the master key $\Psi$, and attribute set $A$, for which it outputs the secret key $P$ for the IoT node and $\Theta$ for the edge /fog node, which would use it to create the partial signature.

*3) Sign$_{outsourced}$:* This phase is executed by the edge/fog node, which takes input the secret key $\Theta$ attribute set $A$, and the predicate $\Gamma$. This is the main part of the algorithm in which the CP-ABE is used to create the partial signature $\sigma_{partial}$. This partial signature is sent to the IoT node, which would calculate the final signature; however, the computational overhead is bared by the edge/fog node.

*4) Sign:* the final signing phase that is performed by the IoT node outputs the signature for the message $M$. It uses the partial signature $\sigma_{partial}$, the secret key $P$ (provided by the attribute authority), and the predicate $\Gamma$.

After the message is signed by the IoT node, the message $M$ and the signature $\sigma$ are forwarded to the next node. The message $M$ includes the identity of the sender, the identity of the forwarding node, and the original data of the IoT node, while signature $\sigma$ is attached to verify if there was a change in the message or if the message is forwarded by an adversary node.

To verify the integrity of the message and authenticate the sender, the verification step is performed, which is also shown in Fig. 5.

*5) Verify:* The verify phase takes input the message M, the signature $\sigma$, the public key $U$, and the predicate $\Gamma$. It validates if the signature is correct or invalid.

If the verification process fails then the secure provenance mechanism can identify the break in provenance chain. It would ask the sender node to resent the data. If the verification step is successful then the data about the message/packet is saved at the receiving node before that data is forwarded to the next node. In this way the messages between the two hops/nodes are ensured to be secure and verified. Before the data/message is forwarded to the next hop/node, the data is stored at the current node with provenance information for each data packet/message.

As there would be a lot of message that passes from an intermediate node, therefore, it would take up a lot of storage and if these intermediate nodes are IoT nodes then they do not have the luxury of high storage capacity. For solving the issue of storage our mechanism incorporates Bloom filters as discussed in the next subsection.

### D. Reducing Storage Requirment using Bloom Filters

The provenance framework stores information of the data whenever it passes from a node. We have utilized Bloom-filters [16] to store the information of the data at each node. Bloom filter does not take a lot of storage; rather it uses an array of bits to store the data. First the secured data packet with sender's digital signature is passed through the hash functions. Three hash functions provide indexes to the 32-bit Bloom filter array, where the bits are turned on (1). Similarly, all the data packets are stored in the Bloom filter. Fig. 6 shows the storage mechanism where i=32 corresponds to the IP address of the sender.



Fig. 4. The Signing Process using Outsourced Attribute based Encryption.



Fig. 5. The Signature Verification Process using Outsourced Attribute based Encryption.

Fig. 6.    The Logging of Provenance Data at the IoT Node.

Fig. 7 shows how the membership is checked for a data packet. When the provenance tracking is performed, a request for provenance checking is send to the neighboring nodes. To check if the data packet was passed through a node, the data packet is again passed through the hash functions. The hash function would provide the indexes and if any index is found to be set to (1), then we would know that the data packet was passed through the node. The 32-bit value at that index would provide the IP address of the sender of this packet. If any of the indexes is found to be 'not-set' (0) then the packet was not passed from this node. Now, the provenance tracking request is forwarded to that node which performs the similar action to identify the node that forwarded the said packet to this node. Following this scheme, the provenance path is constructed.



Fig. 7.    The Query Mechanism for Checking the Membership of the said Data Packet.

## IV.  SIMULATION AND EVALUATION

The main objective of this research was to reduce the computational load from the IoT nodes and reduce the storage requirements for storing provenance information. We have performed simulation for our proposed scheme using a server as a cloud, an Nvidia Jetson tk2 device as an edge node and multiple instances of a Java application on a laptop to simulate as IoT nodes. We have generated random packets between the

IoT nodes, which have utilized the service from the edge node to generate partial signature by the edge /fog nodes and the rest of the signature (with less computational load) by the IoT nodes. We have calculated the storage requirement and computational load at each IoT node to evaluate our proposed scheme.

Fig. 8 shows the storage measurements based on the number of IoT nodes and increased number of packets. We have selected the Bloom filter array with size m = 1024. With the increase in membership, the probability for false positive in the Bloom filter also increase, hence, after 1024/4 = 256 packets, the Bloom filter array is archived, and a new array is allocated to the Bloom filter. As each entry is of 32 bits, therefor, the size of each Bloom filter array is 1024 x 32 bits = 4KB.

Fig. 9 shows the computational load measurement of the proposed scheme for performing key generation and signature with the support of the edge node and without the support of edge node. This computational load also affects the battery power of an IoT device; hence, the battery drainage issue can also be resolve using outsourced signature technique.



Fig. 8.    Required Storage in KBs with Increased Number of Packets for different Number of Nodes.



Fig. 9.    Time Required to Generate Key and Signature Computation with Increased Number of Packets with and without Edge Node Calculating Partial Digital Signatures.

## V. Conclusion

In this paper, we have proposed a secure provenance tracking mechanism which uses the outsourced signing technique using Attribute Based Encryption by offloading the IoT node. The IoT node only calculates the partial digital signature, while the high-overhead computations are performed by the edge node. Our scheme also reduces the storage requirement at the IoT devices by storing the provenance data in a multi-array Bloom filter. We have discussed in detail the mechanism of the storage techniques using Bloom filter, which uses hash function to save the information about the packets received and forwarded by the node. The hash-based searching has the complexity of $O(1)$; therefore, it decreases the provenance tracking time. The results show that our proposed scheme uses less storage and takes less time in terms of computation as compared to the normal secure provenance mechanism.

In future, we aim to enhance the proposed scheme by implementing it on Graphics Processing Units (GPUs) available on mobile devices by using parallel computing, which will reduce the cost in terms of time . We are also investigating for a possible solution using Blockchain technologies, which has inherent characteristics to sim in data provenance.

## Acknowledgment

## References

[1] Joyia, Gulraiz J., Rao M. Liaqat, Aftab Farooq, and Saad Rehman. "Internet of Medical Things (IOMT): applications, benefits and future challenges in healthcare domain." J Commun 12, no. 4 (2017): 240-7.

[2] Chen, Xiaofeng, "Secure Outsourced Attribute-Based Signatures", IEEE transactions on parallel and distributed systems, 2014, Volume: 25 Issue: 12 Page: 3285-3294

[3] Sabah Suhail, Zuhaib Uddin Ahmad, Choong Seon Hong, "Introducing Secure Provenance in IoT: Requirements and Challenges" International Workshop on Secure Internet of Things (SIoT 2016), Sep. 26-30, 2016, Heraklion, Crete, Greece

[4] Peter Buneman, Sanjeev Khanna, and Wang Chiew Tan. Data provenance: Some basic issues. In Proceedings of the 20th Conference on Foundations of Software Technology and Theoretical Computer Science, FST TCS 2000, pages 87–93, London, UK, UK, 2000. Springer-Verlag

[5] Pasquier, Thomas; Lau, Matthew K.; Trisovic, Ana; Boose, Emery R.; Couturier, Ben; Crosas, Mercè; Ellison, Aaron M.; Gibson, Valerie; Jones, Chris R.; Seltzer, Margo (5 September 2017). "If these data could talk". Scientific Data. 4: 170114. doi:10.1038/sdata.2017.114.

[6] Robert Ikeda and Jennifer Widom. Data lineage: A survey. Technical report, Stanford University, 2009.

[7] Y. Cui and J. Widom. Lineage tracing for general data warehouse transformations. VLDB Journal, 12(1), 2003.

[8] Muhammad Naveed Aman, Kee Chaing Chua, Biplab Sikdar, Secure Data Provenance for the Internet of Things, Proceedings of the 3rd ACM International Workshop on IoT Privacy, Trust, and Security, April 02-02, 2017, Abu Dhabi, United Arab Emirates

[9] S. Sultana, G. Ghinita, E. Bertino, M. Shehab, "A lightweight secure provenance scheme for wireless sensor networks", Proc. Int. Conf. Parallel Distrib. Syst. - ICPADS, pp. 101-108, 2012.

[10] Muhammad Shoaib Siddiqui, Syed Obaid Amin, and Choong Seon Hong, "Hop-by-hop Traceback in Wireless Sensor Networks", IEEE Communication Letters, Vol. 16, No.2, pp.242-245, February 2012

[11] Mohsin Kamal and Muhammad Tariq, "Light-Weight Security and Data Provenance for Multi-Hop Internet of Things", in IEEE Access, Volume: 6, 2018.

[12] ProvChain: A Blockchain-based Data Provenance Architecture in Cloud Environment with Enhanced Privacy and Availability, in the Proceedings of the 17th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing, Madrid, Spain. May 2017, Pages 468-477

[13] Asim, Muhammad, Milan Petkovic, and Tanya Ignatenko. "Attribute-based encryption with encryption and decryption outsourcing." (2014).

[14] Bethencourt, J.; Sahai, A.; Waters, B. (2007-05-01). "Ciphertext-Policy Attribute-Based Encryption". 2007 IEEE Symposium on Security and Privacy (SP '07): 321–334. doi:10.1109/SP.2007.11.

[15] M. Ambrosin, A. Anzanpour, M. Conti, T. Dargahi, S. R. Moosavi, A. Rahmani, P. Liljeberg, "On the Feasibility of Attribute-Based Encryption on Internet of Things Devices", IEEE Micro Magazine, vol. 36, no. 6, pp. 25-35, December 2016.

[16] A. Broder, M. Mitzenmacher, "Network applications of bloom filters: a survey", Internet Mathematics, vol. 1, no. 4, pp. 485-509, 2004.

[17] Stallings, William (1990-05-03). Cryptography and Network Security: Principles and Practice. Prentice Hall. p. 165. ISBN 9780138690175.

# Query Expansion based on Explicit-Relevant Feedback and Synonyms for English Quran Translation Information Retrieval

Nuhu Yusuf[1], Mohd Amin Mohd Yunus[2], Norfaradilla Wahid[3]

Faculty of Computer Science and Information Technology, Universiti Tun Hussein Onn Malaysia (UTHM), Parit Raja, Malaysia[1]
Management & Information Technology Department, Abubakar Tafawa Balewa University (ATBU), Bauchi, Nigeria[1]
Faculty of Computer Science and Information Technology, Universiti Tun Hussein Onn Malaysia (UTHM), Parit Raja, Malaysia[2]
Faculty of Computer Science and Information Technology, Universiti Tun Hussein Onn Malaysia (UTHM), Parit Raja, Malaysia[3]

*Abstract*—Search engines are commonly present as information retrieval applications that help to retrieve relevant information from different domain areas. The crucial part of improving the quality of search engine is based on query expansion, which expands the query with additional information to match additional important documents. This paper presents a query expansion approach that utilizes explicit relevant feedback with word synonyms and semantic relatedness. We describe the possibility and demonstrations based on the experimental work pertain to search engines where relevant judgment and word synonyms can improve search quality. In order to show the level of improving the proposed approach, we compared the results obtained from the experiments based on Yusuf Ali, Arberry and Sarwar Quran datasets. The proposed approach shows improvement over other methods.

*Keywords*—*Query expansion; search engine; relevant feedbacks; explicit relevant feedback; synonyms; information retrieval*

## I. INTRODUCTION

Search engines are one of the most successful information retrieval systems that are proposed in order to address information overload and allow users to find relevant information using search queries. Search engines can be examined as an improvement of the query reformulation that provides quality of search results and performance.

Search engines extend to wider areas of usage including desktop [1], federated [2], enterprise [3] and mobile [4][5] to improve the performance of search and to emphasize on the relevancy of the information obtained based on users queries [3]. Although, web search engine related search are the most popular and widely-used among researchers.

Recently, search engines in the form of cross-lingual [6] topics have become widely available for Quran verses retrieval. In contrast to other cross-lingual search engines, the Quran verse search engine contains more different translations within different languages. Yet these search engines performance results are still not encouraging. The main challenge of the Quran search engine for retrieving relevant search results is that the user queries are not sufficient enough to retrieve relevant Quran verses.

The search engines process that generally provides the highest search performance is query expansion [7]. Query expansion is based on the user query assessment about the quality of the search result and expanding such query to retrieve relevance of the results. The major research issue in query expansion area focuses on improving poor precision result values and term selections. Different methods have been proposed to address these issues based on relevant feedbacks. Basically, there are three relevant feedbacks approaches explicit, implicit and pseudo-relevant feedbacks methods.

Explicit relevant feedback [7] required the use of an expert, also known as Assessor in the field Quran translation to judge the relevance of the results retrieved. The relevance results specify either the Quran verse retrieved is relevant or irrelevant to the query. Implicit relevant feedback [8] focuses on user behavior while searching for a document. These could be whether a user selects and view a document or not. If a user views a document, assumed it is relevant and if not means irrelevant. Also among the widely used relevant feedback is the pseudo-relevant feedback method [9] [10] which assumed that the top retrieved documents are relevant to the user queries. Query expansion methods based on relevant feedbacks proof effective in providing relevant search results. However, there is a need to combine these feedbacks, especially explicit feedback with word synonyms to improve the performance of Quran search engines. Existing research papers have focused to improve the search engine performance. Among the papers is Jakub et al. [11] which used relevance assessment files obtained from TREC to expand. Rashid [12] present a current paper on how a query expansion method can improve search performance for Urdu language using relevance assessment. Lucchese et al. [13] focus on how to improve search engine performance by selecting an effective and efficient term. Lavrenko et al. [14] present the need to estimate the relevance model in order to obtain word synonyms.

Afzal and Muktar [15] suggested a quran English WordNet as a solution to short queries, especially when linking to semantic similarities. Their work achieved significant improvement. Although semantic similarity prove effective, Moawad, Alromima and Elgohary [16] stated that absence of semantic resources were identified in many languages and as such alternative semantic approach need to be develop.

Bentrcia, Zidat and Marir [17] examine the possibility of using semantic relatedness instead of only semantic similarities. This paper used only semantic relatedness to obtain relationship between two concepts using "AND" conjunctive. However, Lashkari, Bagheri and Ghorbani [18] identified based on the literatures that depending only on semantic resource wouldn't improve better search performance. Integrations of other methods with semantic was suggested to yield better results. This paper utilizes explicit-relevant feedback with combined WordNet and semantic relatedness to expand the query.

The remainder of this paper is organized as follows: Section 2 will provide a review of related work on query expansion methods based on relevant feedback and how these methods improve search performance; Section 3 describes the proposed approach using explicit relevant feedback and synonyms; Section 4 present the experiments conducted, and finally, Section 5 present the conclusion and future work.

## II. Related Work

### A. Query Expansion Methods

Many query expansion techniques are now available for testing the performance of search engines. According to Azad & Deepak [7], query expansion techniques can be categorized into fourteen (14) techniques. These techniques can be in either global or local analysis query expansion. However, in practice, only a few successfully improve search engine results. The most commons ones found effective are WordNet, thesaurus, explicit feedback, implicit feedback, pseudo-relevance feedback and Wikipedia. Each of the techniques can be applied to different scenario depending on the performance of the search engine you want to improve. Explicit, implicit and pseudo-relevance feedbacks are important in improving the performance of search results using relevance feedbacks from users. Others also play major roles in various search improvements.

### B. Query Expansion Methods based on Relevance Feedbacks

Many query expansion techniques have been used over the years. The explicit feedback technique has shown effective in retrieving the quality results from the search engine. In explicit feedback, a query will be sent to the search engine. The expert will be asked to judge the relevance of the results obtained from that query. Also, an expert must have prior knowledge of the domain are considered. The results that expert judged most relevant will be given to a system to compute the performance search engine. Liu et al [19] proposed a method that combines both explicit and implicit feedbacks into consideration. They emphasize addressing the challenges that may come from the top-N recommendation. Jiang, He and Allan [20] suggested the use of multiple explicit feedbacks for improving the relevant judgment in subsequent information retrieval evaluations. They used both the implicit feedback during user sessions and explicit feedback using at least four criteria. Their evaluations showed that explicit feedback has a relationship with the user experience. Mach et al. [21] evaluate whether including expert judgment can yield positive results in climate change assessment. They invited experts from different domains for their inputs on the issues. Their result shows that using expert judgment can transparently and consistently improve the

quality of the needed results. Lester et al. [22] developed a framework incorporating expert judgment to minimizes inconsistency in the provided results. Their concern is to evaluate analogue structure-activity relationships (SAS). However, the bias in expert judgment [23] may reduce the quality of the results. This is because of the inappropriate use of experts from other domains. Wilson [24] presented the use of multiple experts in the evaluation in order to bias. The study considered from within and outside the organization. The results indicate multiple experts can have more relationship than within. According to Hasanain [25] present how the system can use automatic ranking without considering relevance expert judgment. But, Alvarado-Valencia [26] argue that including expert judgment can ensure system credibility. Verma et al. [27] use various relevant judgment collected from desktop and mobile experts. The result indicates desktop can many nonrelevant documents in addition to few relevance ones.

Implicit feedback technique that deduces what the users' intent to do based on their observed behaviours. Understanding user behaviour can significantly improve the performance of search engine [28]. The implicit feedback is a good technique for measuring similarities [8] of many documents based on user search behaviour. The search behavior may be from online purchase history [29][30][31], browsing history [29][32], search patterns [31][33][34], or even mouse movements [35][34]. All these prove effective in infer intention for a particular search engine. For instance, Kawasaki and Hasuike [29] show browsing history can be utilized to provide a recommendation to users in electronic commerce sites. Ghosh, Rath and Shah [33] present how the complexity of the web search engine can be improved using users search behaviour logs. Xie et al., [34] stated that search behaviour pattern can reveal users intention within the search environment. These search behaviour pattern can see mostly in terms of click time or query reformulations. Kwok [35] also found that mouse click can easily predict the intention of users in searching needed information. Recently, Zhang et al., [36] investigated how implicit feedback collected user search behaviour can improve the quality of search results. They manipulated the relevant documents obtained from the top rank documents list. However, their investigation does not consider the semantic relationship between various search behaviours which can provide effective search quality.

Pseudo-Relevance Feedback, which others refer to as blind relevance feedback, is another query expansion technique that can improve search performance. Pseudo-relevance feedback can address mismatch vocabulary challenges [10]. It can improve search results without users' involvement. Montazeralghaem, Zamani and Shakery [37] reveal that Pseudo-relevance feedback that utilizes a number of top-rank documents will improve search effectiveness. Albishre et al., [9] proposed a pseudo-relevance feedback model to address microblog documents mismatching issues. They considered how a query can be modified in order to improve users' performance. Furthermore, Na and Kim [38] suggest that the performance of pseudo-relevance feedback model can only be reasonably achieved when the document length have been normalized. They expanded the query by adding a few

additional terms to increase the size of the query. That means the size of the query should be neither not too short nor too high. Expanding query with multiple terms [39] can improve the performance of search engines than using single or few terms. They proposed a technique that can combine multiple terms using fuzzy logic. Also, Bayesian [40] with many retrieval models can effectively improve search performance. They considered using models that can allow users to search for information using a query and optimize the results. Furthermore, Singh and Sharan [41] proposed a model that considered multiple selection terms in order to obtain relevant documents. Their model utilizes explicit-relevance feedback in retrieving a number of relevant documents.

## III. PROPOSED APPROACH

Similar to other query expansion, our proposed method utilizes three traditional ranking algorithms. However, after computing search performance, the algorithm that provides better results was used to expand our query. In this case, our proposed query expansion method utilizes the BM25 algorithm. Furthermore, we improve the search results of BM25 with synonyms and we called the new algorithm QEES. Firstly, based on the idea of vector space model, the cosine similarity between term X and Y was computed using the equation (1).

$$cos(X,Y) = \frac{X \cdot Y}{|X||Y|} \qquad (1)$$

The synonyms of QEES can be computed using equation (2):

$$f_{qx} = z_{cd} \cdot \log \frac{z \cdot z_{cd}}{z_c \cdot z_d} \qquad (2)$$

Where $f_{qx}$ represent the modified query based on the similarities of words in WordNet and the original query. The $z_{cd}$ represents a number of times the word in $c$ and $d$ have seen together and $z_c$, $z_d$ represent the number of times that seen each of the word $c$ and $d$. We assume that $c$ and $d$ represent our original query and the WordNet respectively. We also assume that $c = \{c_1, c_2, ..., c_n\}$ and $d = \{d_1, d_2, ..., d_3\}$, where both $\{c_1, c_2, ..., c_n\}$ and $\{d_1, d_2, ..., d_3\}$ represent various words.

Furthermore, we compute the BM25 ranking with synonyms for an expanded query on a document using equation (3) below:

$$Score(D, Q_{exp}) = tf \cdot \frac{f(f_{qx}, D) \cdot (k_1)}{f(f_{qx}, D) + k_1 \cdot \left(1 - b + b \frac{|D|}{avgdl}\right)} \qquad (3)$$

Where $tf$ can also be calculated in equation (4):

$$\log \frac{N - nf_{qx} + 0.5}{nf_{qx} + 0.5} \qquad (4)$$

D represents a document in a corpus and the $Q_{exp}$ represent the expanded query. The total number of the document is represented as N while the number of documents containing the appended word synonyms. Moreover, the appended word synonyms term frequency in document D is represented as $f(f_{qx}, D)$. The length of document D in words and the average length in the text collection are denoted as |D| and avgdl

respectively. The free parameters for search optimization are represented by $k_1$ and $b$.

## IV. EXPERIMENTS

### A. Datasets and Queries

We conduct our experiment on Yusuf Ali, Sarwar and Arberry English version datasets. These three datasets are collected from Tanzil [42]. Yusuf Ali is an English Quran translation dataset by Islamic scholar Abdullahi Yusuf Ali. The scholar was born in India and wrote many Islamic books during his lifetime. This dataset gains acceptance by many English speaking countries. The second dataset is Arberry which is also another English version of Quran translation dataset by Arthur John Arberry. Arthur John Arberry is a non-Muslim scholar but helps to translate the Quran into English. Finally, the Sarwar is third Quran translated dataset used by Shaykh Muhammad Sarwar. He is an Islamic scholar from Pakistan who publishes some Islamic books including the Quran English translation.

It is not necessarily that every quran translation dataset be accepted to other sects. Therefore, as a source of our dataset, we select these datasets because Tanzil [42] is free from any sect or country bias.

Each dataset is given as a single document and contained 6236 verses represented by lines. Therefore, we transformed each verse into a single document which we got 6236 documents are available in each of the datasets. We saved each document with its chapter and verse for easy representation. For instance, 003045 represent chapter 003 verse 045. This mean, the first 3 figures represent the chapter while the last 3 represent the verse.

We believed the search results can only be obtained if the document contains a query keyword. Moreover, we understand that not every query provide better search results based on the studies conducted by Yusuf et al [43]. Therefore, we adopted 36 queries used in Ahmad et al [44] which combine different types of queries. Furthermore, our proposed approach used an English WordNet provided by Princeton University [45]. Version 2.1 of the WordNet has been used to target synonyms words.

### B. Evaluation Metrics

Most of the research papers used precision, recall and mean average precision [1][2] to evaluate quran information retrieval. However, the proposed QEES used Average precision, average recall, mean reciprocal rank and mean average precision metrics for evaluation.

The position at which each ranked document provides the relevant document will determine average precision as in equation (5):

$$AP = \frac{1}{n} \sum_{i-1}^{n} precision(P_i) \qquad (5)$$

Average recall utilizes the position at which each ranked document provides the retrieved document as computed in equation (6):

$$AR = \frac{1}{n} \sum_{i-1}^{n} recall(R_i) \qquad (6)$$

Equation (7) presents a mean average precision which uses rankings from different users' queries and then averages them to obtain average precision.

$$MAP = \frac{1}{Q}\sum_{Z=1}^{Q}\frac{1}{D_Z}\sum_{i=1}^{D_Z}P(d_i) \qquad (7)$$

The mean reciprocal rank computes the reciprocal rank average above queries sets and computed in equation (8):

$$MRR = \frac{1}{N}\sum_{q=1}^{N}\frac{1}{rank_q} \qquad (8)$$

### C. Benchmark

We utilize the explicit-relevant feedback provided by Quran translation experts in Ahmed et al [44]. This document serves as a benchmark for our proposed method. As we obtain the results from the search engine algorithm, they would be compared with the benchmark to separate relevant documents retrieved from the irrelevant ones. Both relevant and irrelevant retrieve will be used to compute search performance.

### D. Comparisons

In this experiment, the proposed QEES method will be compared against different traditional ranking algorithms. Specifically, we will compare our proposed method against original BM25, TF-IDF and Lucene algorithms.

### E. Performance

In this experiment, the proposed QEES method will be compared against different traditional ranking algorithms. Specifically, we will compare our proposed method against original BM25, TF-IDF and Lucene algorithms.

Table III tabulates the average precision and recalls obtained from the three datasets based on the 36 queries used. From the results, we can notice that our proposed QEES perform best on Yusuf Ali dataset.

In terms of the results obtained, Fig. 1 and 2 represent mean average precision (MAP) and mean reciprocal rank (MRR) the results obtained with Yusuf Ali, Sarwar and Arberry datasets.

Fig. 1 shows the MAP in Yusuf Ali dataset by making use of BM25, tf-idf, Lucene and proposed QEES. The proposed QEES perform better in retrieving relevance results. It achieves

**17.75%** improvements as compared to BM25 with 16.23%, tf-idf with 11.19% and Lucene with 5.39%.

Fig. 2 shows the MRR in Yusuf Ali dataset achieved with BM25, tf-idf, Lucene and expanded synonyms. The order of the probability of documents is important, as it is the basis for ranking first correct answers in ranked documents. As the relevant documents are retrieved, only the rank of the relevant documents is considered, other relevant answers are all ignored. When comparing the results, the search engine has **10.58%** performances when applying proposed QEES which is less than the TFIDF with 11.01%.

In summary, the proposed QEES using Yusuf Ali dataset proves effective on the MAP. The performance of various search engine methods proved that relevant results can be obtained for a particular search query.

For Arberry dataset, the results obtained in Fig. 1 have some similarities with the Sarwar datasets on the MAP. However, the BM25 and the proposed QEES slightly have little differences. The proposed QEES is inferior in retrieving relevant results. It improves search performance by achieving 9.28% as compared to the Lucene method with 16.43% and BM25 with 10.07%.

Fig. 2 also shows the MRR on Arberry dataset where 10.04% has been achieved on BM25 as compared to the proposed QEES with **7.53%** results.

The Sarwar dataset results in Fig. 1 show the search performance on four search methods. Although a proposed QEES is significant on MAP as compared to BM25 traditional method, the MRR result with BM25 is higher than the proposed QEES. Such results can be applied to practical situations. Moreover, Fig. 2 shows the BM25 method has 9.47% results as compared to TFIDF with 6.39%, Lucene with 7.95% and proposed QEES with **9.07%**.

Interestingly, as shown in Tables I and II, the results obtained in Yusuf Ali dataset using BM25 and proposed QEES perform significantly on the MAP. Arberry and Sarwar datasets perform better on Lucene using MAP. In terms of MRR, Yusuf Ali performs best with TFIDF. Sarwar and Arberry perform best on BM25 and Lucene respectively on MRR.



Fig. 1.    MAP Search Performance Results on the Three Datasets.

Fig. 2.    MRR Search Performance Results on the Three Datasets.

TABLE I.        MAP Search Performance Results on the Three Datasets

| Datasets | Methods | MAP |
|---|---|---|
| Yusuf Ali | BM25 | 16.23 |
| | TFIDF | 11.19 |
| | Lucene | 5.39 |
| | Proposed Method | **17.75** |
| Sarwar | BM25 | 13.25 |
| | TFIDF | 7.72 |
| | Lucene | 10.76 |
| | Proposed Method | **16.36** |
| Arberry | BM25 | 10.07 |
| | TFIDF | 6.14 |
| | Lucene | 16.43 |
| | Proposed QEES | **9.28** |

TABLE II.        MRR Search Performance Results on the Three Datasets

| Datasets | Methods | MRR |
|---|---|---|
| *Yusuf Ali* | BM25 | 9.91 |
| | TFIDF | 11.01 |
| | Lucene | 5.32 |
| | Proposed Method | **10.58** |
| *Sarwar* | BM25 | 9.47 |
| | TFIDF | 6.39 |
| | Lucene | 7.95 |
| | Proposed Method | **9.07** |
| *Arberry* | BM25 | 10.04 |
| | TFIDF | 7.16 |
| | Lucene | 12.37 |
| | Proposed QEES | **7.53** |

TABLE III.    AVERAGE PRECISION AND RECALL SEARCH PERFORMANCE RESULTS ON THE THREE DATASETS

| Queries | Yusuf Ali | | Sarwar | | Arberry | |
|---|---|---|---|---|---|---|
| | Average Recall | Average Precision | Average Recall | Average Precision | Average Recall | Average Precision |
| 1 | 0.0790 | 0.2111 | 0 | 0 | 0 | 0 |
| 2 | 0.0203 | 0.225 | 0 | 0 | 0 | 0 |
| 3 | 0.0962 | 0.0958 | 0.0962 | 0.3740 | 0.0577 | 0.3333 |
| 4 | 4.5 | 0.3303 | 2.1667 | 0.1561 | 2.5 | 0.2327 |
| 5 | 0.0698 | 0.5828 | 0.0582 | 0.0606 | 0.0233 | 0.0196 |
| 6 | 0.15 | 0.0685 | 0.2 | 0.4306 | 0.2 | 0.0707 |
| 7 | 0.1923 | 0.0859 | 0.1154 | 0.0635 | 0.1154 | 0.1192 |
| 8 | 0.5625 | 0.2173 | 0.4375 | 0.0605 | 0.25 | 0.0484 |
| 9 | 0.1 | 0.0204 | 0 | 0 | 0 | 0 |
| 10 | 0.25 | 0.0981 | 0.1667 | 0.0189 | 0.1667 | 0.0833 |
| 11 | 0.1539 | 0.1780 | 0.0513 | 0.0257 | 0.0641 | 0.0415 |
| 12 | 0 | 0 | 0.0036 | 0.3444 | 0.0030 | 0.0973 |
| 13 | 0 | 0 | 0.125 | 0.5 | 0 | 0 |
| 14 | 0.3333 | 0.1111 | 0.3333 | 1 | 0.5 | 0.0387 |
| 15 | 0.0686 | 0.2937 | 0.0882 | 0.1304 | 0.0588 | 0.2891 |
| 16 | 0.25 | 0.0898 | 0 | 0 | 0.125 | 0.1429 |
| 17 | 0.2143 | 0.0512 | 0.2857 | 0.0379 | 0.2857 | 0.0306 |
| 18 | 0.2778 | 0.0718 | 0.5556 | 0.1472 | 0.4722 | 0.2316 |
| 19 | 0.0233 | 0.0313 | 0.0233 | 0.0110 | 0.0233 | 0.0149 |
| 20 | 0.0714 | 0.0128 | 0 | 0 | 0 | 0 |
| 21 | 0.1177 | 0.2993 | 0.0294 | 0.0746 | 0.0098 | 0.0196 |
| 22 | 0.1818 | 0.1396 | 0.2273 | 0.3506 | 0.1818 | 0.3145 |
| 23 | 0.0714 | 0.3807 | 0.0325 | 0.1358 | 0.0195 | 0.0161 |
| 24 | 0.0645 | 0.5255 | 0.0645 | 0.4556 | 0.0484 | 0.075 |
| 25 | 1 | 0.1235 | 2 | 0.0331 | 1 | 0.0141 |
| 26 | 0.0063 | 0.4119 | 0.0038 | 0.3026 | 0.0025 | 0.2 |
| 27 | 0.0133 | 0.1122 | 0.0067 | 0.0266 | 0.0089 | 0.0796 |
| 28 | 0.0027 | 0.0281 | 0 | 0 | 0 | 0 |
| 29 | 0.0088 | 0.0159 | 0 | 0 | 0 | 0 |
| 30 | 0.0389 | 0.4817 | 0.0111 | 0.1 | 0 | 0 |
| 31 | 0.1765 | 0.3279 | 0.2353 | 0.5410 | 0.1471 | 0.0812 |
| 32 | 0.0068 | 0.0259 | 0.0170 | 0.1144 | 0.0068 | 0.4064 |
| 33 | 0.3125 | 0.4943 | 0.1875 | 0.2576 | 0 | 0 |
| 34 | 0.0325 | 0.0864 | 0.0649 | 0.1213 | 0.1234 | 0.3223 |
| 35 | 0.0417 | 0.0553 | 0.0278 | 0.0175 | 0.0278 | 0.0164 |
| 36 | 0.25 | 0.1064 | 0 | 0 | 0 | 0 |
| Average | **9.7379** | **6.3894** | 7.6142 | 5.8912 | 6.4211 | 3.3221 |

## V. Conclusion and Future Work

The performance of any search engine is a crucial part of the success of any search engine across various domains. Therefore, in this paper, we have presented a query expansion method that utilizes explicit relevant feedback and word synonyms for improving the performance of Quran web search engine. The improvement can be seen in terms of means average precision and mean reciprocal rank performance, average precision and recall.

The proposed QEES new approach to query expansion has many benefits; especially it can be used across different search engine algorithms to rank documents according to relevance. The proposed QEES achieved as per as 1.19723151 on the MAP which is a significant improvement when compared to other methods. Furthermore, the proposed approach can be able to retrieve reciprocal ranks of results for a given query.

For future work, our research experiments have shown how a word that is exactly or nearly the same use to improve search results. We observed that relevant documents have used different terms. In contrast, the search engine must have a distributed representation of term with semantics metadata so that meaning of a word can be better processed using machine learning algorithms such as neural network and therefore a beneficial to improve the performance of query expansion for better results.

## Acknowledgment

## References

[1] M. S. Raje, "Analysis of Desktop Search and Ranking of Their Results Based on Semantics from User Feedback," 2016 12th Int. Conf. Signal-Image Technol. Internet-Based Syst., pp. 241–245, 2016.

[2] M. Dragoni, A. Rexha, H. Ziak, and R. Kern, "A semantic federated search engine for domain-specific document retrieval," Proc. ACM Symp. Appl. Comput., vol. Part F1280, pp. 303–308, 2017.

[3] U. Kruschwitz and C. Hull, "Searching the Enterprise," Found. Trends® Inf. Retr., vol. 11, no. 1, pp. 1–142, 2017.

[4] C. Luo, Y. Liu, T. Sakai, F. Zhang, M. Zhang, and S. Ma, "Evaluating Mobile Search with Height-Biased Gain," Proc. 40th Int. ACM SIGIR Conf. Res. Dev. Inf. Retr. - SIGIR '17, pp. 435–444, 2017.

[5] J. Mao, Y. Liu, N. Kando, Z. He, and M. Zhang, "A Two-Stage Model for User ' s Examination Behavior in Mobile Search," Assoc. Comput. Mach., pp. 273–276, 2018.

[6] G. Chandra and S. K. Dwivedi, "Query expansion based on term selection for Hindi - English cross lingual IR," J. King Saud Univ. - Comput. Inf. Sci., 2017.

[7] H. K. Azad and A. Deepak, Query Expansion Techniques for Information Retrieval: a Survey. 2017.

[8] M. Liu, W. Pan, M. Liu, Y. Chen, X. Peng, and Z. Ming, "Mixed similarity learning for recommendation with implicit feedback," Knowledge-Based Syst., vol. 119, pp. 178–185, 2017.

[9] K. Albishre, Y. Li, and Y. Xu, "Effective pseudo-relevance for Microblog retrieval," Proc. Australas. Comput. Sci. Week Multiconference - ACSW '17, pp. 1–6, 2017.

[10] R. Cummins, "Improved Query-Topic Models Using Pseudo-Relevant Pólya Document Models," Proc. ACM SIGIR Int. Conf. Theory Inf. Retr. - ICTIR '17, pp. 101–108, 2017.

[11] Jakub Dutkiewicz and C. Jędrzejek, "Calculating Optimal Queries from the Query Relevance File," in Proceedings of the 11th International Conference MISSI, 2018, pp. 249–259.

[12] I. Rasheed and H. Banka, "Query Expansion in Information Retrieval for Urdu Language," 2018 Fourth Int. Conf. Inf. Retr. Knowl. Manag., pp. 1–6, 2018.

[13] C. Lucchese, F. M. Nardini, and R. Trani, "Efficient and Effective Query Expansion for Web Search," pp. 1551–1554, 2018.

[14] V. Lavrenko and W. B. Croft, "Relevance based language models," Proc. 24th Annu. Int. ACM SIGIR Conf. Res. Dev. Inf. Retr. - SIGIR '01, pp. 120–127, 2001.

[15] H. Afzal and T. Mukhtar, "Semantically Enhanced Concept Search of the Holy Quran: Qur'anic English WordNet," Arab. J. Sci. Eng., 2019.

[16] I. Moawad, W. Alromima, and R. Elgohary, "Bi-Gram Term Collocations-based Query Expansion Approach for Improving Arabic Information Retrieval," Arab. J. Sci. Eng., vol. 43, no. 12, pp. 7705–7718, 2018.

[17] R. Bentrcia, S. Zidat, and F. Marir, "Extracting Semantic Relations from the holy Quran Based on Arabic Conjunctive Patterns," Comput. Intell., vol. 30, pp. 382–390, 2016.

[18] F. Lashkari, E. Bagheri, and A. A. Ghorbani, "Neural embedding-based indices for semantic search," Inf. Process. Manag., vol. 56, no. 3, pp. 733–755, 2019.

[19] N. N. Liu, E. W. Xiang, and M. Zhao, "Unifying Explicit and Implicit Feedback for Collaborative Filtering," in IEEE 2nd International Conference on Big Data Analysis, 2017, pp. 1445–1448.

[20] J. Jiang, D. He, and J. Allan, "Comparing In Situ and Multidimensional Relevance Judgments," Proc. 40th Int. ACM SIGIR Conf. Res. Dev. Inf. Retr. - SIGIR '17, pp. 405–414, 2017.

[21] K. J. Mach, M. D. Mastrandrea, P. T. Freeman, and C. B. Field, "Unleashing expert judgment in assessment," Glob. Environ. Chang., vol. 44, pp. 1–14, 2017.

[22] C. Lester, A. Reis, M. Laufersweiler, S. Wu, and K. Blackburn, "Structure activity relationship (SAR) toxicological assessments: The role of expert judgment," Regul. Toxicol. Pharmacol., vol. 92, no. January, pp. 390–406, 2018.

[23] J. R. Axt, H. Nguyen, and B. A. Nosek, "The Judgment Bias Task: A reliable flexible method for assessing individual differences in social judgment biases.," J. Exp. Soc. Psychol., no. February, pp. 1–19, 2018.

[24] K. J. Wilson, "An investigation of dependence in expert judgement studies with multiple experts," Int. J. Forecast., vol. 33, no. 1, pp. 325–336, 2017.

[25] M. Hasanain, "Automatic Ranking of Information Retrieval Systems," in The 11th ACM International Conference on Web Search and Data Mining, 2018.

[26] J. Alvarado-Valencia, L. H. Barrero, D. Önkal, and J. T. Dennerlein, "Expertise, credibility of system forecasts and integration methods in judgmental demand forecasting," Int. J. Forecast., vol. 33, no. 1, pp. 298–313, 2017.

[27] N. Craswell, "Study of Relevance and Effort across Devices," in Proceedings of the 2018 Conference on Human Information Interaction&Retrieval, 2018, pp. 309–312.

[28] A. H. Awadallah, "LearnIR : WSDM 2018 Workshop on Learning from User Interactions," in Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining, 2018, pp. 797–798.

[29] M. Kawasaki, "A Recommendation System by Collaborative Filtering Including Information and Characteristics on Users and Items," in Computational Intelligence (SSCI), 2017 IEEE Symposium Series, 2017, pp. 1–8.

[30] M. Zhou, "Micro Behaviors : A New Perspective in E-commerce Recommender Systems," in Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining, 2018, pp. 727–735.

[31] Y. Wen, P. Yeh, T. Tsai, W. Peng, and H. Shuai, "Customer Purchase Behavior Prediction from Payment Datasets," in Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining, 2018, pp. 628–636.

[32] U. Gadiraju, R. Yu, S. Dietze, and P. Holtz, "Analyzing Knowledge Gain of Users in Informational Search Sessions on the Web," ACM SIGIR Conf. Hum. Inf. Interact. Retr. (CHIIR), March 11-15, 2018 New Brunswick, New Jersey, USA, pp. 2–11, 2018.

[33] S. Ghosh, M. Rath, C. Shah, and H. Street, "Searching as Learning : Exploring Search Behavior and Learning Outcomes in Learning-Related Tasks," in Proceedings of the 2018 Conference on Human Information Interaction&Retrieval, 2018, pp. 22–31.

[34] X. Xie and M. Zhang, "Why People Search for Images using Web Search Engines," in Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining, 2018, pp. 655–663.

[35] T. C. K. Kwok, E. Y. Fu, E. Y. Wu, M. X. Huang, G. Ngai, and H. V. Leong, "Every Little Movement Has a Meaning of Its Own: Using Past Mouse Movements to Predict the Next Interaction," in IUI '18 Proceedings of the 23th International Conference on Intelligent User Interfaces, pp. 397–401.

[36] H. Zhang, M. Abualsaud, and M. D. Smucker, "A Study of Immediate Requery Behavior in Search," in Proceedings of the 2018 Conference on Human Information Interaction&Retrieval - CHIIR '18, 2018, pp. 181–190.

[37] A. Montazeralghaem, H. Zamani, and A. Shakery, "Term Proximity Constraints for Pseudo-Relevance Feedback," Proc. 40th Int. ACM SIGIR Conf. Res. Dev. Inf. Retr. - SIGIR '17, pp. 1085–1088, 2017.

[38] S. H. Na and K. Kim, "Verbosity normalized pseudo-relevance feedback in information retrieval," Inf. Process. Manag., vol. 54, no. 2, pp. 219–239, 2018.

[39] J. Singh et al., "Fuzzy Logic Hybrid Model with Semantic Filtering Approach for Pseudo Relevance Feedback-based Query Expansion," in Computational Intelligence (SSCI), 2017 IEEE Symposium Series, 2017.

[40] D. Li and E. Kanoulas, "Bayesian Optimization for Optimizing Retrieval Systems," in Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining, 2018, pp. 360–368.

[41] J. Singh and A. Sharan, "Rank fusion and semantic genetic notion based automatic query expansion model," Swarm Evol. Comput., vol. 38, no. September 2017, pp. 295–308, 2018.

[42] Z.-Z. Hamid, "Quran Translations." Tanzil, 2007.

[43] N. Yusuf, M. A. M. Yunus, and N. Wahid, "A Comparative Analysis of Web Search Query: Informational Vs Navigational Queries," Int. J. Adv. Sci. Eng. Inf. Technol., vol. 9, no. 1, 2019.

[44] F. D. Ahmad, "A Malay language document retrieval system: An experimental approach and analysis," Universiti Kebangsaan Malaysia, 1995.

[45] P. U. About WordNet, "WordNet." Princeton University, 2010.

# Modelling and Implementation of Proactive Risk Management in e-Learning Projects: A Step Towards Enhancing Quality of e-Learning

Haneen Hijazi[1]
Faculty of Information Technology
Hashemite University
Zarqa, Jordan

Bashar Hammad[2]
Department of Mechanical and
Maintenance Engineering
German Jordanian University
Amman, Jordan

Ahmad Al-Khasawneh[3]
Faculty of Information Technology
Hashemite University
Zarqa, Jordan

*Abstract*—**The introduction of e-Learning to higher education institutions has been evolving drastically. However, the quality of e-Learning becomes a central issue in order to provide all stakeholders with the necessary confidence to compete with traditional learning methods. Risk management plays a vital role in the successful implementation of e-Learning projects and in attaining high-quality e-Learning courses. Little research has been conducted about implementing risk management in e-Learning projects. This work proposes a quality assurance framework for e-Learning projects. This framework comprises a proactive risk management model that integrates risk management into the e-Learning process. This integration helps in obtaining high-quality e-Learning courses by preventing negative e-Learning risks from being materialized. The model is verified to evaluate its effectiveness through a Renewable Energy Course that was converted from a traditional face-to-face into e-Learning course. Quantitative and qualitative measures are performed to analyze the data collected through the implementation of the project. The results show that the proposed model is managed to mitigate the majority of probable risk factors leading to high-quality e-Courses development and delivery.**

*Keywords*—*e-Learning; technology-enhanced learning; quality; proactive risk management; risk factors; higher educatio*

## I. INTRODUCTION

The existence of many Higher Education Institutions (HEIs) today could be ascribed to their abilities to keep pace with the continuous technological changes. Some institutions feel overwhelmed by these changes. Some otherwise consider them as inevitable dimension to strengthen their competitive advantages. The major revolution in technology has been the evolution of Information and Communication Technology (ICT) in the recent decades. Indeed, the evolution is not in the technology itself but rather in its applications in knowledge, information sharing, and education. One of the most effective ICT applications is what is now referred to as e-Learning.

In its simplest definitions, e-Learning means doing learning activities electronically through the Internet [1]. It is considered a key part of distance education [2]. Some definitions restrict e-Learning to the delivery of the e-Content over the Internet. Broader definitions widen the concept to cover the interaction among participants too, delivered by

different communication technologies, mainly the Internet [3]. This delivery could be fully online, or a hybrid approach that integrates electronic learning and traditional classrooms in what is so-called blended learning [4].

e-Learning has got through Higher Education (HE) drastically over the last few years. HEIs have recently recognized the importance of e-Learning in reducing operating cost and increasing students' satisfaction. Indeed, these issues are necessary but not sufficient for a university to achieve the desired competitive advantage. As quality has been playing an increasingly important role in the educational system [5], universities need to guarantee a high quality e-Learning to compete strongly.

Despite that e-Learning activities highly penetrate HE, quality of e-Learning has been an issue of debate. It was difficult to define what quality means to e-Learning courses. Several conceptual models and approaches rose recently, but the actual practice of quality of e-Learning in HEIs is still poor. Moreover, most of these approaches focus solely on the courses quality and the learning outcomes. Indeed, an overall detailed process-oriented quality assurance framework must exist and be followed during the e-Learning process to ensure the quality of the entire course, not only the output [6].

e-Learning projects deal with design, implementation, and utilization of social and information technological systems [7]. These systems involve several software applications (e.g. Learning Management System (LMS) and e-Content). However, there is no theoretical basis for project management that is specific to e-Learning [8]. Hence, most e-Learning projects management approaches follow software project management methodologies and inherit their characteristics because developing educational software shares several aspects with software development [9]. Mainly, it is the learning component what differentiates e-Learning projects from other types of projects [10]. Building upon this, the successful implementation of an e-Learning project requires balancing between project schedule, budget, and quality. Since the e-Learning process is mainly characterized by its quality, competitive universities should put the quality of the e-Learning projects at the forefront despite other challenges.

Both success factors and possible failures together should be taken into account to increase the probability of project success. These possible failures are called risks. An e-Learning project inherits all types of risk factors encountered in Information System projects beside many other risks that are specific to e-Learning [11]. These factors might negatively affect the quality of the e-Learning course being developed and/or the learning outcomes, project schedule and/or resources. Many risk factors are associated with e-Learning projects [12]. These factors can be personal or dispositional, learning style, instructional, situational, organizational, content suitability, and technological [13]. In order to minimize their negative impact, these factors must be managed carefully [14]. Successful management of e-Learning risk factors would improve the quality of the e-Learning process and, consequently, the competitive advantage of the institution. Hence, Risk management should be the core competence in e-Learning projects.

Risk Management is an important part in project management and very crucial for project's success. Indeed, this fact applies to all types of projects. Risk management involves predicting risks that might negatively affect the project schedule, budget or quality, and taking measures to avoid or mitigate the impacts arising from those risks [15]. Risk management in e-Learning projects could be defined as the set of principles, practices, procedures, methodologies, and tools aimed at identifying, analyzing and handling risk factors that could negatively affect the content development and delivery process and hinder the e-Learning project from achieving its desired outcomes. In e-Learning projects, risk management is a critical discipline that helps in reducing uncertainty, avoiding rework, improving content quality, making e-Learning process more reliable, decreasing learner's dissatisfaction and increasing the overall success chances. Several approaches to risk management exist. Reactive risk management does not apply mitigation strategies till the occurrence of the risks. Reactive risk management is expensive in terms of time and cost required to make necessary changes in the purpose of managing risk at the time of its occurrence. In contrast, proactive risk management provides information to stakeholders on how to best use resources to prevent the occurrence of unwanted events [11]. The latter aims at avoiding risks before they materialize; hence, it can be referred to as preventive risk management [16].

The best way to manage risks in e-Learning projects is to select the most suitable instructional design methodology and consider it during the development process as a mean to manage risks. Deciding upon the model that best fits a project is influenced by how risky the project is; the types of these risks and the degree to which each model supports risk management [17]. e-Learning projects are risky; they are vulnerable to several risks during the development and delivery phases. ADDIE (Analyze, Design, Develop, Implement and Evaluate) is the most popular well-known instructional design model. It is a prescriptive sequential instructional design model. Recently, opponents raise doubts around ADDIE model due to its strict linear implementation from the analysis phase to the evaluation [18]. Indeed, ADDIE e-Learning projects suffer from major risk factors that could not be handled using

pure implementation of ADDIE. In ADDIE model, no feedback from stakeholders is incorporated until the last phase of the project [19]. Hence, the major risk factor is the late change in requirements. Any late change in requirements would either require a large amount of rework which would cost extra time and money [19], or it may lead to an unsatisfied user (i.e. low-quality process). Either case, project failure is inevitable. Another risk factor is that overlapping is not allowed, in other words, practitioners cannot move to the next phase until the previous phase is completely finished. Moreover, using ADDIE, no deliverables are made available to learners until the last phases of the project when all deliverables are ready [19]. Clearly, these factors will negatively affect the project especially if the project suffers from time contention. A recent trend is to abandon ADDIE and to move towards Agile approaches [9, 20]. Agile is a lean approach to project management that enables building releasable yet good quality products in short time periods [19]. Agile is an iterative [21], team-based, collaborative approach. At the end of each iteration, a working deliverable is made available to users, feedback from the users is sought at the end of each iteration and changes, if exist, are incorporated in successive iterations [22]. Moreover, iterations may overlap. Clearly, agile model avoids the major risks of the ADDIE model.

Risk management has not been extensively performed in e-Learning projects [11]. In e-Learning, the project mainly passes through two phases; content development and learning delivery. Risks need to be managed carefully during both phases so that the project can achieve its expected outcomes. To achieve its outcomes, an e-Learning project should guarantees a high-quality content and a high- quality learning. Assuring this, the project surely leads to the major outcome of the learning process; high-quality student.

Implementing e-Learning in engineering education is challenging. Pedagogy, infrastructure, policy, strategy, quality, and management are the major challenges [23]. However, modern technologies have ushered in an era of change in engineering education [23]. In Engineering education, e-Learning involves the use of ICT to deliver virtual classrooms, conducting laboratory experiments, administering Virtual Learning Environment (VLE), developing professional e-Content that is rich in animations and visualizations that are used to demonstrate material, concepts, diagrams, processes, circuits, components and functioning [13]. However, an adequate application of e-Learning in engineering education would facilitate the learning process and lead to high quality learning.

Despite being a developing country, Jordan has made great strides in the fields of ICT but with a humbling experience in e-Learning in HEIs. This is due to several barriers including resistance, technological infrastructure hinders, quality assurance issues and slow change of learning structures and processes. Hashemite University (HU) is a public university in Jordan. The experience of HU in e-Learning is not recent. However, all previous HU practices in e-Learning were blended learning; none of the offered courses were carried out fully online. Recently, HU started to recognize the importance of the fully online courses in reducing operating cost and

increasing students' satisfaction especially with the economic challenges, resources constraints and the geographical location of the university [24]. With all these in mind, the quality of the outcomes of the fully online learning process has become a major priority of the university. Hence, the need to develop, deliver and evaluate a fully-online pilot course was a necessity. HU has started to move towards online courses and several e-Learning projects have been started. HU intended to follow a quality framework in order to assure the quality of the e-Learning process and the outcomes of the e-Courses.

However, the literature lacks such a process-oriented quality assurance framework. Hence, in this paper, we propose a process-oriented quality assurance framework to ensure the quality of an entire course. This framework is risk-oriented; it embeds a proactive risk management model for managing e-Learning projects. According to the knowledge of the authors, the literature has not discussed how proactive risk management in e-Learning projects would enhance the quality of e-Learning courses. To validate this framework, a project was started in 2014 to develop and deliver a pilot course in Renewable Energy (RE). The course was first offered online in summer 2015-2016 then continuous piloting and monitoring were maintained over a period of 6 semesters. RE course was selected because it is a hot topic in Engineering, relatively advanced technical course, and never been developed and delivered electronically in the region.

This study aims at enhancing the quality of e-Learning in higher education by modelling and implementing proactive risk management in e-Learning projects. The main research issues that this paper aims to investigate are:

- To describe how proactive risk management can be implemented in e-Learning projects to enhance their quality.

- To identify the major risk factors associated with managing the implementation of an e-Learning project.

- To devise the major risk mitigation strategies associated with each risk factor.

- To compile a case study of an online course to collect and extract information about the applicability of e-Learning at HU.

- To describe how can we measure the quality of an e-Learning course.

The rest of the paper is organized as follows: Section II reviews related work. Section III introduces the proposed framework, section IV introduces an implementation of the model through the RE course, section V displays, discusses and analyses the results of the RE course evaluation and presents some limitations of this work, and section VI concludes the work and suggests future work.

## II. RELATED WORK

Several current projects and research aim at enhancing quality of e-Learning in HEIs. In [25], Bralić and Divjak proposed a blended learning model that integrates Massive Open Online Course (MOOC) into a traditional classroom. Their model was based on learning outcomes and used to evaluate the effectiveness of integrating a MOOC with classroom-based teaching. A model was proposed by Casanova and Moreira in [26] for teachers in HE to reflect and discuss the quality of Technology-Enhanced Learning (TEL) in their blended learning programs. They argued that HEIs need to be more critical with regard to the use of TEL, and to support it as a counterpart to traditional learning. The experience of the North Carolina Central University for an Introductory Biology course over four terms was discussed in [27]. In this research, Hollowell, Brooks, and Anderson discussed the impact of the application of quality course design standards on the design and student outcomes. Atoum, Al-Zoubi, Abu Jaber, Al-Dmour and Hammad [28] presented a new approach for delivering e-Learning courses in Jordanian universities. The researchers introduced a national quality assurance system for TEL that aims at improving, developing and implementing accreditation standards for quality assurance of TEL courses and study programs at a national level. English language course was implemented and delivered at a national level as a pilot study according to a strict quality assurance framework. In their project [5], Mazohl and Makl introduced scientific description of a practicable quality framework for blended learning. The proposed framework focuses on the quality of courses, the course itself, the quality in the organizations delivering blended learning courses, the learners' needs and the environmental conditions. Based on this framework, a pilot course was developed and tested at the University of Helsinki, Finland.

In [29], Gómez-Rey, Barbera and Fernández-Navarro explored the quality of the online learning experience based on the Sloan-C framework and the Online Learning Consortium's (OLC) quality scorecard. The researchers found that the OLC index has ignored the opinions of the learners in evaluating quality of online programs. Hence, they proposed an alternative way of measuring the quality of online learning programs using teachers and students' perceptions and satisfaction. Misut and Pribilova [6] proposed and verified a quality assurance method of e-Learning – ELQ based on Kirkpatrick model which includes four levels of evaluation: reaction, learning, behavior, and results. Ghislandi, Raffaghelli, and Yang [30] introduced an approach that takes into account the participants' engagement as insiders of a quality learning culture. In [31], Bremer described how the AKUE model could be used to improve the quality of e-Learning and the eContent development. The AKUE model involves four phases: analysis, conception, implementation, and evaluation. In [32] Ossiannilsson and Landgren introduced a conceptual framework to enhance quality of e-Learning in HE based on experiences from three international benchmarking projects. The framework suggests that various aspects of accessibility, flexibility, interactivity, personalization, and productivity should be implemented at all levels in order to meet students' expectations. In [33] Lin and Chen reported that a successful e-Learning system should take both system and information quality into account. They combined Technology Acceptant Model (TAM) with Information system Success Model (ISM) by considering system quality, quality of platform information, and course information. Sarsa and Soler [34] studied the relations among the variables of e-Learning quality by means of five conceptual

maps that ease the visualization of these relations. Marshall in [35] summarized the outcomes of multiple international e-Learning Maturity Model (eMM) assessments which aimed at improving e-Learning quality in the organizations.

Very simple trials in the literature have discussed the role of risk management in e-Learning projects. Vesper, Kartog˘lu, Herrington and Reeves [11] employed two risk assessment strategies in the formative evaluation of a task-based e-Learning program developed by the World Health Organization (WHO). The first strategy used an expert reviewer and the second used a risk assessment expert facilitator. Both strategies aimed at identifying probable risks early and controlling them. Reference [36] examined the students' risk perception while using the aLF. The aLF is a LMS that was developed by Spanish National University of Distance Education (UNED), Vázquez-Cano and García found that risks are concentrated in two dimensions: "basic risks" and "own and beyond students' circumstances risks". Barik and Karforma [1] presented the risks that might face different e-Learning system stakeholders. They also suggested tools and techniques to minimize those risks. The identified risks and techniques were related to integrity, security and reliability of an e-Learning system. In [3], Mahmud and Gope discussed several technological, psychological, socio-cultural and economic factors that would affect successful implementation of e-Learning in HE in Bangladesh. They concluded by recommending measures to resolve these issues with government and the private sectors.

In [37], Surcel and Reeiu presented a series of the problems of designing and implementing an e-Learning strategy, objectives, planning and didactic process management. They classified risks into risks associated to the professors and risks associated to the students. Finally, they proposed general controls to manage these risks. Allen and Hardin in [38] presented a model of management that encompasses the Instructional System Design (ISD) process. They also presented a process for evaluating the risk factors of the project and how to manage changes throughout the project that may threaten the project's success. Andersson [39] identified 37 major challenges for e-Learning in developing countries. The work used data from the eBIT program in the University of Colombo, Sri Lanka. These factors were discussed and solutions were suggested. Angelou and Economides in [40] presented a real-option methodology for controlling risks in e-Learning infrastructure business field and choosing the optimum ICT investment's deployment strategy. In [12], Ifinedo investigated the risks associated with implementing an e-Learning information system project in Estonia. As a rank-order list of the typical project risk factors encountered in this project was produced.

Despite that HEIs have recently recognized the importance of e-Learning in engineering education; few attempts were found in the literature. Bandaya, Ahmed, and Jan [13] discussed the application of e-Learning in engineering education. The research investigates e-Learning practices in the

Engineering institutions of the state of Jammu and Kashmir as a case study. Rodríguez, Granados, and Muñoz [41] presented the intimate relationship between the e-Learning method and the studies of Engineering in Spain through teaching examples on several subjects of different Engineering studies. Benchicou, Aichouni and Nehari [42] reported the results of an empirical study that measures the readiness of HEIs in Algeria towards the application of the e-Learning in engineering education. An important barrier for implementing e-Learning in engineering education is the need for remote experimentations. Chandra and Samuel in [43] implemented a user-friendly system that allows students to carry out laboratory experiments from remote locations. Hence, despite the advent of e-Learning in all education fields, subtle improvements are required in the Engineering field.

In [44], Rooij and Williams stated that ADDIE is not enough for project management in instructional design and proposed research opportunities for closing the gap between instructional design education and practice.

## III. PROPOSED FRAMEWORK

Recent research discusses the influence of Software Engineering methodologies and practices over instructional design methodologies to provide high quality e-Learning [9, 20, 45, 46, 47]. In this proposed framework, two Software Engineering concepts were combined with the instructional design methodology. These two concepts are Risk Management and Agility.

In this research, the authors propose a framework for e-Learning projects that utilizes a proactive approach to risk management. Fig. 1 depicts an overview of the embedded model. This framework is risk-oriented wherein probable risks and risk factors are identified early and the whole development process is guided by the identified risks.

Typically, risk management process is integrated into the development and delivery process, and risks are avoided during the execution of the e-Learning process. In this framework, avoidance is imposed utilizing two aspects:

- First: adopting an instructional design model that best fits e-Learning projects

- Second: devising and implementing avoidance strategies that handle probable risk factors

In the proposed framework, a hybrid "Agiled-ADDIE" approach is used. Best practices from Agile is blended with ADDIE. e-Learning practitioners pass by all phases of ADDIE but with an order that is subject to continuous feedback from the different stakeholders. Using this approach, an effective collaboration and communication between all stakeholders, developers and learners is assured. This communication is the major constituent of Agile approach. The proposed framework consists of five stages: risk identification, planning, production, delivery and evaluation. These stages are:

- Risk Identification: In order to implement a proactive risk management framework, risk factors should be identified early before proceeding into the actual e-Learning development process. Hence, the first stage in this framework is risk identification stage. In this stage, the project manager sets a detailed list of risk factors that threaten the e-Learning project. This list is constructed based on project documentation, reviews of previous similar e-Learning projects, available checklists, and the project manager's experience. This initial list is refined later in the planning stage. Once the initial set of risk factors has been identified, they have to be managed. In this proactive framework, e-Learning practitioners proceed into the eContent development and delivery process phases, activities, risk mitigation strategies with an eye towards the identified risks and preventing them from being materialized [48].

- Planning: In this stage, the project is initialized; resources, risks and the course are planned. Team members are hired, tools are selected, budget and schedule are planned, sources of material are decided upon, and learners' analysis is carried on. The initial list of risk factors identified in the previous stage is refined here. Brainstorming sessions which involve all team members beside learners' analysis may come out with new important risk factors or may lessen the severity of any of the previously identified ones in the context of the project. Most importantly, a set of avoidance strategies is devised for each factor that mitigates the risk before being materialized. The devised strategies are practiced later in the proper phase of the project. Also, in this stage, the course objectives, outcomes, outlines and the assessment criteria are set. The organizational structure of the course and the basic unit of development (i.e. referred to as module) are also decided upon.

- Production: Developing eContent shares common aspects with software development, especially in the design and production stages [9]. Hence, Agility would be successful in this stage. The main goal of this method is to minimize the risk of incomplete or bad quality output. Using Agile, modules are produced iteratively. A module is produced at each iteration and enhanced in the successive iteration. Each module typically goes iteratively through the following phases:

- Design: In this phase, a set of learning objectives for each module besides the sequencing in which they should be achieved are formulated. The module outlines and the general look and feel of the module are created.

- Content preparation: In this phase, references and sources of information are selected, the material is collected and refined and the final content is written. The external support material may also be introduced.

- Storyboard design: A document is created in this phase that describes all elements of the final product including text elements, images, audio elements, animations, and interactions.

- Development: Interface layout and course outlines are created using the authoring tool. Media, interactive component, self-assessments, and quizzes are developed and then imported into the authoring tool.

- Publishing: The module is produced in a shareable format that can be handled by LMS in this phase.

- Review: The review in this phase is a formative assessment activity. Once the module is developed, and before it is delivered to the learners, a review should be conducted at the end of the iteration to ensure the quality of the module and continuous improvement. This review should mainly involve the Subject Matter Expert (SME) and the developer in order to make any required changes early in the successive iteration. Once the review has no negative feedback, then the module is ready to be delivered to the learner.

- Delivery: In this stage, the module is deployed into the LMS and made available to learners. This stage also includes managing and facilitating learners' activities such as virtual class-rooms, assignment, and quizzes. It is worth to mention that a module can be delivered even if others are not ready. This stage also includes quality assurance activities. During the delivery of each module, formative assessment is conducted. This assessment uses feedback from students during the learning process activities in order to evaluate the quality of the eContent, assess student's reception, improve weakness areas and strengthen the e-Learning course.

- Evaluation: Once all modules have been delivered and the course has finished, the course content and the instructional delivery should be evaluated. Feedback from students is used to perform the summative assessment. This assessment uses quantitative and qualitative analysis. Quantitative analysis uses students' satisfaction surveys to evaluate students' reaction towards the course and students' results and measure the knowledge they acquired through the course. Qualitative analysis uses interviews and open-ended questionnaires. The aim of the summative analysis at the end of the e-Learning course is to ensure that the course has achieved its expected outcomes and learning objectives and that the proposed risk management strategies have proven its effectiveness.

In real practice, the development stages described above remain applicable throughout all e-Content development projects. However, their differences emerge with regard to the anticipated project challenges and the proposed mitigation strategies.

Fig. 1.   The Proposed Proactive Risk Management Model for e-Learning Projects.

## IV. IMPLEMENTATION OF THE PROPOSED MODEL

In order to validate the effectiveness of the proposed risk-oriented quality assurance model, one pilot course in RE was designed and delivered at HU. Continuous piloting and monitoring of the course have been spanning over six semesters in purpose of examining the comprehension of students. The project started in early 2014 and was first delivered on summer 2015-2016. The intention of the project was to create an e-Content that could be used efficiently by HU students in the online RE course. The course included static content, media elements, animations and interactive components to help students understand course topics via activities such as e-Content packages, exercises, design problems, self-assessments, quizzes, online assignments, peer discussions and virtual classrooms. All of these activities were hosted at Moodle LMS.

The RE e-Course was initially created for the undergraduate Mechatronics Engineering students at HU. The sample of the study consisted of 186 undergraduate students distributed on six semesters as in Table I.

Using the risk management method described above, the project team members were able to systematically and proactively identify risks related to RE course and determine various ways to reduce their effects. Due to the risk-oriented nature of the framework, the team used the framework as dynamic rather than static approach. The process evolves as new risks arise and other risks disappear. The following subsections describe in details how the various stages of the model were implemented during the development and delivery of the RE course.

### A. Risk Identification Stage

In this stage, the project manager identified the major sources of risk (i.e. risk factors) that are specific to the RE course. An initial list of risk factors was identified through brainstorming sessions that involved all stakeholders. This list was refined later in the planning stage. The identified risk factors were categorized into four categories. These categories are content, process, technology and human risks. Content risks are those factors that are related to the courseware preparation, design, and development. Process risks are the risks that are involved in the course delivery and learning process. Technology risks are related to technological infrastructure of the online educational system including both hardware and software issues. Human risks involve risks related to the end users of the e-Learning course (i.e. students and tutors). A final refined list of risk factors is displayed in Table II. In this list, 43 risk factors were identified and categorized into the four categories. Each factor was given a number that uniquely identifies it (i.e. Risk Identifier (RID)).

### B. Planning Stage

In this stage, the planning covers four dimensions; risks, resources, course and evaluation planning.

Risk Planning: A set of mitigation strategies were devised for each of the identified risk factors. These strategies, listed in Tables III-XI were defined to be practiced later in the proper phases. The strategies were proposed after conducting brainstorming sessions that involved all project team members. Each member was asked to employ risk-based thinking and devise mitigation strategies that would be used to manage the previously defined risk factors. Moreover, these brainstorming sessions came out with new information about the already identified risk factors. As a result, a list of refined risk factors was produced as shown in Table II. These factors were the harvest of the project managers' experience, team members brainstorming sessions, and predefined ready-made checklist and taxonomies.

Resource planning: involved planning of team and roles, tools, students' analysis, budget and schedule, and divided into the following three aspects:

- Team Building: Team members were hired, roles were identified and assigned. The team involved and instructional designer (ID) who is responsible for defining the instructional, delivery and evaluation strategies, SME; the source of knowledge and responsible for content preparation, e-Content developer who is responsible for developing media components, assembling course elements and installing the courseware onto LMS, course administrator who manages learners accounts, online tutor who supports and motivates students learning activities during the course, and a technical support specialist who provides technical support for all stakeholders during all phases.

- Technology Tools: A decision was made about the tools needed to be used to create and deliver the e-Learning content. We used Adobe Photoshop for creating bitmap images, Adobe Illustrator for vectoral images, Adobe Flash for creating animations, Trivantis Lectora as an authoring tool, Sony Sound Forge for sound file editing, Moodle as a LMS, SQL Server as a Database Management System (DBMS), Google forms for conducting surveys, Microsoft (MS) PowerPoint for making presentations and MS Word for creating tutorials and documents.

- Students Analysis: Students come to the course with different backgrounds, abilities and varying levels of understanding, computer skills and technical experience. Hence, the course should be designed in a way that satisfies the needs of all these students. Moreover, analyzing students' backgrounds may reveal new risk factors that could threaten the development process, or lessen the severity of other factors. For these purposes, a pre-survey was conducted in the planning phase. The results and the detailed analysis of the pre-survey are introduced in the following section.

TABLE I.        SAMPLE DISTRIBUTION

| Semester | Enrolled students |
|---|---|
| Summer 2015-2016 | 24 |
| First 2016-2017 | 30 |
| Second 2016-2017 | 31 |
| Summer I 2016-2017 | 40 |
| Summer II 2016-2017 | 28 |
| First 2017-2018 | 33 |

TABLE II. RE COURSE PROJECT RISK FACTORS

| Category | RID | Risk Factor Description |
|---|---|---|
| Content Risks | 1 | Course is difficult to navigate |
| | 2 | Course is unusable |
| | 3 | Content is inaccessible |
| | 4 | Content is rigid and not interesting |
| | 5 | Course is not visually attractive |
| | 6 | Lack of interactivity |
| | 7 | Difficulties to work with several types of media content |
| | 8 | Requirements change |
| | 9 | Inadequate educational resources |
| | 10 | Low quality media content |
| | 11 | Loading delay |
| | 12 | Lack of consistency |
| | 13 | Low quality content |
| | 14 | Course structure is not understandable |
| | 15 | Content is difficult to understand |
| | 16 | Student cannot identify what should he know from each module |
| | 17 | Use of foreign language |
| | 18 | Content developers are not familiar with the content domain |
| Process Risks | 19 | Content cannot be deployed into LMS successfully |
| | 20 | Timing and sequencing of activities are unclear |
| | 21 | Course does not fulfill its stated objectives and learning outcomes |
| | 22 | Inadequate assessment |
| | 23 | Unclear assessment policy |
| | 24 | Violation of assessment procedures |
| | 25 | Lack of direct face to face interaction with the tutor |
| | 26 | Students miss collaborative work |
| | 27 | Miscommunication between team members |
| | 28 | Delivery delay |
| | 29 | Poor technical assistance |
| | 30 | Unauthorized access |
| | 31 | Violation of law |
| | 32 | Students are confused about what and how to learn. |
| Technology Risks | 33 | Problems in Internet connection |
| | 34 | Unreliable technical, hardware and software infrastructure |
| | 35 | Browser incompatibility |
| | 36 | Large number of concurrent connections to server |
| | 37 | Student does not have computers |
| | 38 | Update and upgrade risks |
| Human Risks | 39 | Tutor is inexperienced in e-Learning technologies |
| | 40 | Tutor's resistance to online learning |
| | 41 | Students' resistance to online learning |
| | 42 | Students lack the required computer skills |
| | 43 | Non-interactive tutor |

Course Planning: In course planning, high level objectives, course outlines, organizational structure, and assessment criteria were defined. The organizational structure of the RE course was hierarchal. The e-Content was composed of seven topics. The "Topic" is considered the basic unit of development. These Topics are: Introduction to RE, Photovoltaic Systems, Hydropower Energy, Geothermal Energy, Solar Thermal Energy, Wind Energy and Energy Economics. All topics are organized in similar ways. Each Topic is divided into sections and each section is divided into lessons. In order to evaluate students' knowledge and understanding of the material, it was agreed upon post assessment quizzes after each topic, mid-term and final exam. In addition, design projects were carried out online.

Evaluation Planning: It was decided to use both formative and summative assessments. Formative assessment was planned to be achieved through reviews conducted in the production phase and through feedback from students during the learning activities in order to collect information for the purpose of improving the e-Learning material being delivered. Summative assessment was decided to be at the end of the course in order to measure the effectiveness of the e-Learning process and the proposed framework. Most studies rely on user satisfaction or the acquired knowledge (or both) in order to evaluate an e-Learning process [49]. In this study, we use students satisfaction surveys at the end of the semester to find out how satisfied they are in the course and to make improvements based on their feedback. Moreover, students' final grades are analyzed at the end of the semester in order to measure the quality of the results and the acquired knowledge. Student's final grades consist of midterm and final exams in addition to the quizzes and projects taken during the course.

During each activity in each phase, the set of the mitigation strategies proposed in the planning phase are practiced cautiously by team members based on their roles in order to mitigate the negative effects of the corresponding identified risk factors. Lists of mitigation strategies to be practiced in each phase are introduced in Tables III-XI. Each strategy is uniquely defined using a strategy identifier (SID). Beside each of the identified strategies, the target risk factors IDs (TRIDs) that aim to mitigate are indicated.

For planning phase, a list of the proposed mitigation strategies involved in the planning phase is displayed in Table III.

*C. Production Stage*

The actual implementation of the e-Content was carried out at this stage. The development of the e-Content was accomplished module by module. Agile development was a major risk mitigation strategy to avoid requirements change and delivery delay risks as appears in Table IV. According to the framework, each module iteratively passed by 6 phases; design, content preparation, storyboard design, development, publishing, and review. The first phase of the production stage is the module design. A set of mitigation strategies were suggested in this phase in order to mitigate risk factors, these strategies are displayed in Table V. Once the module was designed, content preparation started. Several risk mitigation strategies were followed whilst preparing content in order to

mitigate risk factors. These strategies are depicted in Table VI. Creating storyboards is very important; they are themselves a mitigation strategy. The creation process also involves several risk mitigation strategies. These strategies are displayed in Table VII. In the development phase, the actual implementation of the courseware was carried on. The majority of content risks could be avoided by following several risk management strategies in this phase. These strategies are displayed in Table VIII. In the publishing phase, modules are produced in a shareable format to be handled by the LMS. The strategies that could be followed are displayed in Table IX. The last phase in the production is to test the courseware module and review it before it is delivered to students. The strategies that could be followed in this phase are displayed in Table X.

TABLE III.  RISK MITIGATION STRATEGIES IN THE PLANNING STAGE

| SID | Strategy | TRID |
|---|---|---|
| S1 | Allocate a variety of software and development kit | 7 |
| S2 | Ensure that the choice of software tools can easily support different file and media formats | 7 |
| S3 | Use courseware authoring tool | 7 |
| S4 | All tools are licensed | 31 |
| S5 | Hire professional media specialists (graphic designer, animator) | 7 |
| S6 | Team members are selected from university (insource) | 27 |
| S7 | Leverage talents from the organization (i.e. university students) to help in developing media content | 7 |
| S8 | Select a tutor with good computer skills and comfortable working online with students | 39,43 |
| S9 | Provide training on how to make the best use of online facilities for e-Learning support | 39,40, 41,42 |
| S10 | Training sessions on multimedia development | 7 |
| S11 | Impose computer skills classes as a prerequisite for the course | 42 |
| S12 | Breaking course into several modules (modular structure) | 8 |
| S13 | Describe the general organizational structure of the course | 14 |
| S14 | Assessment procedure is determined early in this phase | 23 |
| S15 | Mid and final summative exams were decided upon as formal Assessment | 22 |
| S16 | Quizzes are suggested to be used as part of the formal assessment process | 22 |

TABLE IV.  RISK MITIGATION STRATEGIES IN THE PRODUCTION STAGE

| SID | Strategy | TRID |
|---|---|---|
| S17 | Agile development | 8,28 |

TABLE V.  RISK MITIGATION STRATEGIES IN THE PRODUCTION STAGE - DESIGN PHASE

| SID | Strategy | TRID |
|---|---|---|
| S18 | Objectives and outlines are identified before each module | 16,32 |
| S19 | Objectives of the module should match the expected course outcomes | 21 |
| S20 | Include a short description for each module | 16 |
| S21 | Pre-quiz is selected as an evaluation strategy | 16 |
| S22 | Self-assessment is selected as an evaluation strategy | 21 |
| S23 | post-quiz is selected as an evaluation strategy | 21 |

TABLE VI.  RISK MITIGATION STRATEGIES IN THE PRODUCTION STAGE - CONTENT PREPARATION PHASE

| SID | Strategy | TRID |
|---|---|---|
| S24 | Instructional content relates directly to objectives | 21 |
| S25 | Select material from high-quality textbooks and articles | 9 |
| S26 | Use internal social media network to ask around for resource material | 9 |

| S27 | Provide links to extra material | 9,15 |
|---|---|---|
| S28 | Keep paragraphs short | 5,15 |
| S29 | Use common formal language "i.e. English" | 17 |
| S30 | Use simplest words and elaborations | 15,17 |
| S31 | Integrate real-life examples and problems into the course | 5,15 |
| S32 | There is a summary at the end of each module | 16 |
| S33 | Glossary is used to define key terms and abbreviations | 16 |

TABLE VII.  RISK MITIGATION STRATEGIES IN THE PRODUCTION STAGE - STORYBOARD DESIGN PHASE

| SID | Strategy | TRID |
|---|---|---|
| S34 | Script is prepared jointly (i.e. content provider and content developer) | 13,18 |
| S35 | Script describes the organizational structure of the module and exact table of contents | 1,14 |
| S36 | Script should describe in detail media elements and desired interactions | 13 |
| S37 | Do not include more than two paragraphs or 7 bullets items to a page. | 5, 15 |
| S38 | Use bullets, tables, callouts, interspersed images to organize concepts | 15,4,5 |
| S39 | Punctuation and capitalizations are used appropriately | 13 |
| S40 | Avoid monolithic: chunk information into small pieces | 15 |
| S41 | Rely on the power of interactivity | 15, 6 |
| S42 | Employ humour to emphasize a point in relevant, light-hearted way | 4 |
| S43 | Script lesson is proofread in terms of content by SME | 13 |
| S44 | Script lesson is proofread in terms of language by language editor and word processor | 13 |

TABLE VIII.  RISK MITIGATION STRATEGIES IN THE PRODUCTION STAGE - DEVELOPMENT PHASE

| SID | Strategy | TRID |
|---|---|---|
| S45 | The eContent design allow student to pause and resume the course without losing their place | 1 |
| S46 | Animations and navigational elements does not distract focus of attention | 1 |
| S47 | Hyperlinks are clearly identified. | 1 |
| S48 | All hyperlinks work and direct the student to the proper location | 1 |
| S49 | Only reasonable hyperlinks and navigation elements are provided (students are not overwhelmed with hyperlinks) | 1 |
| S50 | Hyperlinks are not introduced early at paragraphs need to be read completely or near important points. | 1 |
| S51 | Navigation is allowed using "back" and "forward" buttons, table of contents, and navigation path. | 1 |
| S52 | Navigation buttons are easily identifiable, and perceive a good level of affordance | 1 |
| S53 | At any time student can identify his location with respect to course using a navigation path | 1 |
| S54 | Develop a course map that enables student get an at-a-glance view of course content | 1 |
| S55 | Clear instructions are designed to prevent possible runtime errors | 2 |
| S56 | Course layout inspires student what to do in each page | 2 |
| S57 | e-Learning activities are labelled and numbered in a way guides learners through the course. | 1,2 |
| S58 | Error messages and direction are expressed in understandable language | 2 |
| S59 | Needed information is provided on the same screen to minimize recall | 2 |
| S60 | Horizontal and vertical scrolling is avoided | 2 |
| S61 | Use icons related to actions | 2 |
| S62 | Avoid colour combinations that are problematic for colour blind people | 3 |
| S63 | Use large enough thick fonts | 3 |
| S64 | Combine sound narration to highlight certain points or to provide certain comments on animations | 3 |

| S65 | Use Alt tags to describe images | 3 |
| S66 | Text is provided for all non-text elements | 3 |
| S67 | Use Font colours visible against background colour | 3,5 |
| S68 | (i.e. Main body text uses black font against white background) | 3 |
| S69 | Main body text uses Sanserif font (i.e. Verdana) | 3 |
| S70 | Avoid unnecessary colours, instead use bold and italic to emphasize | 3 |
| S71 | Use appropriate line spacing | 3 |
| S72 | Use play/pause to control sound | 4,5,15 |
| S73 | Incorporate various types of multimedia (text, graphics, audio, video animations) | 4 |
| S74 | Employ narrative storytelling | 4 |
| S75 | Use appealing, simple, informative and helpful voice | 4, 5 |
| S76 | Use bullets, tables, callouts, interspersed images | 4, 6 |
| S77 | Use thought-bubble callouts that appeared alongside our content with characters' faces | 5 |
| S78 | Balance between text and graphics | 5,12 |
| S79 | Use decorative fonts only for headings | 5 |
| S80 | Paragraph are justified | 4,6 |
| S81 | Include quizzes with feedback at the end of each module | 6 |
| S82 | Use of animations, navigations, and learning objects. | 6 |
| S83 | Use of interactive infographics | 10 |
| S84 | Sound files are recorded and edited in .wav format | 10,11,33 |
| S85 | Use PNG format for images (i.e. lossless data compression format) | 10 |
| S86 | Most animations have framerate of 24 fps | 3 |
| S87 | No animations faster than 30 fps | 10 |
| S88 | Animations narration sound files compressed as speech 22 kHz | 10 |
| S89 | Animations use lossless compression PNG images | 10 |
| S90 | Use text to speech software | 12 |
| S91 | Use the same voice through different modalities such as videos, animations, infographics, text, etc. | 12 |
| S92 | Each module has the same layout | 12 |
| S93 | Keep fonts types and formats consistent through the course | 12 |
| S94 | Do not use more than 3 fonts | 12 |
| S95 | Overall design is uniformed | 12 |
| S96 | Colours are used consistently | 12 |
| S97 | Navigation is consistent throughout the course | 12 |
| S98 | Animations and videos are consistent in quality, size and type | 12 |
| S99 | Images are consistent in quality, size and type | 31 |
| S100 | References are stated at the end of each topic | 31 |
| S101 | Graphics, videos, animations, books copyrights reserved | 33 |
| S102 | Make e-Learning content smaller | 33 |
| S103 | Videos and images that require a high speed Internet are avoided | 33 |

TABLE IX.    RISK MITIGATION STRATEGIES IN THE PRODUCTION STAGE - PUBLISHING PHASE

| SID | Strategy | TRID |
| --- | --- | --- |
| S104 | Publish module into sharable content object reference model (SCORM) format | 19 |
| S105 | Publish illustration videos onto a YouTube channel | 15 |

TABLE X.    RISK MITIGATION STRATEGIES IN THE PRODUCTION STAGE - REVIEW PHASE

| SID | Strategy | TRID |
| --- | --- | --- |
| S106 | Involve the SME and the developer in the review | 13 |
| S107 | Test the content on different mobile devices | 3 |
| S108 | Test the content on the most common browsers | 35 |
| S109 | Testing was done with a student account | 3, 30 |
| S110 | All combinations of assessments were tested | 22 |

## D. Delivery Stage

In the delivery stage, e-Content modules are made available to the learners and the learning process takes place. The list of strategies that could be followed to handle risk factors (mainly process factors) is described in Table XI.

TABLE XI.    RISK MITIGATION STRATEGIES IN THE DELIVERY STAGE

| SID | Strategy | TRID |
| --- | --- | --- |
| S111 | A help video was produced and its link is included on the main page. | 1, 2 |
| S112 | Hardware and software requirements are identified early | 3, 37 |
| S113 | Direct download links are provided for the required plugins (PDF reader, SWF player) | 3 |
| S114 | Design a study guide and course syllabus | 20 |
| S115 | Identify mandatory assignments and submission dates early | 20 |
| S116 | Each module is given a sufficient time in the syllabus | 20 |
| S117 | Feedback from students are sought periodically | 21 |
| S118 | Assessment includes several types of questions | 22 |
| S119 | Feedback is given for each question | 22 |
| S120 | Assessments were challenging and realistic | 22 |
| S121 | Every learning objective is assessed | 22 |
| S122 | Exams, quizzes, graded assignments are clarified early | 23 |
| S123 | Exams and quizzes were held online but not distant | 24 |
| S124 | Exams and quizzes held at university lab (broadband connection) | 33 |
| S125 | Quizzes questions and choices (if there is) were shuffled | 24 |
| S126 | Quizzes and exams had "start and end time" | 24 |
| S127 | Student names and login information are imported from university registration systems | 24 |
| S128 | Conduct virtual office hours using synchronous tools (i.e. chat rooms, instant messaging) | 4,6,25, 26 |
| S129 | Use asynchronous communication tools (i.e. email, forums, social networks groups) | 4,6,25, 26 |
| S130 | Bulletin board is used for general announcements from instructors | 25 |
| S131 | Assignments are submitted and graded with tutor feedback via Moodle | 25 |
| S132 | Tutors are given administrative privileges regarding content management. | 29 |
| S133 | Contact information of the technical support is given to both tutors and students | 29 |
| S134 | Assign a separate topic in the discussion forum for reporting on technical problems | 29 |
| S135 | Accounts were set carefully with privileges based on roles | 30 |
| S136 | Only registered students are enrolled and can access the course | 30 |
| S137 | Once the course is finished, all enrolment are cancelled | 30 |
| S138 | Use firewall to control access | 30 |
| S139 | References are stated at the end of each topic | 31 |
| S140 | Assign a tutor that facilitates the learning process | 32 |
| S141 | Technical assistance team | 32 |
| S142 | Configure the maximum worker threads server configuration option | 36 |
| S143 | Guarantee the availability of services using redundant Moodle server | 34 |
| S144 | Regular data backups are taken from the course in case of a breakdown of certain components | 34 |
| S145 | A high technical specification lab is dedicated to course access | 34, 37 |
| S146 | Regular backup of the course | 38 |
| S147 | Avoid updates or upgrades at critical times | 38 |
| S148 | The version of the course is indicated | 38 |
| S149 | Orientation day to motivate students to online learning | 41,42 |

### E. Evaluation Stage

Evaluation should be done at all stages; before the start date of the course, during the delivery, and after its completion. The post-course evaluation (i.e. summative evaluation) is the most important and the most challenging one. It is achieved by conducting a comprehensive survey at end of the course.

- Satisfaction Survey

At the end of the course, feedback from all project stakeholders was sought. The evaluation was conducted on two types of users, namely students and development team. Being the focal aspect of the e-Learning process, students' feedback is given a high concentration in the evaluation process. Students' feedback is very important. Throughout this feedback, students can describe their learning experience in the course. They can describe the content, material, activities, course design, delivery process, assessment methodology, etc. From their comments, pros and cons of the course are revealed so that they can be considered in the future to improve the design, deliverables and the delivery of the course. Herein, student's satisfaction surveys were used; a questionnaire-based approach that relies solely on students and how they were satisfied with the course. In this context, the main purpose of this questionnaire was to measure how the proposed framework and more specifically the proposed risk management strategies were effective in mitigating the identified risks. It was believed that if students exhibit high satisfaction towards the different course components, then the risk factors had been successfully managed by the proposed framework. Moreover, this questionnaire highly supports the internal quality assurance mechanism used at HU.

As mentioned before, post-course evaluation is the most challenging. This could be ascribed to two reasons. First, this survey should be comprehensive in a way that covers all key aspects of both design and delivery processes. The second is that this comprehensive nature of the survey would result in a quite lengthy survey that overwhelms students. In order to overcome these barriers and to encourage students to react to survey, the questionnaire was conducted online (using Google Forms) at the end of each semester. The questionnaire consists of 38 5-Likert scale questions (1 = Strongly Disagree, 2 = Disagree, 3 = Neutral, 4 = Agree, and 5 = Strongly Agree). The questions are formulated in a way that covers all risk categories and factors. The scale consisted of four subscales that measure course content, learning process, technology use, and people involved. Each question corresponds to one or more risk factors and measures whether each of these risk factors was mitigated well.

In purpose of evaluating students' satisfaction, the arithmetic mean score was used. The mean is the average of all responses for each item. The mean is often used to report central tendency of Likert items. It gives the best overall statistic of the typical rating given by survey respondents since it takes all data into account [50]. Mean score values above 4 are considered strong, between 3.5 and 4 are considered solid, and scores below 3.5 should be of concern [51]. Moreover, to measure the variability of students' responses, the standard deviation for each item was calculated. An item with a standard deviation greater than 1 indicates a wide variety in students' responses [51].

In order to ensure the comprehensive and the assimilation of students, the questionnaire was conducted among the students at the end of each semester to find out how satisfied they were in the course. Students' responses in all semesters are analyzed in the following section.

- Grades Analysis

It was believed that without measuring learning effectiveness, the quality of the e-Content and the e-Learning process could not be evaluated. Students' results and their grades were used as indicators of learning effectiveness and indicate the level of student achievement in the course. Student's final grades were drawn electronically from the university registration system. It consisted of the midterm, final exams, quizzes, and projects/assignments taken during the course. Exams and quizzes were administered to consist of good quality question items that cover all learning outcomes defined in the course. Moreover, information obtained from this evaluation can be used to enhance the pedagogical quality of the content.

## V. RESULTS, ANALYSIS AND DISCUSSION

In this section, the results of the pre-survey that aimed at identifying new risk factors related to students or updating the severity of the already defined ones are introduced in details. More important, the effectiveness of the proposed framework in developing and delivering high quality e-Learning is verified based on satisfaction surveys and grades analysis. The results of implementing the framework in developing and delivering RE course in six semesters during the period 2015-2017 are introduced.

### A. Pre-Survey

Students backgrounds were analyzed through a pre-survey that aimed at identifying other risk factors related to the learners. The pre-survey was conducted online and made available to the students via LMS. The pre-survey targeted a random sample of the Mechatronics Engineering students from different levels. A total number of 121 responses were collected. The pre-survey covered five dimensions: language and technical skills, e-Learning experience, feelings and doubts, hardware and software platforms, and disabilities.

Fig. 2 exhibits students' backgrounds and skills. When questioned about their fluency in English only 55 students (45.45%) said that they are fluent in English. This percentage would emphasize the "use of foreign language" risk factor. Regarding computer literacy, 83 students (68.60%) acknowledged that they have good computer skills. This fair percentage would also emphasize the risk factor "Students lack the required computer skills". A total number of 75 students (61.98%) used to access information from the web, 85 students (70.25%) feel at ease with online technology, 55 students (45.45%) used to use forums, and 52 students (42.98%) used to participate in chat rooms. The relatively low percentage of using forums and chat rooms would negatively affect students' interaction with the tutor and other students through the synchronous and asynchronous communication mechanisms.

Students were asked if they already had any e-Learning experience before, Fig. 3. Only 15 students (12.40%) enrolled full online course. More specifically, 48 students (39.67%) used forums in other courses, 39 students (32.23%) used chatrooms, 56 students (46.28%) practiced self-assessment program, 80 students (66.12%) practiced online exams, and 52 students (42.98%) used VLE to download material and resources. Overall, most e-Learning activities are practiced with a low percentage (excluding online exams) as appears in Fig. 3. This would negatively affect students' usage of course activities. Regarding VLE, 100 students (82.64%) had used Moodle VLE, and only 21 students (17.36%) had used a VLE other than Moodle. Hence, a decision was made to adopt Moodle as a VLE for this course.

Students also were asked about their feelings, doubts, and worries about the course, Fig. 4. A total number of 28 students (23.14%) preferred to take this course in a traditional classroom than totally online, 35 students (28.93%) were uncertain and 58 students (47.93%) preferred online delivery. From this, we can conclude that still there is some resistance from students towards online learning (i.e. risk factor). Deeply, 31 students (25.62%) were worried about the absence of the face-to-face tutor, 42 students (34.71%) were uncertain and 48 students (39.67%) had no worries about this issue (another risk factor). 29 students (23.97%) were worried about being isolated from other colleagues, 23 students (19.01%) were uncertain and 69 students (57.02%) were not worried. This fair percentage could be ascribed to the various currently available communication technologies which make this factor less severe. Moreover, 45 students (37.19%) were worried about the assessments and grading policy, 35 students (28.93%) were uncertain and 20 students (33.88%) had no worries. This would increase the importance of mitigating this risk factor.

Additionally, students were questioned about the hardware and software platform they have, Fig. 5. A total number of 115 students (95.04%) students have personal computers. This high percentage would eliminate this risk factor. Moreover, 51(44.35%) have windows 7 installed, 12 students (10.43%) have windows 8, 42 students (36.52%) have windows 10, and 10 students (8.70%) have other operating systems. Also, 40 students (34.78%) have 1-2 GB RAM and 75 students (65.22%) have RAMs larger than 2GB. Small size RAMs may negatively affect course navigation, delay playing media content, etc. 97 students (84.35%) uses Chrome browser, 5 students (4.35%) IE and 7 students (6.09%) Firefox and 6 students (5.22%) used other browsers. Accordingly, students use different Internet Browsers; browsers incompatibility is another risk factor. Regarding plugins, 109 students (94.78%) have a PDF reader installed on their pcs, and 55 students (47.83%) have a SWF player installed. This might prevent the student from playing animations and other media contents (i.e. risk factor). With respect to the Internet connection, 103 students (85.12%) have Internet connection 24/7. Among them, 28 students (27.18%) have up to 2 Mbps Internet connection, 49 students (47.57%) have up to 8 Mbps and 26 students (25.24%) has up to 16 Mbps. Problems in connection is a risk factor. Moreover, Low connection speed may cause loading delay (another risk factor).



Fig. 2.   Students' Backgrounds and Skills.



Fig. 3.   Students' e-Learning Experience.



Fig. 4.   Students' Feelings and Doubts.



Fig. 5.   Technical Infrastructure.

Lastly, regarding disabilities, 25.62% students have poor vision. None of them are either visually impaired, color blind, deaf or have movement disabilities. People with disabilities may not properly access the course (risk factor).

*B. Satisfaction Surveys*

Satisfaction surveys involve quantitative and qualitative assessments. Quantitative assessment targets the learners in purpose of measuring how satisfied they were with the course. Each quantitative assessment item is associated with a risk factor. Based on students' reaction towards each item, we can measure how the model and the proposed risk management strategies were effective in mitigating each risk factor. Qualitative assessment targets other stakeholders (i.e. mainly development team) to measure how specific-related risk factors are mitigated. Below, satisfaction surveys results are introduced based on risk factors categories. The last part of the survey (i.e. Global items) involves general question items to measure how students are satisfied with the course in general.

- Content Risks

*1) Quantitative analysis:* Concerning content risks (see Table XII), students reacted positively towards most items with a mean average value equals 4.0. Items 1-5, 7-11, and 13 have mean values greater than or equal 4, then, according to [51] these responses are described as "strong". Items 6, 12, and 14 have mean values between 3.5 and 4 then they are considered "solid". In order to assess the variety in responses, standard deviation was calculated for each item. As appears in Table XII, all standard deviation values are around 1 which indicates little variations in responses. Moreover, the most frequent answer for all items were either Strongly Agree or Agree (mode =5 or 4). This implies that content risk factors were mitigated very well. The solid mean values for items 6, 12, and 14 indicates that risk factors 6, 15, and 17 were mitigated but special attention should be paid in the future to increase students' satisfaction towards items 6, 12, and 14. In other words, new mitigation strategies should be followed to increase course interactivity and simplify the content and language. From another point of view, low values in language simplicity could be due to students' low fluency in English language in general (only 45.45% of the students said they are fluent). Moreover, the loading delay encountered could be ascribed to the relatively low Internet connection speed (74.75% of department's students have download speeds up to only 8 Mbps).

*2) Qualitative analysis:* Regarding risks 7, 8, 9, 18, and 19, development team were enquired to assess how these risks were mitigated. They said that they faced no difficulties in working with different types of media content due to the availability of tools, professionals and the training sessions they had. When questioned about the requirements change, they reported that they had few challenges since the content was implemented incrementally in iterations with a continuous review with the SME; changes were adopted early. Content developers reported that they had very few difficulties in collecting material; they employed social networks, and used high quality books and web links to write high quality content. When asked about their familiarity with the domain, content developers said that they are unfamiliar with RE domain but they could overcome this risk by preparing the script jointly with the content providers and SMEs. Developers also reported that they had no problems in deploying the e-Content into the LMS since they had published into as shareable format (i.e. SCORM). Hence, also risk factors 7, 8, 9, 18, and 19 were mitigated successfully.

- Process Risks

*1) Quantitative analysis:* Regarding process risks (see Table XIII), with mean values greater than 4, students exhibited strong satisfaction towards items 15-26 (except 22) regarding feel of isolation (item 22), a solid mean value of 3.8 was obtained. Indeed, 70.37% did not feel isolated from other colleges which is slightly larger than the (57.02%) who were worried about this issue before taking the course. More mitigation strategies should be taken to increase this percentage. The items mark little variances since almost all items values are around 1. Also, the mode was either 4 or 5. This implies that risk factors 20-26, 28-29, and 31-32 were mitigated very well.

*2) Qualitative analysis:* Concerning risks 27 and 30, development team were enquired to assess how these risks were mitigated. They said that they faced very little communication problems because they all work in the same institution (i.e. HU). Regarding access authorization, they said that they did the testing using student accounts to make sure that they are given the right permissions. Hence, also risk factors 27 and 30 were mitigated successfully.

TABLE XII. STUDENTS' SATISFACTION SURVEY RESULTS - CONTENT RISKS ITEMS

| Item | Description | Mean | SD | Mode | TRID |
|---|---|---|---|---|---|
| 1 | I was able to navigate the course easily | 4.1 | 1.1 | 5 | 1 |
| 2 | The course was easy to use | 4.1 | 1.1 | 5 | 2 |
| 3 | I could easily access course content and activities | 4.3 | 1.1 | 5 | 3 |
| 4 | The course was interesting | 4.1 | 1.2 | 5 | 4 |
| 5 | I liked the design of the course | 4.0 | 1.1 | 5 | 5 |
| 6 | The course was highly interactive | 3.8 | 1.0 | 4 | 6 |
| 7 | Media (images, sounds, animations, videos) quality was high | 4.0 | 1.1 | 5 | 10 |
| 8 | There was no large delay in loading media elements | 4.0 | 1.1 | 5 | 11 |
| 9 | Course design and components were consistent | 4.0 | 1.0 | 4 | 12 |
| 10 | The course was free of syntactical, grammatical errors and mistaken information | 4.0 | 1.1 | 5 | 13 |
| 11 | The organization of the course into units and subunits was clear | 4.2 | 1.0 | 5 | 14 |
| 12 | The content was easy to understand | 3.9 | 1.0 | 4 | 15 |
| 13 | The objectives and outlines for each module were clear | 4.2 | 1.0 | 5 | 16 |
| 14 | The language was simple and clear | 3.8 | 1.2 | 5 | 17 |

TABLE XIII.   STUDENTS' SATISFACTION SURVEY RESULTS - PROCESS RISKS ITEMS

| Item | Description | Mean | SD | Mode | TRID |
|---|---|---|---|---|---|
| 15 | Timing and sequencing of topics and other activities were clear | 4.1 | 1.1 | 5 | 20 |
| 16 | The aims and objectives of the course were achieved | 4.1 | 1.1 | 5 | 21 |
| 17 | Assignments and exams were related to the materials taught in the class | 4.3 | 1.0 | 5 | 22 |
| 18 | Assessments were satisfactory | 4.2 | 1.0 | 5 | 22 |
| 19 | Grading criteria were outlined in the syllabus | 4.4 | 1.0 | 5 | 23 |
| 20 | Assessment procedure can NOT be easily violated (i.e. cheating, unauthorized access, impersonate, etc.) | 4.0 | 1.1 | 5 | 24 |
| 21 | I would prefer to take this course online rather than in a traditional classroom | 4.0 | 1.3 | 5 | 25 |
| 22 | I did not feel isolated from my colleagues | 3.8 | 1.2 | 4 | 26 |
| 23 | The unit's modules were delivered onto Moodle on time | 4.3 | 1.0 | 5 | 28 |
| 24 | The technical support for this course was effective | 4.1 | 1.0 | 5 | 29 |
| 25 | References and copyrights were declared clearly | 4.1 | 1.1 | 5 | 31 |
| 26 | I was able to know what and how to learn | 4.2 | 1.0 | 5 | 32 |

- Technology Risks

*1) Quantitative analysis:* Regarding technology risks (see Table XIV), students reacted positively towards items 27-28 with mean values of 3.7 and 3.9 respectively (solid values) and an average value of 3.8. Also, the mode for both was 5. This implies that risk factors 34, 36 were mitigated very well. Regarding Browser incompatibility risk (factor 35), students were asked to mention the browsers they used to navigate the course. A percentage of 72.66% of students used Internet Explorer, 91.37% used Google chrome, 15.83% used safari, 10.07% used Firefox, and 24.46% used android Internet browser. Hence, the course could be viewed using different browsers; avoiding risk 35.

*2) Qualitative analysis:* Concerning update and upgrade risks (Factor 38), development team reported that they did not face any update problem since they did not perform any update or upgrade during the delivery phase of the course in order to avoid these issues. The vast majority of the students (as from the pre-survey) reported that they have personal computers which would eliminate Factor 37. Moreover, a computer lab was dedicated to facilitate students' access to course content.

TABLE XIV.   STUDENTS' SATISFACTION SURVEY RESULTS - TECHNOLOGY RISKS ITEMS

| Item | Description | Mean | SD | Mode | TRID |
|---|---|---|---|---|---|
| 27 | I could access the course even if the main server is down | 3.7 | 1.13 | 5 | 34 |
| 28 | No significant drop in performance occurs when large number of students access the server concurrently (i.e. exams) | 3.9 | 1.12 | 5 | 36 |

- Human Risks

*1) Quantitative analysis:* Obviously human risks were mitigated well (see Table XV). Students reacted positively towards items 29-33 mean values above 4. Moreover, the most frequent answer (mode) for all items was Strongly Agree. This implies that the risk factors 39-43 were mitigated very well.

- Global Items

The last part of the survey (see Table XVI) reflects global measures that assess the overall satisfaction of the RE course project. Concerning global items 34-38, a percentage of 83.33% of students said that the course saved their time with a mean of 4.3. On the other hand, 68.52% said that the course saved their money with a mean of 3.9. 77.16% recommend other students to take this course online in contrast with only 47.93% who had preferred, in the students pre-analysis, to take this course totally online. Moreover, 73.46% would like to take other online courses in the future. Finally, overall, 83.95% were satisfied with the course in general with a mean of 4.2. This implies that students generally were highly satisfied with the course.

TABLE XV.   STUDENTS' SATISFACTION SURVEY RESULTS - HUMAN RISKS ITEMS

| Item | Description | Mean | SD | Mode | TRID |
|---|---|---|---|---|---|
| 29 | The tutor was experienced in e-Learning issues | 4.1 | 1.01 | 5 | 39 |
| 30 | The tutor was motivator and e-Learning supportive | 4.2 | 0.97 | 5 | 40 |
| 31 | I was motivated to take this course online | 4.1 | 1.16 | 5 | 41 |
| 32 | I already have the required computer skills | 4.3 | 1.05 | 5 | 42 |
| 33 | The tutor was accessible and prepared to teach the course online | 4.2 | 1.04 | 5 | 43 |

TABLE XVI.   STUDENTS SATISFACTION SURVEY - GLOBAL ITEMS

| Item | Description | Mean | SD | Mode | Percentage |
|---|---|---|---|---|---|
| 34 | Taking the course online saved my time | 4.3 | 1.1 | 5 | 83.33% |
| 35 | Taking the course online saved my money | 3.9 | 1.2 | 5 | 68.52% |
| 36 | I recommend other students to take this course online | 4.1 | 1.2 | 5 | 77.16% |
| 37 | I would like to take other courses online in the future | 4.0 | 1.3 | 5 | 73.46% |
| 38 | Overall, I would rate this course Excellent | 4.2 | 1.0 | 5 | 83.95% |

*C. Grades Analysis*

In order to validate the effectiveness of the proposed model, the quality of the e-Content and the e-Learning process was evaluated. Students' grades were used as a measurement for the quality of the e-Course content and process. In this section, students' results along six semesters (during 2015-2017) of the electronic delivery of the course are introduced in Table XVII. In contrast, students' results along the preceding three semesters (during 2014-2015) of the traditional delivery

of the course are also introduced Table XVIII. Then, a comparison between the results of the electronic delivery versus the traditional delivery is introduced in Fig. 6.

Table XVII shows that the students achieved good results in the electronic delivery form of the course. In compare with the traditional delivery of the course, Fig. 6 shows that e-Course students performed similarly or even better in most categories. The percentage of students who failed in RE e-Course along the semesters was 3.76% in contrast with 9% in the traditional course. Moreover, the percentage of students who got high grades (i.e. above B) in the e-Learning course is slightly larger than the traditional. For instance, a percentage of 9.13% of students got B in the e-Course whilst only 4% got B in the traditional. Another percent value of 8.06% of students got B+ compared to 8% in the traditional. A percentage of 11.29% of online students got A- in contrast with 10% of the traditional, and 5.91% got A in the online, while 5% got the same mark in the traditional . These values indicate a high quality electronic course and therefore validate the effectiveness of the proposed proactive risk management model leveraged in the e-Course development and delivery.



Fig. 6.   e-Course Students Results Versus Traditional Results.

TABLE XVII.   .RE E-COURSE STUDENTS GRADES

| Grade | Range | Sum. 2015-2016 | First 2016-2017 | Second 2016-2017 | Sum.I 2016-2017 | Sum.II 2016-2017 | First 2017-2018 | Percentage |
|---|---|---|---|---|---|---|---|---|
| A+ | 90.00-100.00 | 2 | 2 | 2 | 1 | 0 | 1 | 4.30% |
| A | 86.00-89.99 | 2 | 3 | 2 | 2 | 2 | 0 | 5.91% |
| A- | 82.00-85.99 | 2 | 3 | 5 | 2 | 6 | 3 | 11.29% |
| B+ | 78.00-81.99 | 5 | 2 | 4 | 0 | 2 | 2 | 8.06% |
| B | 74.00-77.99 | 2 | 6 | 1 | 4 | 1 | 3 | 9.14% |
| B- | 70.00-73.99 | 1 | 3 | 6 | 2 | 1 | 5 | 9.68% |
| C+ | 66.00-69.99 | 3 | 4 | 3 | 5 | 4 | 4 | 12.37% |
| C | 62.00-65.99 | 2 | 4 | 1 | 7 | 1 | 3 | 9.68% |
| C- | 58.00-61.99 | 2 | 2 | 3 | 5 | 3 | 5 | 10.75% |
| D+ | 54.00-57.99 | 1 | 1 | 1 | 6 | 4 | 1 | 7.53% |
| D | 50.00-53.99 | 1 | 0 | 2 | 4 | 1 | 6 | 7.53% |
| F | 0.00-49.99 | 1 | 0 | 1 | 2 | 3 | 0 | 3.76% |
| Total | | 24 | 30 | 31 | 40 | 28 | 33 | |

TABLE XVIII. RE TRADITIONAL COURSE STUDENTS' GRADES

| Grade | Range | First 2014-2015 | Second 2014-2015 | First 2015-2016 | Percentage |
|---|---|---|---|---|---|
| A+ | 90.00-100.00 | 2 | 2 | 1 | 5.00% |
| A | 86.00-89.99 | 2 | 2 | 1 | 5.00% |
| A- | 82.00-85.99 | 5 | 1 | 4 | 10.00% |
| B+ | 78.00-81.99 | 3 | 2 | 3 | 8.00% |
| B | 74.00-77.99 | 0 | 1 | 3 | 4.00% |
| B- | 70.00-73.99 | 2 | 3 | 8 | 13.00% |
| C+ | 66.00-69.99 | 2 | 5 | 4 | 11.00% |
| C | 62.00-65.99 | 2 | 3 | 7 | 12.00% |
| C- | 58.00-61.99 | 1 | 3 | 3 | 7.00% |
| D+ | 54.00-57.99 | 1 | 3 | 5 | 9.00% |
| D | 50.00-53.99 | 0 | 3 | 4 | 7.00% |
| F | 0.00-49.99 | 0 | 6 | 3 | 9.00% |
| Total | | 20 | 34 | 46 | |

Delivering the RE e-Course online for 6 semester achieved promising results. The authors believe that continuing delivering the course for many other semesters with a larger sample size and constant continuous improvement through the semesters will yield better results. The study revealed that technical issues related to hardware infrastructure, network connection, and server availability still would negatively affect communication and interactivity despite the good levels achieved. Moreover, the questionnaire was not checked for validity, this issue was beyond the current study.

## VI.   CONLUSION AND FUTURE WORK

Improving education has been a prime priority for HEIs that seek to employ technology in learning to generate knowledgeable students. The success of e-Learning depends mainly on the quality of the course. Hence, in this paper, a proactive risk management framework that aims at ensuring high-quality e-Learning has been introduced. The framework was implemented in an online RE e-Course. In order to validate the effectiveness of the framework in developing the course, both users' satisfaction studies and students' grades analysis were conducted. This framework embeds an iterative approach to the development of instructional design products. Also, the framework is adaptable; it can be tailored to suit any e-Learning project according to its objectives, characteristics, audience and probable risks. Hence, the model proposed in this paper will serve as a tool for HEIs to define their customized models based on courses contexts and risks. The development of the RE e-Course was guided by the proposed framework. Its learning objectives and activities were designed and implemented according to good design principles and best practices in the literature with an eye towards avoiding e-Learning risk factors. Throughout the RE course case study, this paper pointed out 43 e-Learning risk factors that need to be addressed and 148 risk management strategies needed to address these factors. Hence, the paper also sets the foundations to overcome these factors and to improve the e-Learning approach in HE.

Satisfaction surveys analysis used Qualitative data and quantitative measurements; including mean, mode and standard deviation. These surveys were used to assess how each of these risks was mitigated. These surveys revealed that each of these

risks has been mitigated to a certain degree. Moreover, students' grades analyses were conducted to assess the quality of the learning. The study revealed that participants who concluded the course were highly satisfied and achieved good results compared with traditional course. This implies the effectiveness of the proposed framework in developing and delivering high-quality e-Learning courses. The approach in this framework is not the only way to develop an effective e-Learning content. Rather, it is a reasonable approach to ensure e-Learning project success and high-quality e-Courses based on the results we obtained. One dimension of the future work is to implement other online courses with a larger audience based on the framework to ensure generalizability of the framework. An issue is concerned with the validity of satisfaction surveys is left for future studies. Regarding items with less satisfaction values, further improvements are required to properly address the related risks. Furthermore, the correlation between students' satisfaction and the knowledge gain could be examined in the future.

## REFERENCES

[1] Barik, N., & Karforma, S. (2012, January). Risks and remedies in e-learning system. International Journal of Network Security & Its Applications, 4(1), 51-59. doi:10.5121/ijnsa.2012.4105

[2] Jovanovic de Bozinoff, M., & Taskosic, M. (2014). E-learning risks management as competitive advantage in institutions of higher education. Modern Computer Applications in Science and Education, (pp. 164-170). Cambridge.

[3] Mahmud, K., & Gope, K. (2009). Challenges of implementing e-learning for higher education in least developed countries: a case study on Bangladesh. Proceedings of the 2009 International Conference on Information and Multimedia Technology (pp. 155-159). Jeju Island, Korea: IEEE Computer Society, Washington, DC, USA. doi:DOI 10.1109/ICIMT.2009.27

[4] Dziuban, C., Graham, C. R., Moskal, P. D., Norberg, A., & Sicilia, N. (2018). Blended learning: the new normal and emerging technologies. International Journal of Educational Technology in Higher Education, 15(3), 1-16.

[5] Mazohl, P., & Makl, H. (2017). Quality assurance in blended learning - a quality framework. In P. G. Mazohl, H. M. Mas, F. Breitenecker, A. Koerner, & S. Winkler, Blended learning Quality – Concepts Optimized for Adult Education (pp. 23-37).

[6] Misut, M., & Pribilova, K. (2015). Measuring of quality in the context of e-learning. social and behavioral sciences 177, (pp. 312 – 319).

[7] Hoppe, G., & Breitner, M. H. (2004). Business models for e-learning. Multikonferenz Wirtschaftsinformatik, (pp. 3-18).

[8] Woodill, G., & Pasian, B. L. (2015). E-Learning project management: a review of the literature. In B. Pasian, & G. Woodill, Plan to Learn: Case Studies in e-Learning Project Management (pp. 4-10).

[9] Durdu , P. O., Yalabik , N., & Cagiltay , K. (2009). A distributed online curriculum and courseware development model. Educational Technology & Society, 12(1), 230–248.

[10] Baruque, L. B., & Brazil, A. L. (2014). Managing e-learning content development risks. Proceedings of E-Learn: World Conference on E-Learning in Corporate, Government, Healthcare, and Higher Education. 1, pp. 145-152. New Orleans, LA, United States: Association for the Advancement of Computing in Education. doi:DOI 10.13140/RG.2.1.3368.4003

[11] Vesper, J. L., Kartog˘lu, Ü., Herrington, J., & Reeves, T. C. (2016). Incorporating risk assessment into the formative evaluation of an authentic e-learning program. British Journal of Educational Technology, 47(6), 1113–1124.

[12] Ifinedo, P. (2005). Uncertainties and risks in the implementation of an e-learning information systems project in a higher-learning environment: viewpoints from Estonia. Journal of Information & Knowledge Management, 4(1), 1-10.

[13] Bandaya, M. T., Ahmed, M., & Jan, T. R. (2014). Applications of e-learning in engineering education: a case study. Socil and Behavioral Sciences 123, (pp. 406 – 413).

[14] Alqrainy, S., & Hijazi, H. (2014). Managing risks in the system analysis and requirements definition phase. International Journal of Computer Applications, 23-29.

[15] Galvão, T. A., Neto , F. M., Campos, M. T., & Júnior, E. d. (2012). An approach to assess knowledge and skills in risk management through project-based learning. International Journal of Distance Education Technologies, 10(3), 17-34.

[16] Khdour, T., & Hijazi, H. (2012). A step towards preventive risk management in software projects. Proceedings of the International Conference on Software Technology and Engineering (ICSTE) (pp. 471-478). Phuket, Thiland: ASME. doi:10.1115/1.860151_ch75

[17] Hijazi, H., Khdour, T., & Alarabeyyat, A. (2012). A review of risk management in different software development methodologies. International Journal of Computer Applications, 45(7), 8-12. doi:10.5120/6790-9113

[18] Bratt, S. (2013). Agile instructional development framework: strategies for increasing learner and instructional designer collaboration. World Conference on Educational Multimedia, Hypermedia and Telecommunications. Victoria, British Columbia.

[19] Willeke, M. H. (2011). Agile in academics: applying agile to instructional design. 2011 Agile Conference, IEEE, (pp. 246-251).

[20] Cvetanovic, S. (2014). Development and enhancement of learning objects for e-learning systems using light agile method. The Fifth International Conference on e-Learning. Belgrade, Serbia.

[21] Albeanu, G. (2009). Agile CMMI for e-learning software development. Proceedings of the 5th International Scientific Conference eLSE "e-Learning and Software for Education", (pp. 135-142). Bucharest.

[22] Raja, W., & Nirmala, K. (2016). Agile development methods for online training courses web application development. International Journal of Applied Engineering Research, 11(4), 2601-2606.

[23] Uhomoibhi, J. O. (2006). E-learning and engineering education for sustainable development. 9th International Conference on Engineering Education. San Juan, PR.

[24] Fayyoumi, E., Idwan , S., AL-Sarayreh, K., & Obeidallah , R. (2015). E-learning: challenges and ambitions at Hashemite University. Int. J. Innovation and Learning, 17(4), 470-485.

[25] Bralić, A., & Divjak, B. (2018). Integrating MOOCs in traditionally taught courses: achieving learning outcomes with blended learning. International Journal of Educational Technology in Higher Education, 15(2), 1-16.

[26] Casanova, D., & Moreira, A. (2017). A model for discussing the quality of technology-enhanced learning in blended learning programmes. International Journal of Mobile and Blended Learning (IJMBL), 9(4), 1-20.

[27] Hollowell, G. P., Brooks, R. M., & Anderson, Y. B. (2017). Course design, quality matters training, and student outcomes. American Journal of Distance Education.

[28] Atoum, A., Al-Zoubi, A., Abu Jaber, M., Al-Dmour, M., & Hammad, B. (2017). A new approach for delivering e-Learning courses in Jordanian universities. Advances in Social Sciences Research Journal, 4(8), 1-13.

[29] Gómez-Rey , P., Barbera, E., & Fernández-Navarro, F. (2016). Measuring teachers and learners' perceptions of the quality of their online learning experience. Journal of Distance Education, 146-163.

[30] Ghislandi, P., Raffaghelli, J., & Yang, N. (2013). Mediated Quality: An approach for the e-Learning quality in higher education. International Journal of Digital Literacy and Digital Competence, 4(1), 56-73.

[31] Bremer, C. (2012). Enhancing e-learning Quality through the Application of the AKUE Procedure Model. Journal of Computer Assisted Learning, 28(1), 15–26.

[32] Ossiannilsson, E., & Landgren, L. (2012). Quality in e-learning – a conceptual framework based on experiences from three international benchmarking projects. Journal of Computer Assisted Learning, 28(1), 42–51.

[33] Lin, T.-C., & Chen, C.-J. (2012). Validating the satisfaction and continuance intention of e-learning systems: combining TAM and IS success models. 10(1).

[34] Sarsa , J., & Soler, R. (2012). E-Learning quality: relations and perceptions. International Journal of Information and Communication Technology Education, 8(2), 46-60.

[35] Marshall, S. (2012). Improving the quality of e-learning: lessons from the eMM. Journal of Computer Assisted Learning, 28(1), 65–78.

[36] Vázquez-Cano, E., & García, M. S. (2015). Analysis of risks in a learning management system: a case study in the Spanish National University of Distance Education (UNED). New Approaches In Educational Research, 4(1), 62-68.

[37] Surcel, T., & Reeiu, A. (2009). The risk management on developing the e-learning strategy. Proceedings 5th International Scientific Conference: e-Learning and Software for Education.

[38] Allen, S., & Hardin, P. C. (2008). Developing instructional technology products using effective project management practices. Journal of Computing in Higher Education, 19(2), 72-97.

[39] Andersson, A. (2008). Seven major challenges for e-learning in developing countries: case study eBIT, Sri Lanka. International Journal of Education and Development using Information and Communication Technology, 4(3), 45-62.

[40] Angelou, G. N., & Economides, A. A. (2007). E-learning investment risk management. Information Resource Management Journal, 20(4), 80-104.

[41] Rodríguez, J. C., Granados, J. J., & Muñoz, F. M. (2013). Engineering education through e-Learning technology in Spain. International Journal of Artificial Intelligence and Interactive Multimedia, 2(1).

[42] Benchicou, S., Aichouni, M., & Nehari, D. (2010). E-learning in engineering education: a theoretical and empirical study of the Algerian higher education institution. European Journal of Engineering Education, 35(3), 325-343.

[43] Chandra, A. P., & Samuel, R. S. (2010). E learning in engineering education: design of a collaborative advanced remote access laboratory. International Journal of Distance Education Technologies (IJDET), 8(2), 14-27.

[44] Rooij, V., & Williams, S. (2010). Project management in instructional design: ADDIE is not enough. *British Journal of Educational Technology, 41*(5), 852-864. doi:10.1111/j.1467-8535.2009.00982.x

[45] Martens, A., & Harrer, A. (2008). Software engineering in e-learning systems. In L. A. Tomei, Encyclopedia of Information Technology Curriculum Integration (pp. 782-789).

[46] Chimalakonda, S., & Nori, K. V. (2012). A software engineering perspective for accelerating educational technologies. IEEE 12th International Conference on Advanced Learning Technologies (ICALT).

[47] Arimoto, M. M., Barbosa, E. F., & Barroca, L. (2015). An agile learning design method for open educational resources. Frontiers in Education Conference (FIE), IEEE. El Paso, TX, USA.

[48] Hijazi, H., Alqrainy, S., Muaidi, H., & Khdour, T. (2014). A framework for ıntegrating risk management into the software development process. Research Journal of Applied Sciences, Engineering and Technology, 8(8), 919-928. doi:13329-RJASET-DOI

[49] Moreira, I. C., Ventura, S. R., Ramos, I., & Rodrigues, P. P. (2013). Learner's satisfaction within a breast imaging e-Learning course for radiographers. Proceedings of the IEEE 26th International Symposium; Computer-Based Medical Systems (pp. 215-220). Porto, Portugal.: IEEE.

[50] Tullis, T., & Albert , W. (2013). Measuring the user experience, second edition: collecting, analyzing, and presenting usability metrics. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc.

[51] Winer, L., Di Genova, L., Vungoc, P. - A., & Talsma, S. (2012). Interpreting end - of - course evaluation results. Montreal: Teaching and Learning Services, McGill University.

# Design and Development of an Industrial Solver for Integrated Planning of Production and Logistics

Yassine EL KHAYYAM[1]*, Brahim HERROU[2]

LTI, FST of Fez, University Sidi Mohamed Ben Abdallah, Fez, Morocco

*Abstract*—Faced with an increasingly hard competition, an increasingly unstable economic environment and ever-increasing customer requirements, companies should optimize costs and lead times not only at their level but also at the entire supply chain to which they belong. In such situation, an integrated supply chain management is necessary. In this paper, we discuss one of the essential building blocks of the integrated supply chain management, which is the integrated planning of the sup-ply chain. We introduce a new method for integrated planning of production and logistics, which is the MLRP (Manufacturing and Logistics Requirement Planning). This method allows supply chain planners to determine in advance, for the entire planning horizon, the manufacturing orders, the supplier's commands and the transport orders as well as vehicles routing for distribution of finished products and for the collection of components/raw materials. We will also discuss in this article the design and development of the solver which execute the MLRP method, this solver is the SMLRP that will be used to implement the proposed method on the different encountered industrial cases.

*Keywords*—MRP; VRP; production planning; transports planning; integration; MLRP

## I. INTRODUCTION

The evolution of information and communication technologies (ICT) has been a key factor in the industrial revolution lived in recent decades. Due to its two-dimensional advantages (the processing information dimension and the real-time information sharing/transferring dimension), ICT have contributed to the creation of global supply chains aimed at optimizing the costs and lead times throughout these supply chains. In this study, we are interested in planning systems as the fundamental element of any industrial management system and particularly supply chain management. The supply chain planning matrix, Fig. 1, considers that any supply chain could be divided into several internal supply chains, each of these supply chains is composed by four main processes the supply, the production, the distribution and the sales [1].

This article is both an implementation and a generalization of the MLRP system "Manufacturing and Logistic Requirement Planning" that we introduced in [2]. The generalization aims to extend the proposed model to determine simultaneously production requirements and logistics requirements whether upstream logistics (related to replenishment) or downstream logistics (related to distribution). The implementation, meanwhile, consists in designing and developing a solver allowing determining for each period belonging to the planning horizon the different orders of production and replenishment as well as the logistic

*\*Corresponding Author*

necessary means for the replenishment and / or distribution; it also allows generating vehicle routes and presenting them on geographical maps.

In this article, we first present the concepts used in the proposed approach, namely the MRP production systems and the vehicles routing problems, then we expose the new method of integrated planning of production and logistics requirements, the MLRP "Manufacturing and Logistic Requirement Planning". Next, we describe the MLRP solver from a software point of view (inputs / outputs, use case diagram and sequence diagram) and we conclude the article with a case study of a supply chain to expose the different features of the MLRP solver as well as its adaptability to different industrial configurations



Fig 1. Supply Chain Planning Matrix.

## II. MLRP: MANUFACTURING AND LOGISTIC REQUIREMENT PLANNING

In this section we present the MLRP method as an integrated planning solution for determining requirements of production and logistics, we will first introduce the different concepts that will be used by the MLRP which are the MRP production systems and vehicle routing problems.

### A. MRP Planning Systems: Concepts and Evolutions

Based on the bill of materials and the forecasted needs of finished products, on the delivery or manufacturing lead times of the various components. The MRP systems determines in advance the components manufacturing orders as well as the supply orders that must be transmitted to the production units and to the different suppliers [3], Fig. 2 presents the MRP process flow diagram [4].

Fig 2.    MRP Process Flow Chart.

As shown in Fig. 2, the MRP computation is based primarily on data collection: bill of material, order backlog and / or master production plan, quantities of available items and lead times for obtaining these items. Then from the top level, for each bill of material level, each item and at each period considered p, repeat:

*1)* Calculate the gross requirement GRp at the beginning of the period p, the gross requirement for an item Y of level n in the bill of material is the product of the estimated order of the higher-level item X in the bill of material and the assembly coefficient cm (Y) of the considered item:

$$GR_pY_n = PO_p(X, n-1) * cm(Y_n/X) \qquad (1)$$

*2)* Calculate the available items, AIp, at the beginning of the period according to equation (2), where FSp is the forecasted stock expected after transactions made during the period. OLp is the production order in progress or the supply or-der expected at the beginning of the period p.

$$AI_p = FS_{p-1} + OL_p \qquad (2)$$

*3)* Calculate net requirements, NRp, at the beginning of the period according to equation (3)

$$NR_p = Max(0, GR_p - AI_p) \qquad (3)$$

*4)* Define the proposed orders to satisfy net requirements by specifying items quantities and the launch dates.

*5)* Calculate the forecasted stock, FSp, at the end of the period according to equation (4), where OLp is a supply order under delivery and OPp is a production order.

$$FS_p = FS_{p-1} + OL_p + OP_p - GR_p \qquad (4)$$

As defined, MRP is very useful for determining accurately required orders quantities and required orders start dates of each component / raw materials for each period within the planning horizon. It is considered as a kernel of any MRP

production system, however it has some limitations, among which we cite [4]:

*a)* Lead times are considered deterministic, which is not the case in reality, production/supply lead times are subject to be changed because of the various hazards that may affect production units or suppliers.

*b)* The quality of the product supplied by production units or suppliers is considered perfect (no non-compliances or rejects), which is not the case in reality, in this situation production orders and supply orders must take into consideration the quality of the supplied products.

*c)* Capacities of production units or suppliers are not taken into consideration by the MRP when it is computed.

Several studies have been initiated a bout the MRP systems to overcome these limitations, as regards the failure to take into account lead times variability and products quality, several approaches have been proposed that address uncertainty in MRP systems based on stochastic inventory control [5], and fuzzy logic [6]. Otherwise, production capacities, supplier constraints, cost minimization and demand variability are taken into account as part of mathematical models. [7] Provides a mathematical model based on integer linear programming that takes into consideration, the scheduling of manufacturing / supply orders, production capacities, changes in production plans, storage conditions and storage costs.

In parallel, other research work focuses on dealing with MRP and transport management issues in an integrated way instead of being processed sequentially, [8] proposes a conceptual model that integrates both the aspects related to production planning and those related to transports management into a single model, The MRP IV (Fig. 3). [10] Proposes a linear programming model based on the MRP IV framework proposed by [8].



Fig 3.    MRP IV Framework.

Except the works of Mula et al. [8 and 9] that presents a conceptual model, MRP IV, which integrates transports planning and the MRP model, as well as those of Diaz-Madronero [10] relating to the proposal of a linear programming model which optimize simultaneously production and transports costs based on the MRP IV. The transports planning integration into the MRP model is not sufficiently addressed in the literature, to our knowledge, there is a lack of concrete models and information systems that ensure integrated planning of production and logistics.

### B. Vehicles Routing Problems: Concepts and Evolution

Vehicle routing problems were initially introduced by Dantzig and Ramser [11] to plan deliveries of fuel to gas stations, the basic version of vehicle routing problem, Fig. 4, consists in determining the optimal routes for a set of vehicles located in a depot to serve a set of customers [12].

Resolving vehicle routing problems is an important field of operational research and logistics management [13]; these problems are considered NP-hard because they cannot be solved in polynomial time [14].



Fig 4. Vehicle Routes Example.

Vehicle routing problems could be modeled as a complete graph G (V, E) with V= {$v_0$, $v_1$, $v_2$... $v_n$} is a set of nodes, v0 is the warehouse and the other vi are the destinations to be served by a fleet of m vehicles with a limited capacity $C_k$. E= [{$v_i$, $v_j$}, (i, j) ∈ A={(i, j) : i, j=0,1,2, …, n, i ≠j}] is the set of edges connecting the different nodes. Equations (5) - (9) represents mathematically the vehicle routing problem with limited capacity of vehicles "CVRP" (P1).

$$Minimize \quad \sum_{i=0}^{n}\sum_{j=0}^{n}\sum_{k=1}^{m} C_{ij}x_{ijk} \qquad (5)$$

$$\sum_{i=1}^{n} x_{i0k} = \sum_{j=1}^{n} x_{0jk} = 1 \quad \forall \in \{1,2,3,…,m\} \qquad (6)$$

$$\sum_{i=1}^{n}\sum_{k=1}^{m} x_{ijk} = \sum_{j=1}^{n}\sum_{k=1}^{m} x_{ijk} = 1 \quad \forall i,j \in \{1,2,3,…,n\} \qquad (7)$$

$$\sum_{i=1}^{n}\sum_{j=1}^{n} d_j x_{ijk} \leq C_k \quad \forall k \in \{1,2,3,…,m\} \qquad (8)$$

$$\sum_{k=1}^{m}\sum_{(i,j)\in NxN} x_{ijk} \leq |N| - 1 \quad \forall N \subset \{1,2,3,…,n\} \qquad (9)$$

Equation (5) represents the objective function that minimizes necessary costs to serve all customers from warehouse. Cij is the necessary cost to run through the edge (i, j), this one could be the amount of consumed fuel, the consumed time or the traveled distance between i and j. $x_{ijk}$ is a binary variable which is equal to 1 if the vehicle k run through the edge (i, j) and 0 if it isn't.

Equation (6) expresses the fact that each vehicle leaving the warehouse to serve one or some customers must return to the warehouse. Equation (7) adds the constraint that node demand must be served once; a single vehicle must serve each node.

Equation (8) ensures that the capacity of each vehicle will not be exceeded. Equation (9) ensures the elimination of sub tours that do not start and / or do not end at the warehouse.

The basic version of vehicle routing problem has been transformed into several versions by adding one or more constraints on the initial problem variables. These variables could be the vehicles capacities, the traveled distance, the beginning or arrival time, the delivery or recovery lead-time of the products from warehouse or from customers, there are various types of vehicle routing problems, some of which are below [15]:

- **MDVRP**: This variant takes into consideration several warehouses that can provide products instead of a single warehouse.

- **VRPTW**: A time window constraint for serving customers is added to the classic vehicle routing problem.

- **VRPB** deals with two customer's subgroups, some will have to be delivered from warehouse and the others will forward products to the warehouse.

- **OVRP**: This is to deal with the problems of vehicle routes in which the routes are not closed circuits, the vehicles do not return to the warehouse.

- **DVRP**: They are problems of vehicle routes in which the customer requests are known before (quantities) but can be formulated during the transport operation, these problems are considered dynamic routes problems.

- **Stochastic VRP** [16]: In this VRP variant, no information, about customer requests, is available before starting routes. Vehicle routes are planned based on probability distributions of customer requests.

The vehicle routing problem was addressed in the literature by two types of methods [17], Fig. 5, exact methods and approximate methods. The exact methods allow finding the optimal solution by exploring all possible solutions. However, the increase of the studied problem dimension (number of served clients and / or number of logistics means) makes exhaustive exploration of all solutions impossible in a sufficiently small duration.

The approximate methods, on the other hand, make it possible to find acceptable solutions but do not guarantee that the solution found is the optimal solution. There are two types of approximate methods, heuristics and metaheuristics. Heuristics are by definition a way to guide an algorithm to reduce the problem complexity; they are specific to a given problem. We distinguish three heuristics specific to the vehicle routing problems:

- Constructive heuristics are iterative algorithms in which at each iteration a partial solution is completed. The most popular constructive heuristic is the "savings" heuristic [18], it starts from an initial solution where each destination is served by a vehicle then we try to merge routes by computing for each pair of customers (vi, vj) the savings made by going from vi to vj instead of returning to the warehouse. Savings are then ordered and the corresponding customers to the largest saving are grouped together on the same route.

- The improvement heuristics try to improve a solution by proceeding to exchanges customers within routes. The exchanges could be carried out either within the same route or between customers being part of different routes [19].

- The two-phase methods consist of breaking down the vehicle routing problem into two sub-problems, one relating to the clients clustering and the other relating to determining an optimal route for each subgroup. According to the order in which the sub problems are treated, there are two methods, the Cluster First-Route second method and the Route first- Cluster second method.

Metaheuristics are advanced and powerful heuristics that could be applied to any optimization problem; they are divided into two categories: single-solution or local search metaheuristics and metaheuristic population-based solutions. Among the metaheuristics proposed in the literature, we cite below examples of local search metaheuristics, the simulated annealing and the taboo search, and example of metaheuristic population-based solutions, the evolutionary algorithms [20]:

- The taboo search: this method was formalized in 1986 by Glover [21], its principle is based on the mechanism inspired by human memory. It performs updates to an initial solution during successive iterations; during each iteration, the method constitutes a set of initial solution neighbor's by performing a single elementary movement. Then it evaluates the objective function value corresponding to the different neighbors obtained and substitutes the initial solution by the best solution founded even if the latter is bad than the initial solution, this helps to avoid local minimums. To avoid going back to a solution already obtained in previous iterations, this method uses a list of taboo movements that it avoids when forming the neighborhood, and it inserts in this list the movement corresponding to the solution obtained during each iteration.

- Simulated annealing: this method was inspired by metallurgist's techniques that are used to obtain a material with well-ordered molecules in solid state; Annealing involves heating a material to a very high temperature and then slowly lowering the temperature. Simulated annealing [22] applies this process to an optimization problem solution; the objective function is assimilated to the material energy that is subsequently minimized by introducing a fictive temperature. For each iteration, a basic modification is performed to the solution, if this modification implies a decrease of the objective function ($\Delta E \leq 0$) it is accepted, otherwise, it will be accepted with a probability equal to $e^{\frac{-\Delta E}{T}}$, T is a constant temperature until reaching the thermodynamic equilibrium. Once the equilibrium is reached, this temperature is reduced and the whole process is repeated until a reduced temperature is reached (cooled system).



Fig 5.   VRP resolution methods.

- Evolutionary algorithms are inspired by biological evolution of species; genetic algorithms [23] are the most popular evolutionary algorithms. They start from an initial solution population that they try to improve gradually over several generations by applying in a repetitive way the selection and reproduction principles. The selection principle consists to select the most suitable individuals for survival and reproduction (comparing the value of the corresponding objective function for each individual), and the principle of reproduction consists in mixing, recombining and changing (mutation) characteristics of solutions (parents) to form new solutions (descendants) with new potentialities.

### C. MLRP: The Integrated Planning System for Production and Logistics Requirements

We introduced in [2] the MLRP as a method for determining production and up-stream logistics requirements (for components and raw materials). MLR takes into account customer requests, restructured bill of materials (composition of the finished product adapted to the MLRP), the lead times (transport and production), dimensions and weights of the components, volume and the maximum weight of the available means of transport and transport costs per trip.

Our goal is to extend the initial version of the MLRP, so that it can simultaneously plan production needs and either upstream logistics needs (related to supply) and downstream logistics needs (related to distribution). The new version of the MLRP will also allow vehicle routes scheduling for distribution and replenishment. We will use the same restructured bill of material presented in [2] as well as the generator already developed in order to create step-by-step the bill of material from the highest level (finished product) to the lowest level of the components. The only change introduced in the bill of material is the addition of the raw material/component supplier identifier when this one is supplied and not manufactured or assembled.

Fig. 6 shows an example of bill of material adapted to the MLRP of a product A in which cm is the assembly coefficient, LT_p is the lead-time of production or supply, LT_l is the component transport lead-time. V and W are component volume and weight and finally SS and S_t0 are the security stock and the initial stock of a component.



Fig 6.    Bill of Material Example Adapted to the MLRP.

Fig. 7 illustrates the MLRP algorithm, this figure describes the MLRP first phase: the determination of manufacturing, supply and transport orders.

Once the transport orders have been determined, MLRP plans for each period of the planning horizon, the vehicle routes required for upstream logistics (recuperation of raw material and / or components) and for downstream logistics (distribution of finished products).

MLRP vehicle routing is based on the problem (P1) described by the equations (5) to (9) which describe the objective function and the constraints relating to the problems of routing heterogeneous vehicles with limited capacities (HCVRP), which is the case in the majority of the industrial problems.

As defined, (P1) does not consider nodes whose demand exceeds all available vehicle capacities (equation 9). To overcome this limitation, we proceed first to treat these nodes one by one trying to satisfy them by one or more trips using the available means of transport. The quantity Qit relative to node i (customer or supplier) and to period t, can be writ-ten according to equation (12) in which αk is the number of vehicles of type k or the number of trips that the type k vehicle will make. The parameter Cte is the rest of the Euclidean division of Qit on the vehicle's capacities (of all the possible combinations of the different vehicle types), Fk is the cost to make a trip using the Vk vehicle. We are looking in this situation for the combination of vehicles that generates the least cost, equation (10), while minimizing the rest of the Euclidean division, equation (11), the optimization in question can be formulated according to the problem (P0).

$$Minimize \sum_{k=1}^{K} \alpha_k F_k \qquad (10)$$

$$Minimize \ Q_{it} - \sum_{k=1}^{K} \alpha_k C_k \qquad (11)$$

$$Q_{it} = \sum_{k=1}^{K} \alpha_k C_k + Cte \qquad (12)$$

After the resolution of P0, we retain the different combinations selected for the nodes having initially requests higher than all the available capacities, we denote by VCSj the vehicle combination selected for the node j. then we replace the requests at these nodes by the rest of the Euclidean division of the initial request on the selected combination. Next, we solve the problem (P1 ') obtained by replacing the dj in problem (P1) by dj'.

$$d'_j = d_j - \sum_{k=1}^{K} \alpha_k C_k \qquad k\epsilon \ VCS_j(P0) \qquad (13)$$

Another limitation encountered in the vehicle routing for the supply chain upstream part is the heterogeneity of transported components, "we do not carry a single type of product as in the case of the distribution of finished products", in this situation, several cases are possible:

Fig 7.    MLRP Algorithm.

*a)* The different components / raw materials cannot be shipped together; in this case, the transport requirements are determined separately for each node by determining the combination of vehicles corresponding to the minimum transport cost. This approach is already detailed as part of our works introducing the MLRP [2].

*b)* The different components / raw materials can be shipped together, in this case we convert, equation (14), all the quantities at the nodes taking into consideration the equivalence between a component of a node and another component that we choose as reference (we choose the one with the smallest density and the smallest weight). Then we apply successively the two optimization steps defined by (P0) and (P1 ').

$$d''_j = k * d_j \quad \& \quad k = Min\left(\frac{v_b}{v_{ref}}, \frac{w_b}{w_{ref}}\right) \tag{14}$$

*c)* Some components can be embedded together and others cannot, in this case, we apply the approach a) for the first category and the approach b) on the second grouped category.

## III. MLRP SOLVER DESCRIPTION

In order to describe the MLRP solver (SMLRP) which ensure integrated planning of production and logistics requirements, we present in this section a static view and a dynamic view of the developed system. First, we represent the static view by the input / output diagram, and then we represent the dynamic view by a sequence diagram. Fig. 8 shows the inputs and outputs of the SMLRP, Outputs are the various manufacturing, supply and transport orders spread over the different periods of the planning horizon. The SMLRP also generates, for all periods, necessary vehicle routes for the supply chain upstream and downstream transports. The SMLRP inputs are bill of material adapted to the MLRP, customer orders spread over the planning horizon, types of available logistic means, locations of stakeholders (customers, suppliers, production unit).

Fig. 9 shows sequence diagram of the operation "Op_MLRP" which is the main operation of the SMLRP, this operation implements the diagram previously detailed in Fig. 7.

The sequence diagram shows the way in which "Op_MLRP" interacts with the other SMLRP objects to provide integrated planning for production and logistics requirements:

- OrderBacklogManagerImpl: This is a class that allows performing operations on the customer orders (database persistence, aggregations ...).

- ComponentManagerImpl: This manages bill of material components.

- MrpManagerImpl: This provides methods that insert and update manufacturing and transportation orders for each component at different times in planning horizon.

- Vc_mapManagerImpl: This object allows providing the capacity of each type of vehicle according to bill of material components.

- StakeholderManagerImpl: The Op_MLRP object interacts indirectly with this object through the operations planDistributionRouting and planRelenishmentRouting, these two operations plan vehicles routes by invoking JSPRIT, a dedicated java library for VRP.

- RouteManagerImpl: Ensures persistence of the JSPRIT provided solutions.

The JSPRIT library is an open source library developed in java, it implements local search metaheuristics to solve the travel salesman problem and vehicle routing problems. This library solves particularly problems of heterogeneous vehicles routing with limited capacities (HCVRP) which is the reference problem of the SMLRP, to which we converge each time by carrying out the necessary transformations for the various confronted situations (Section 2C).

We implemented the SMLRP Solver using a Java / J2EE architecture, the user accesses the solver through a web application that allows him to introduce the finished product bill of material, the master production plan and the logistic means. The solver presents production and transport orders in a summary table and plot the vehicle routes by period and by transport type (distribution or replenishment) on a geographical map. The SMLRP graphic user interfaces are presented in the case study exposed below.



Fig 8.    The MLRP inputs/outputs.

Fig 9. MLRP sequence diagram.

## IV. CASE STUDY

To expose the various SMLRP functionalities, we study in this section an example of a manufacturing supply chain of manual pallets trucks, Fig. 10. We will focus on a company that is a part of this supply chain, this company specializes in iron and steel metal structures manufacturing. This company manufactures the pallet truck metal components, acquires the other components from its suppliers, assembles all the components and packs the finished product to transport it to its customers.



Fig 10.  Hand Pallet truck.

To demonstrate the integrated planning of production and logistics requirements provided by the SMLRP solver, we cite below the inputs / outputs presented in Fig. 8, relating to the studied problem.

### A. The SMLRP Inputs for the Case Study

*1) The MLRP adapted bill of material for the case study*: As already mentioned in the previous paragraphs, the MLRP adapted bill of material, includes a panoply of information about the finished product structure, the component dimensions and component suppliers. Fig. 11 represents the hand pallet truck bill of material (HPT), it contains the assembly coefficients and the hierarchy of the different components, the data related to the production lead times ($LT\_p$), the initial stocks ($S\_t0$), the safety stocks ($SS$), the volumes ($V$) and the weights ($W$) of the components that need to be transported, as well as their transport lead time ($LT\_l$). Below the dimensions of the items that will be transported:



Fig 11.  Hand pallet truck bill of material.

- The dimensions (in mm) of the packed pallet truck are 1230x1250x1000, which is the equivalent of 1537 liters; its weight is 87 kg.

- The dimensions (in mm) of the packaged cylinder (HPT04-Cylinder) are 100x128x350, which is the equivalent of 4.5 liters; its weight is 11 kg.

- The dimensions (in mm) of the pallet scale (HPT02-WM) are 100x100x75, which is the equivalent of 0.75 liters; its weight is 1.2 kg.

- The dimensions (in mm) of the steering wheel (HPT06-Wheel) are 200x200x70, which is the equivalent of 2.8 liters; its weight is 4.5 kg.

- The dimensions (in mm) of the pallet truck rollers (HPT06-Roller) are 82x82x70, which is the equivalent of 0.47 liters; its weight is 1.05 kg.

*2) Stakeholders*: All the companies belonging to the studied supply chain are located in Morocco. The company EP is located in Casablanca, it has a main warehouse in the same place as its production unit, and its suppliers are:

- S1 is located in Berrechid, it provides the steering wheels.

- S2 is located in Kenitra, it provides the scale pallets.

- S3 is located in Tetouan, it provides the cylinders.

- S4 is located in Bouskoura, it provides pallet truck rollers.

The EP customers are scattered all over the country, the SMLRP solver allows introducing, with precision, stakeholders GPS coordinates. Fig. 12 represents the GUI, which allow introducing data related to different stakeholders (type, GPS coordinates and component provided ...).

*3) Types of available logistic means*: There are three available vehicle types RK, RMa and RMi. The RK can support a weight of 635kg, has a useful volume of 2600 liters, while the RMa can support a weight of 2100kg, and has a useful volume of 10800 liters, and finally the RMi can support a weight of 5500kg and has a useful volume of 40000 liters. Fig. 13 shows the GUI which allow entering data related to the different logistic means (capacities and costs).



Fig 12.  Stakeholders GUI.

**Logistics information:**

Type(*):
Used Volume(*):
Max Weight(*):
Cost/trip(*):

| Type ▲ | Used_Volume | Max_weight | Cost |
|---|---|---|---|
| RMi | 40000.0 | 5500.0 | 4500.0 |
| RK | 2600.0 | 635.0 | 700.0 |
| RMa | 10800.0 | 2100.0 | 1500.0 |

[riod] [Add Logistic Type] [Delete Logistic Type] [execute MRP] [Close]

Fig 13.   Logistic data GUI.

TABLE I.        EP CUSTOMER ORDERS

| | | Period | | | |
|---|---|---|---|---|---|
| | | **5** | **8** | **9** | **10** |
| Customer | **C1** | 20 | 30 | | |
| | **C2** | | 15 | | 40 |
| | **C3** | 30 | 5 | 10 | |
| | **C4** | | 25 | | 50 |
| | **C5** | 5 | | 20 | |
| | **C6** | 30 | | | 10 |
| | **C7** | 20 | | 40 | |
| | **C8** | 30 | 45 | | 50 |
| | **C9** | | | 20 | |

*4) Customer orders by period*: Table I shows the EP customer orders spread over several periods; these data are introduced in the SMLRP.

### B. The SMLRP Outputs for the Case Study

Once the studied problem data are inserted into the SMLRP solver, we execute MLRP algorithm in order to get the production, replenishment and transport orders planned over the planning horizon. The solver presents these results as a table (Fig. 14). The lines corresponding to the "Production ORD" are production orders for the components manufactured in-house, and sourcing orders from suppliers for components provided by other companies. The lines corresponding to the "Logistic ORD" are the transport orders. The solver also provides per period and depending on the operation concerned by the transport (finished product distribution or components collection) the necessary vehicles and the associated vehicle routes.

Fig. 15 shows as example the distribution route relating to the period 9, in this route, the solver indicates that it is necessary to provide five RMi vehicles, four of these vehicles will make a round trip (deposit-client) for customers C3, C5, C7 and C9. As the capacity of the RMi is 27 HPT, C3 and C7 customers cannot be delivered by one trip, the fifth RMi vehicle will make a route deposit-C3-C7- deposit in order to deliver the rest of the products for C3 and C7.

Fig. 16 shows as example the component replenishment route relating to the period 5, a single RMi vehicle is used in this route to recover the necessary components by performing the depot-S3-S2-S1-S4-depot course, the vehicle filling percentage is 56%.

| Period | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **HPT00** | | | | | | | | | | | | |
| Gross Req | 0 | 0 | 0 | 0 | 0 | 135 | 0 | 0 | 120 | 120 | 150 | 0 |
| On Hand Inventory | 30 | 30 | 30 | 30 | 30 | 30 | 20 | 20 | 20 | 20 | 20 | 20 |
| Net Req | 0 | 0 | 0 | 0 | 0 | 125 | 0 | 0 | 120 | 120 | 150 | 0 |
| Planned order receipt | 0 | 0 | 0 | 0 | 0 | 125 | 0 | 0 | 120 | 120 | 150 | 0 |
| Production ORD | 0 | 0 | 0 | 125 | 0 | 0 | 120 | 120 | 150 | 0 | 0 | 0 |
| Logistic ORD | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **HPT01-Fork** | | | | | | | | | | | | |
| Gross Req | 0 | 0 | 0 | 125 | 0 | 0 | 120 | 120 | 150 | 0 | 0 | - |
| On Hand Inventory | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 0 | - |
| Net Req | 0 | 0 | 0 | 125 | 0 | 0 | 120 | 120 | 150 | 0 | 0 | - |
| Planned order receipt | 0 | 0 | 0 | 125 | 0 | 0 | 120 | 120 | 150 | 0 | 0 | - |
| Production ORD | 0 | 125 | 0 | 0 | 120 | 120 | 150 | 0 | 0 | 0 | 0 | - |
| Logistic ORD | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | - |
| **HPT02-WM** | | | | | | | | | | | | |
| Gross Req | 0 | 0 | 0 | 125 | 0 | 0 | 120 | 120 | 150 | 0 | 0 | - |
| On Hand Inventory | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 0 | - |
| Net Req | 0 | 0 | 0 | 125 | 0 | 0 | 120 | 120 | 150 | 0 | 0 | - |
| Planned order receipt | 0 | 0 | 0 | 125 | 0 | 0 | 120 | 120 | 150 | 0 | 0 | - |
| Production ORD | 125 | 0 | 0 | 120 | 120 | 150 | 0 | 0 | 0 | 0 | 0 | - |
| Logistic ORD | 0 | 0 | 125 | 0 | 0 | 120 | 120 | 150 | 0 | 0 | 0 | - |
| **HPT03-Hand** | | | | | | | | | | | | |
| Gross Req | 0 | 0 | 0 | 125 | 0 | 0 | 120 | 120 | 150 | 0 | 0 | - |
| On Hand Inventory | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 0 | - |
| Net Req | 0 | 0 | 0 | 125 | 0 | 0 | 120 | 120 | 150 | 0 | 0 | - |
| Planned order receipt | 0 | 0 | 0 | 125 | 0 | 0 | 120 | 120 | 150 | 0 | 0 | - |
| Production ORD | 0 | 0 | 125 | 0 | 0 | 120 | 120 | 150 | 0 | 0 | 0 | - |
| Logistic ORD | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | - |
| **HPT04-Cylinder** | | | | | | | | | | | | |
| Gross Req | 0 | 0 | 0 | 125 | 0 | 0 | 120 | 120 | 150 | 0 | 0 | - |
| On Hand Inventory | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 0 | - |
| Net Req | 0 | 0 | 0 | 125 | 0 | 0 | 120 | 120 | 150 | 0 | 0 | - |
| Planned order receipt | 0 | 0 | 0 | 125 | 0 | 0 | 120 | 120 | 150 | 0 | 0 | - |
| Production ORD | 125 | 0 | 0 | 120 | 120 | 150 | 0 | 0 | 0 | 0 | 0 | - |
| Logistic ORD | 0 | 0 | 125 | 0 | 0 | 120 | 120 | 150 | 0 | 0 | 0 | - |
| **HPT05-supportC** | | | | | | | | | | | | |
| Gross Req | 0 | 0 | 0 | 125 | 0 | 0 | 120 | 120 | 150 | 0 | 0 | - |
| On Hand Inventory | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 0 | - |
| Net Req | 0 | 0 | 0 | 125 | 0 | 0 | 120 | 120 | 150 | 0 | 0 | - |
| Planned order receipt | 0 | 0 | 0 | 125 | 0 | 0 | 120 | 120 | 150 | 0 | 0 | - |
| Production ORD | 0 | 0 | 125 | 0 | 0 | 120 | 120 | 150 | 0 | 0 | 0 | - |
| Logistic ORD | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | - |
| **HPT06-wheel** | | | | | | | | | | | | |
| Gross Req | 0 | 0 | 0 | 250 | 0 | 0 | 240 | 240 | 300 | 0 | 0 | - |
| On Hand Inventory | 0 | 0 | 0 | 0 | 20 | 20 | 20 | 20 | 20 | 20 | 0 | - |
| Net Req | 0 | 0 | 0 | 270 | 0 | 0 | 240 | 240 | 300 | 0 | 0 | - |
| Planned order receipt | 0 | 0 | 0 | 270 | 0 | 0 | 240 | 240 | 300 | 0 | 0 | - |
| Production ORD | 270 | 0 | 0 | 240 | 240 | 300 | 0 | 0 | 0 | 0 | 0 | - |
| Logistic ORD | 0 | 0 | 270 | 0 | 0 | 240 | 240 | 300 | 0 | 0 | 0 | - |
| **HPT07-Roller** | | | | | | | | | | | | |
| Gross Req | 0 | 0 | 0 | 500 | 0 | 0 | 480 | 480 | 600 | 0 | 0 | - |
| On Hand Inventory | 0 | 0 | 0 | 0 | 40 | 40 | 40 | 40 | 40 | 40 | 0 | - |
| Net Req | 0 | 0 | 0 | 540 | 0 | 0 | 480 | 480 | 600 | 0 | 0 | - |
| Planned order receipt | 0 | 0 | 0 | 540 | 0 | 0 | 480 | 480 | 600 | 0 | 0 | - |
| Production ORD | 540 | 0 | 0 | 480 | 480 | 600 | 0 | 0 | 0 | 0 | 0 | - |
| Logistic ORD | 0 | 0 | 540 | 0 | 0 | 480 | 480 | 600 | 0 | 0 | 0 | - |
| **HPT08-sideBars** | | | | | | | | | | | | |
| Gross Req | 0 | 0 | 0 | 500 | 0 | 0 | 480 | 480 | 600 | 0 | 0 | - |
| On Hand Inventory | 0 | 0 | 0 | 0 | 40 | 40 | 40 | 40 | 40 | 40 | 0 | - |
| Net Req | 0 | 0 | 0 | 540 | 0 | 0 | 480 | 480 | 600 | 0 | 0 | - |
| Planned order receipt | 0 | 0 | 0 | 540 | 0 | 0 | 480 | 480 | 600 | 0 | 0 | - |
| Production ORD | 0 | 540 | 0 | 0 | 480 | 480 | 600 | 0 | 0 | 0 | 0 | - |
| Logistic ORD | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | - |

Fig 14.   Executed MLRP for the EP Company.



Fig 15.   Distribution route for period 9.

Fig 16.   Component replenishment route relating to the period 5.

## V. CONCLUSION

In this paper, we presented the implementation and extension of the integrated planning method for production and logistics needs, the MLRP, which already have been introduced in our previous works. This method allows planning managers to determine in advance, for the entire planning horizon, production orders, orders to be forwarded to suppliers and quantities to be transported upstream and downstream of the supply chain. By dint of the Java / J2EE technology we have implemented the SMLRP solver allowing to execute the MLRP algorithm and we have also benefited from the progress lived in recent decades in the field of vehicles routing problems to integrate the generation of vehicle routes within the MLRP algorithm execution.

### REFERENCES

[1] Stadtler H. et al, Supply Chain Management and Advanced Planning, Springer Texts in Business and Economics, DOI 10.1007/978-3-642-55309-7__4, © Springer-Verlag Berlin Heidelberg 2015.

[2] Y. El Khayyam, B. Herrou, (2018). Integrated production and logistics planning. International Journal of Online Engineering. 14. 84-96. 10.3991/ijoe. v14i12.9068.

[3] L.Boyer et al, Précis d'organisation et de gestion de la production, Éditions d'Organisation, 1982.

[4] P. Vrat, Just in Time, MRP, and Lean Supply Chains, Materials Management, Springer Texts in Business and Economics 2014.

[5] Barba-Gutiérrez, Y., Adenso-Díaz, B.: Reverse MRP under uncertain and imprecise demand. The International Journal of Advanced Manufacturing Technology 40(3), 413–424 (2009)

[6] J. Mula et al, MRP with flexible constraints: A fuzzy mathematical programming approach. Fuzzy Sets and Systems 157(1), 74–97 (2006).

[7] R. Karni, Integer linear programming formulation of the material requirements planning problem. Journal of Optimization Theory and Applications 35(2), 217–230 (1981).

[8] J.Mula et al, A Conceptual Model for Integrating Transport Planning: MRP IV, IFIP International Federation for Information Processing pp. 54–65, 2012.

[9] N. Absi, S. Kedad-Sidhoum, The multi-item capacitated lot-sizing problem with setup times and shortage costs. European Journal of Operational Research 185(3), 1351–1374 (2008).

[10] M. Díaz-Madroñero et al, A mathematical programming model for integrating production and procurement transport decisions, Applied Mathematical Modelling 52 (2017) 527–543.

[11] R. Dantzig, J. Ramzer (1959). The truck-dispatching problem. Management Science 6(1):81–91

[12] A.O. Adewumi & O.J. Adeleke, A survey of recent advances in vehicle routing problems, International Journal of System Assurance Engineering and Management (2018) 9: 155. https://doi.org/10.1007/s13198-016-0493-4

[13] C. Wujun, Y. Wenshui (2017), A Survey of Vehicle Routing Problem, MATEC Web of Conferences 100, 0100.

[14] JK. Lenstra, AHG. Rinnooy Kan (1981) Complexity of vehicle routing and scheduling problems. Networks 11(2):221–227.

[15] A. Dixit, & A. Mishra & A. Shukla, (2019). Vehicle Routing Problem with Time Windows Using Meta-Heuristic Algorithms: A Survey: Theory and Applications, ICHSA 2018. 10.1007/978-981-13-0761-4_52.

[16] R. Maini & R. Goel, (2017). Vehicle routing problem and its solution methodologies: a survey. International Journal of Logistics Systems and Management. 28. 419. 10.1504/IJLSM.2017.10008188.

[17] B.I. Sahbi et al (2011), Synthèse du problème de routage de véhicules, Collection des rapports de recherche de Telecom Bretagne.

[18] G.Clarke and JW Wright. Scheduling of vehicles from a central depot to a number of delivery points. Operations research, 12(4) :568-581, 1964.

[19] T.G. Crainic and F. Semet. Recherche opérationnelle et transport de marchandises. In Optimisation combinatoire. 3, Applications. Hermès Science : Lavoisier, 2006.

[20] P. Siarry (ed.), Metaheuristics, Springer International Publishing Switzerland 2016, DOI 10.1007/978-3-319-45403-0_1.

[21] Glover, F.: Future paths for integer programming and links to artificial intelligence. Computers and Operations Research 13(5), 533–549 (1986).

[22] Kirkpatrick, S., Gelatt, C., Vecchi, M.: Optimization by simulated annealing. Science 220(4598), 671–680 (1983)

[23] Goldberg, D.E.: Genetic Algorithms in Search, Optimization and Machine Learning. Addison- Wesley (1989)

# Novel Joint Subcarrier and Power Allocation Method in SWIPT for WSNs Employing OFDM System

Saleemullah Memon[1]

School of Information & Communication Engineering
Beijing University of Posts and Telecommunications
Beijing, China

Kamran Ali Memon*[2]

State Key Laboratory of Information Photonics and Optical
Communications (IPOC), School of Electronic Engineering
Beijing University of Posts and Telecommunications
Beijing, China

Zulfiqar Ali Zardari[3]

Faculty of Information Technology
Beijing University of Technology, Beijing, China

Muhammad Aamir Panhwar[4]

School of Electronic Engineering
Beijing University of Posts and Telecommunications
Beijing, China

Sijjad Ali Khuhro[5]

School of Computer Science and technology
University of Science and technology of china
China

Asiya Siddiqui[6]

School of Information & Communication Engineering
Beijing University of Posts and Telecommunications
Beijing, China

*Abstract*—In recent research trends, simultaneous wireless information and power transfer (SWIPT) has proved to be an innovative technique to deal with limited energy problems in energy harvesting (EH) technologies for wireless sensor networks (WSNs). In this paper, a method of subcarrier and power allocation for both EH and information decoding (ID) operations is proposed under orthogonal frequency division multiplexing (OFDM) systems, with an improved the quality of service (QoS) parameters. This proposed method assigns one group of subcarriers for ID and remaining group of subcarriers is assigned for EH, despite of applying any splitting schemes. We achieved maximum EH under the ID and power constraints with an effective algorithm for the first time incorporating the dual decomposition technique which deals with power and subcarrier allocation problem jointly. The obtained simulation outcomes in relation to power allocation ratio, subcarrier allocation ratio and energy harvested (EH) at the destination node proved better when compared to the schemes that contain water filling, time switching (TS) or power splitting (PS) approaches under different target transmission rates and transmitter and receiver distances.

*Keywords—Simultaneous wireless information & power transfer (SWIPT); Energy harvesting (EH); Information decoding (ID); power allocation; subcarrier allocation*

## I. INTRODUCTION

The concept of simultaneous transformation of information and energy in the field of wireless communication has recently received an upsurge of interest for the researchers. Through this simultaneous wireless information and power transfer (SWIPT), mobile phone users have access to more than just energy but data as well at the same time, which provides huge prospects. However, some ultimate strategic variations are needed in wireless communication networks to bring SWIPT for efficient applications [1]. The rates of transmitting information and the reliability of reception are suitably applied to estimate the wireless systems performance [2].

In [2], an impractical receiver is considered, which possesses the capability to achieve energy harvesting (EH) and information decoding (ID) concurrently. Furthermore, the agreement between harvesting energy and information rate level are essential to determine the system pursuance when the receiver carryout harvested energy using radio frequency (RF) signal [2]. In [3] and [4], however, two constructive SWIPT power splitting (PS) and time switching (TS) strategies were considered. The PS scheme divides the signal received by receiver into two parts with distinct powers, one part is assigned for ID operation and the other part is assigned for EH operation. Whereas in the TS strategy, the receiver automatically switches within one transmission time to EH or ID mode. In addition to the receiver circuit, the PS design of conventional communication systems [5] does not require any further changes. In [6], the major problem of efficient energy optimization in the multi-input-single-output (MISO) downlink scheme based on the PS scheme was considered. This method used langrangian relaxation in conjunction with the dinkelbach method, which maximized the energy efficiency under the constraints of EH and signal-to-interference-plus-noise-ratio (SINR).

Preliminary research is based on point-to-point multi-input-multi-output (MIMO) connections that actually communicate with each other via two devices with numerous antennas. However, current focus has indeed moved to multiuser MIMO systems, where a group of single-antenna users are concurrently assisted by the base station equipped with multi-antennas [7]. In the classic work of [8], SWIPT was investigated by using TS and PS schemes in massive MIMO

*Corresponding Author

which can allow multi-way relay networks (MWRNs) and energy restricting amplify and forward relays.

Orthogonal frequency division multiplexing (OFDM) is a commonly used multi-carrier method in different wireless standards. The overall performance of OFDM based SWIPT mechanisms has been studied in many existing papers [9]-[10]. In order to analyze the performance of SWIPT, [9] selected an OFDM based multi user single-antenna system with a PS technique at the receiver. The results conclude that the energy efficiency of the system can be greatly improved by using RF EH in the limited interference regime.

Multi-antenna receivers are useful for improving the system potential capacity rather than improving the energy efficiency of the system [9]. Authors in [11] focused on SWIPT for broadband wireless systems (BWSs) involving OFDM and beamforming transition. This can develop a number of parallel sub channels to simply the mechanism of the resource allocation. Carriers in [11] are used for EH proposed for that specific user by employing a fixed subcarrier allocation. However, when it comes to multi-user wireless systems with OFDM, authors in [12] researched on the phenomenon of multi distributed users receiving a broadcasted message from a single fixed access point. The authors in [12] considers two multiple access schemes in their analysis, notably time division multiple access (TDMA) and multiple orthogonal frequency division access (OFDMA). TDMA employs TS approach where information receiver for a user works for scheduled period whereas its energy receiver operates for all the periods. However, OFDMA assumes that available subcarriers share the same PS ratio at the receiver. Author in [10] presents the analysis of SWIPT receivers with a single OFDM user channel with an upper limit for the rate of energy tradeoff system under consideration. Simulation results conclude the trade-off between the achievable rate and the harvested energy which serves to be a well-known trade-off between the rate and energy of the system.

In WSNs, small nodes are jot usual near to each other which may adversely affect each other due to urban, climatic or geographical obstacles. As a result, Line of sight (LOS) communication becomes impracticable, therefore [13] and [14] suggests to use idle nodes for creating intermediate hops. Relaying consumes some extra energy to serve the purpose of cooperative transmission. Relaying protocols work in two fashions i.e. a) Decode then forward or regenerate and b) Amplify then forward or transparent relaying as reported in [15] and [16]. [17] analyzes a three-node scheme and comes up with an efficient resource allocation in terms of throughput and energy usage. The same is repeated for multiple access relay arrangement, as found in [18] and [19]. [20] extends it with multi antenna use and utilizes harvested energy not for storing in batteries but for the operational needs. Furthermore, [20] presents a tradeoff for the EH time and the time period relay node is communicating. Author in [21] introduces the methods to reduce interference in RF-EH operation for such mentioned relaying systems.

Lot of research had been focused on SWIPT in an OFDM comprised of PS or TS schemes, where destination node requires a time or power divider to distinguish the signal for ID operation and EH operation. In our work, we have considered SWIPT in joint power and subcarrier allocation-based architecture in an OFDM system, where have not considered any separate splitter/divider at the receiving end. Particularly, the OFDM based subcarriers are partitioned in two portions. One portion is allocated for one group that is used for EH and the other potion is allocated for another group that is used for ID of the received signal. Hence, the only job of receiving node is to consider which group is assigned for the operation of EH, at that time other group will be assigned for the operation of ID. Thus, at the receiving node no splitter/divider is needed anymore. The key contributions of this work can be labelled as:

- EH is maximized under the allocation of joint power and subcarrier criteria, ensuring the achievable target transmission rates.

- The complexity of receiver is minimized in proposed work and no splitter is used as compared to the previous work of SWIPT in OFDM schemes.

- Simulation results validate the dominance of our proposed joint power and subcarrier allocation scheme when compared to the schemes that contain water filling or PS/TS approaches.

The remaining organization of this paper is given as follows. Section 2 presents the system model of our proposed work. Section 3 explains the problem formulation of EH optimization and joint optimal power and subcarrier allocation problems. Section 4 provides the solution of the optimization problems. Section 5 presents the performance of our proposed work. Finally, Section 6 concludes the overall research work and contains future work as well.

## II. SYSTEM MODEL

A wireless link based on OFDM, consists of a single antenna transmitter (Tx) and a single antenna receiver (Rx), as revealed in Fig. 1. We represented all the subcarriers as I = {1, 2, ... , J} and the total bandwidth of the system is divided equally into J subcarriers.



Fig. 1. Wireless SWIPT Architecture.

We have assumed that channel power gain is perfectly known at the transmitter and is indicated as $g_J$. Let $P_t$ denotes the total transmitted power to all J subcarriers and $p_J$ indicates the power allocated to subcarrier J. We have assumed a slow fading situation where all the constant channels coefficients are considered. A bandpass filter is assumed at the receiving node. On each J subcarrier, the received signal corrupted by the noise $n_J$ is designed as an additive-white-Gaussian-noise (AWGN) having 0 mean with a variance of $\sigma_J^2$ and is given as $n_J \sim \mathcal{CN}(0, \sigma_J^2)$.

### III. ENERGY HARVESTING OPTIMIZATION PROBLEM

Based on the SWIPT framework, the signal received at Rx is partitioned into two groups, one group is for the process of ID and the other is for the process of EH operation over all the subcarriers, as given below.

$$I = K_{ID} + K_{EH} \tag{1}$$

where $K_{ID} \leq I$ and $K_{EH} \leq I$ denote the subcarriers used for ID and EH respectively. The EH at the Rx node is given as

$$E = \xi \sum_{J \in K_{EH}} (p_J g_J + \sigma_J^2) \tag{2}$$

where $\xi$ represents the conversion efficiency of energy at Rx, it is assumed that $\xi = 1$, for convenience. For the OFDM link, transmission rate can be achieved as

$$R_T = \sum_{J \in K_{ID}} \log_2(1 + p_J g_J / \sigma_J^2) \tag{3}$$

In order to satisfy the constraints of power and target date rate requirement, we aim to maximize the EH under joint subcarrier and power allocation. Thus, the optimization problem can be formulated as

$$\max \quad E$$
$$\text{s.t} \quad R_T \geq \Psi \tag{4}$$
$$\sum_{J \in K_{EH}} p_J + \sum_{J \in K_{ID}} p_J \leq P_t$$

where $\Psi$ denotes the minimum target transmission rate requirement.

### IV. OPTIMIZATION SOLUTION

The EH optimization problem mentioned above is a nonconvex mixed integer problem; hence it is infeasible to find the direct solution of such problems because of high complexity. However, considering large number of subcarriers, the timing sharing condition can be applied [22]-[23], which makes the duality gap to zero. Thus, the optimization problem can be solved using dual decomposition method. The Lagrange dual function is given as

$$\Gamma(\lambda) = \max_{\{\mathbb{P}, \mathbb{K}\}} \mathcal{L}(\mathbb{P}, \mathbb{K}) \tag{5}$$

where

$$\mathcal{L}(\mathbb{P}, \mathbb{K}) = E + \lambda_1(R_T - \Psi) +$$
$$\lambda_2(P_t - \sum_{J \in K_{EH}} p_J - \sum_{J \in K_{ID}} p_J) \tag{6}$$

$\mathbb{P} = \{p_J\}$ and $\mathbb{K} = \{K_{ID}, K_{EH}\}$ denote the power and subcarrier allocation set respectively and $\lambda = (\lambda_1, \lambda_2)$ denotes the dual variable vector having positive value subject to rate and power constraints. The dual optimization problem can be obtained as

$$\min_\lambda \quad \Gamma(\lambda) \tag{7}$$
$$\text{s.t} \quad \lambda \geq 0$$

As the dual function is a convex as proved in [11], so in order to minimize the lagrange dual function $\Gamma(\lambda)$, we employ sub-gradient based method that will assure the convergence. The sub-gradient method is given as

$$\Delta\lambda_1 = R_T - \Psi \tag{8}$$
$$\Delta\lambda_2 = P_t - \sum_{J \in K_{EH}} p_J - \sum_{J \in K_{ID}} p_J \tag{9}$$

$\Gamma(\lambda)$ cannot be obtained without optimal $\mathbb{P}$ and $\mathbb{K}$ at given values of $\lambda$. We define a two-step process in order to obtain optimal $\mathbb{P}$ and $\mathbb{K}$. In the first step we achieve the optimal $\mathbb{P}$ while fixing, in the later step we find the optimal $\mathbb{K}$.

#### A. Optimal $\mathbb{P}$ while fixing $\mathbb{K}$

For a fixed $\mathbb{K}$, partially derivate (6) with respect to variables of optimization problem $p_J$, can be expressed as follow

$$\frac{\partial \mathcal{L}(\mathbb{K}, \mathbb{P})}{\partial p_J} = \xi g_J - \lambda_2, \quad J \in K_{EH} \tag{10}$$

$$\frac{\partial \mathcal{L}(\mathbb{K}, \mathbb{P})}{\partial p_J} = \frac{\lambda_1 g_J}{p_J g_J + \sigma_J^2} - \lambda_2, \quad J \in K_{ID} \tag{11}$$

After applying the Kuhn-Karush-Tucker (KKT) conditions, the partial derivatives of the Lagrange tend to zero, hence the optimal $p_J$ ($J \in K_{EH}$) for a given $\lambda$, can be expressed as

$$p_J^* = p_{max} = \xi g_J > \lambda_2 \tag{12}$$

$$p_J^* = p_{min} = \xi g_J \leq \lambda_2 \tag{13}$$

The optimal $p_J$ ($J \in K_{ID}$) for a given $\lambda$, can be expressed as

$$p_J^* = \left(\frac{\lambda_1}{\lambda_2} - \frac{\sigma_J^2}{g_J}\right)^+ \tag{14}$$

where, $P_{max}$ and $P_{min}$ are maximum and minimum constraints of power respectively on each subcarrier.

#### B. Optimal $\mathbb{K}$

After achieving optimal P, optimal K can be obtained by substituting the optimal (12), (13) and (14) in (6), after rearranging, the Lagrange can be rewritten as

$$\mathcal{L}(\mathbb{K}) = \sum_{J \in K_{ID}} \left(\lambda_1 \log_2(1 + p_J^* g_J / \sigma_J^2) - \xi(p_J^* g_J + \sigma_J^2)\right) +$$
$$\sum_{J=1}^{J} (\xi(p_J^* g_J + \sigma_J^2) - \lambda_2 p_J^*) - \lambda_1 \Psi +$$
$$\lambda_2 P_t \tag{15}$$

From (15), we can be seen that only right-side portion of the equation contains the subcarriers for ID, hence separating (15) as,

$$F_J^* = \lambda_1 \log_2(1 + p_J^* g_J / \sigma_J^2) - \xi(p_J^* g_J + \sigma_J^2) \tag{16}$$

The resultant set of optimum subcarriers $K_{ID}$ that will maximize the Lagrange dual function. Consequently, the optimum $K_{ID}$ can be obtained as

$$K_{ID}^* = \arg\max_{K_{ID}} \sum_{J \in K_{ID}} F_J^* \qquad (17)$$

From the set subcarriers, the remaining optimal $K_{EH}$ can be obtained as

$$K_{EH}^* = I - K_{ID}^* \qquad (18)$$

Hence, the primal optimization variables $\mathbb{P}$ and $\mathbb{K}$ are successfully achieved with the help of specified dual variables. Therefore, the mentioned optimum problem in (4) can now be completely solved by the process of updating the values of primal dual variables.

## V. Simulation Results and Discussions

This section presents the simulation results and improvements achieved for joint resource allocation in an OFDM based SWIPT architecture as compared to previous research works. We consider power allocation ratio, subcarrier allocation ratio and energy harvested (EH) at the destination node as the evaluation matrices. For the fading channel of frequency selection, we have used 6 taps and kept central frequency at 1.9 GHz. This research in this paper is limited to small scale fading scenarios. The primary role is the LOS signal then the Rician fading channel is designed. Particularly, for all the subcarriers the channel modeling is given as

$$f(J) = \sqrt{\frac{N}{1+N}}\, \tilde{f} + \sqrt{\frac{1}{1+N}}\, \hat{f}(J) \qquad (19)$$

The channel power-gain is represented as $g_J = |f(J)|^2$. The limitations and complete list of parameters used during our simulations are presented in the Table I.

Fig. 2 shows the harvested energy (EH) vs. power transmitted ($P_t$) at different target transmission rates ($\Psi$). We can observe that as the requirements of $\Psi$ increases in result less energy is harvested. The reason is that with high target rate more power is needed for decoding the ID process wherever the transmitted power is fixed. Consequently, less power will remain for EH operation.

TABLE I. Variables and Values of the Parameters Used

| LOS deterministic component | $\tilde{f}$ |
|---|---|
| Rayleigh fading component | $\hat{f}(J)$ |
| Rician fading channel | $f(J)$ |
| No. of subcarriers (I) | 32 |
| Rician factor (N) | 3 |
| Noise spectrum density | -45dBm |
| Energy Conversion Efficiency ($\xi$) | 100% ($\xi$=1) |
| Tx & Rx distance, d | 4 m |
| Target transmission rate ($\Psi$) | 5 bps/Hz |



Fig. 2. EH vs. $P_t$ at different $\Psi$.



Fig. 3. Subcarrier/Power Allocation Ratio vs. Tx and Rx Distance.

Fig. 3 shows the variations in allocation of power ratios and allocation of subcarrier ratios in relation to separation distance between Tx and Rx, where the value of transmitted power is fixed to 0.5 W. We can observe from Fig. 3 that as the distance between the Tx and Rx increase, the more subcarriers and powers are allocated for decoding the information process and less subcarriers and powers are allocated for harvesting the energy process at the same time. This is because, when the distance increases, the channel between the Tx and Rx deteriorates, consequently more resources are assigned for ID operation to meet with the fixed target rate. Thus, few subcarriers and less powers are left for EH operation.

Fig. 4. Subcarrier and Power Allocation at Ψ=5bps/Hz.

Fig. 4 display the joint resource (power/subcarrier) allocations at target transmission rate $\Psi = 5$bps/Hz when $P_t = 0.5$ W and d=0.6m. We can conclude that most portion of the resources are assigned for the operation of EH. The reason is that low target rate values require very low portions of resource allocation in order to take part for the ID process.

## VI. CONCLUSION AND FUTURE WORK

We have proposed a joint power and subcarrier allocation method based on SWIPT framework under OFDM architecture, where a separate splitter is not considered at the receiving node. Particularly, the OFDM subcarriers are partitioned in two portions. One portion is allocated for one group that is used for EH and the other potion is allocated for another group that is used for ID of the received signal. The receiving node only needs to consider which group is allocated for EH, then the other group will be allocated for ID. EH is maximized under the allocation of joint power and subcarrier criteria, though ensuring the achievable target transmission rate. Simulation results show that advantages of our proposed SWIPT OFDM strategy and reveal the dominance of our proposed joint power and subcarrier allocation scheme compared to the schemes that contain water filling or PS/TS approaches.

Furthermore, for future work, researchers can investigate some appropriate techniques to permit the operation of vacant subcarriers for EH and ID. In this situation, multiple frequency bands would be needed so that SWIPT receiver can be able to perform the suitable operations.

### REFERENCES

[1] T. Perera, D. N. Jayakody, S. K. Sharma, S. Chatzinotas, and J. Li, Simultaneous Wireless Information and Power Transfer (SWIPT): Recent Advances and Future Challenges, vol. PP. 2018.

[2] L. R. Varshney, Transporting Information and Energy Simultaneously. 2008.

[3] L. Liu, R. Zhang, and K. Chua, Wireless Information Transfer with Opportunistic Energy Harvesting, vol. 12. 2012.

[4] L. Liu, R. Zhang, and K. Chua, Wireless Information and Power Transfer: A Dynamic Power Splitting Approach, vol. 61. 2013.

[5] X. Zhou, Training-Based SWIPT: Optimal Power Splitting at the Receiver, vol. 64. 2014.

[6] Q. Shi, C. Peng, W. Xu, M. Hong, and Y. Cai, Energy efficiency optimization for MISO SWIPT systems with zero-forcing beamforming, vol. 64. 2015.

[7] lu lu, G. Li, A. Swindlehurst, A. Ashikhmin, and R. Zhang, An Overview of Massive MIMO: Benefits and Challenges, vol. 8. 2014.

[8] R. Zhang and C. K. Ho, MIMO Broadcasting for Simultaneous Wireless Information and Power Transfer., vol. 12. 2011.

[9] D. W. K. Ng, E. S. Lo, and R. Schober, Wireless Information and Power Transfer: Energy Efficiency Optimization in OFDMA Systems, vol. 12. 2013.

[10] P. Grover and A. Sahai, Shannon meets Tesla: Wireless information and power transfer. 2010.

[11] K. Huang and E. G. Larsson, Simultaneous Information and Power Transfer for Broadband Wireless Systems, vol. 61. 2012.

[12] X. Zhou, R. Zhang, and C. K. Ho, Wireless Information and Power Transfer in Multiuser OFDM Systems, vol. 13. 2013.

[13] M. Peng, Y. Liu, D. Wei, W. Wang, and H.-H. Chen, Hierarchical cooperative relay based heterogeneous networks, vol. 18. 2011.

[14] B. Zhou, H. Hu, S.-Q. Huang, and H.-H. Chen, Intracluster Device-to-Device Relay Algorithm With Optimal Resource Utilization, vol. 62. 2013.

[15] J. Nicholas Laneman and G. W. Wornell, Energy-Efficient Antenna Sharing and Relaying for Wireless Networks, vol. 1. 2000.

[16] J. Nicholas Laneman, D. N. C. Tse, and G. W. Wornell, Cooperative Diversity in Wireless Networks: Efficient Protocols and Outage Behavior, vol. 50. 2005.

[17] H. (Henry Chen, Y. Li, J. Luiz Rebelatto, B. F. Uchoa-Filhoand, and B. Vucetic, Harvest-Then-Cooperate: Wireless-Powered Cooperative Communications, vol. 63. 2014.

[18] A. Rajaram, D. N. Jayakody, and V. Skachek, Store-then-cooperate: Energy harvesting scheme in cooperative relay networks. 2016.

[19] W. Huang, H. (Henry Chen, Y. Li, and B. Vucetic, On the Performance of Multi-Antenna Wireless-Powered Communications with Energy Beamforming, vol. 65. 2015.

[20] I. Krikidis, G. Zheng, and B. Ottersten, Harvest-use cooperative networks with half/full-duplex relaying. 2013.

[21] A. Nasir, X. Zhou, S. Durrani, and R. Kennedy, Relaying Protocols for Wireless Energy Harvesting and Information Processing, vol. 12. 2012.

[22] Y. Liu and X. Wang, Information and Energy Cooperation in OFDM Relaying: Protocols and Optimization, vol. 65. 2015.

[23] W. Yu and R. Lui, Dual Methods for Nonconvex Spectrum Optimization of Multicarrier Systems., vol. 54. 2006.

### AUTHOR'S PROFILE

**Saleemullah Memon** received his B.E degree in electronic engineering from Quaid-e-Awam University of Engineering, Science and Technology (QUEST), Pakistan, in 2017. He is currently pursuing his MS degree at Key Laboratory of Universal Wireless Communication (Ministry of Education), Beijing University of Posts and Telecommunications (BUPT), China. His current research interests include wireless communication and networks, simultaneous wireless information and power transfer (SWIPT) in cooperative relaying networks and MIMO systems.

**Kamran Ali Memon** received his Bachelor Degree in Electronics Engineering (2009) from Mehran University of Engineering Technology, Jamshoro Pakistan and Master's Degree in Communication (2015) from Quaid e Awam UEST Nawabshah Pakistan. He worked as a Lecturer/Assistant Professor in QUEST Pakistan for 08 years. Currently he is working toward his PhD at State Key Laboratory of Information Photonics and Optical Communications (IPOC), Beijing University of Posts and telecommunications, China. His research interests include optical and Wireless communications, PONs, Radio over fiber and WSNs.

# Optimization of a Three-Phase Tetrahedral High Voltage Transformer used in the Power Supply of Microwave

Mouhcine Lahame[1], Mohammed Chraygane[2], Hamid Outzguinrimt[3,] Rajaa Oumghar[4]

Materials, Systems and Information of Technology Team (MSTI)

ESTA, Ibn Zohr University

Agadir, Morocco

*Abstract*—**This article deals with the optimization of a three-phase tetrahedral-type high voltage transformer, sized to supply three voltage-doubling cells and three magnetrons per phase. The optimization method used is based on an algorithm implemented in Matlab/Simulink to study the influence of transformer geometrical parameters on the electrical operation of the power supply. This study will allow to find reduced volume of transformer respecting the current constraints imposed by the magnetrons manufacturer. The choice of optimal solution is done by calculation of magnetrons powers in order to respect the nominal operation.**

*Keywords—Optimization; tetrahedral; voltage-doubling; transformer; magnetrons*

## I. Introduction

In the development study of power supplies for microwave applications, we always seek to find more powerful and optimal solutions in terms of installation space, cost of manufacture and maintenance. In this context, this work defines a method to optimize the three-phase HV transformer used in this kind of power supplies with nine magnetrons (three for each phase).

Unlike the old power supplies already developed, whether single-phase magnetron 800Watts-2450Mhz or three-phase of three magnetrons 2400Whats-2450Mhz [1-3]. This new technology of nine magnetrons is optimized compared to that previously developed [4]. It offers an identical microwave power of 7200Watts-2450MHz. So it can use less optimized power supply to size a more powerful industrial microwave.

The design of this power supply is based on a three-phase HV transformer with magnetic shunts of tetrahedral type having a shell type structure. Each phase feeds three doublers cells giving a voltage which is suitable for the operation of the three magnetrons that delivers the microwaves [5]. The magnetic shunts of the transformer ensure the stability of the current in each magnetron in order to not exceed the values recommended by the manufacturer Imax<1.2A and Iavg $\leq$ 300mA.

This paper is divided in two sections. In the first, we present the model and the results obtained by simulation under Matlab/Simulink. In the second section, we study the influence of each geometrical parameter of the transformer on the operation of the power supply. This study will allow us to

define the optimized algorithm based on the simultaneous variation of these parameters. This leads to various solutions that respect the criteria recommended by the magnetron manufacturer. The choice of the best solution is validated by the calculation of the volume as well as the comparison of the results obtained with those of the non-optimized power supply, taking into account the operation of magnetrons in full power.

## II. Modeling and Simulation of the Three-Phase High Voltage Power Supply with Three-Magnetrons per Phase

### A. Description and Modeling of Power Supply

The general model of the three-phase HV power supply constitutes of a magnetic shunt transformer, doublers cells, and three magnetrons for each phase. The three-phase HV transformer is represented by three identical models of single-phase transformers coupled in star as shown in Fig. 1. Unlike conventional transformers, this special transformer contains intermediate magnetic shunts between the side columns and the central column, which allows to ensure the stabilization of the current in the magnetrons by the saturation of its magnetic circuit.



Fig 1. Three-Phase Power Supply of Three Magnetons per Phase having a Tetrahedral-Type Transformer.

From the different electrical and magnetic equations already developed of this transformer [4], each phase is modeled as a quadruple in π, composed of three non-linear inductances on the primary, secondary and shunts sides. These inductances have a section S, a length ℓ and a characteristic $\phi(i)$ which can be determined from the relation $L(i) = \frac{n_2\phi(i)}{I}$ and also the curve B(H) of the material transformer [6][7]. The different equations that determine the current and flux (I, φ) for each inductor are expressed as follows:

For primary inductance $L'_p$:

$$\begin{cases} \phi_p = n_2.S_1.B \\ i'_p = \frac{H.\ell_p}{n_2} \end{cases} \quad (1)$$

For secondary inductance $L_s$:

$$\begin{cases} \phi_s = n_2.S_2.B \\ i_s = \frac{H.\ell_s}{n_2} \end{cases} \quad (2)$$

For shunt inductance $L'_{sh}$ :

$$\begin{cases} \phi_{sh} = n_2.S_{sh}.B \\ i'_{sh} = \frac{H.\ell_{sh}}{n_2} \end{cases} \quad (3)$$

Fig. 2 shows the different geometrical parameters of a single tetrahedral transformer phase.

- The width of the core: a = 75 mm
- The width of the magnetic circuit : b = 25 mm
- Number of stacked sheets of the shunt : $n_3$=18
- Number of primary windings : $n_1$= 224
- Number of secondary windings: $n_2$ =2400
- Height of shunts: h = 0.5×$n_3$ mm
- Primary and secondary core surface: $S_1$= $S_2$= a×b
- Surface of shunt: $S_{sh}$ = b×h
- Thickness of the air gap: e = 0.75 mm
- $\ell_p$=4.5×a (correspond to the path ABCD)
- $\ell_s$=4.5×a (correspond to the path DAFE)
- $\ell_{sh}$=(2.5×a -2×e) (correspond to the path AD)

Each magnetron is presented by a model describing its operation which contains a diode with dynamic resistance R = ΔU/I = 350 Ohm and a threshold voltage E = 3800V. Fig. 3 shows the general model of the power supply.

### B. Simulation of the Model

The equivalent model of the three-phase HV power supply is implemented under Matlab/Simulink. The primary of the transformer is powered by a nominal voltage of 220/380V with a phase shift of 2π/3. Each non-linear inductance is modeled by Simulink blocks showing their operation. One of these blocks is used to interpolate the B(H) curve with the ANFIS neuro-fuzzy method [8-10]. Fig. 4 and 5 give the different

currents/voltages curves of magnetrons, diodes, capacitors and transformer secondary obtained in a previously study [4]. These curves will be the comparison tool between the optimized and non-optimized power supplies.



Fig 2. Geometry of a Single Tetrahedral Transformer Phase.



Fig 3. Electrical Diagram Equivalent of Three-phase High Voltage Power Supply for Three Magnetrons per Phase.



269 | P a g e

Fig 4. Result Obtained by Simulation of Magnetrons Currents, Diodes and Secondary.



Fig 5. Result Obtained by Simulation of Magnetrons Voltages, Capacitors and Secondary.

## III. OPTIMIZATION OF THE TETRAHEDRAL TRANSFORMER USED IN THE POWER SUPPLY

The optimization stage is based on the model developed on Matlab/Simulink of the three-phase HV power supply with three-magnetron per phase. This model will allow us to study with respect to the reference transformer case (non-optimized transformer) the sensitivity of each geometrical parameter to the nominal operation of the power supply [11-13]. This study will give us an idea of how we can simultaneously vary all the parameters in order to meet the following criteria:

- Have the various possible optimal solutions that offer a reduced volume without risk of exceeding the limits recommended by the magnetron manufacturer.

- Among the obtained solutions, find a better optimal one that respects the full power operation of the magnetrons.

### A. Influence of each Single Transformer Parameter on the Magnetron Current

The $\pi$ quadruple model of the three-phase tetrahedral HV transformer contains non-linear inductances that depend on the geometrical parameters. Therefore, the variation of such a parameter modifies the overall operation of the equivalent circuit of the power supply. The simulation results of the model permit to plot the variation of the maximum and average magnetron current in terms of the selected transformer parameters as shown in Fig. 6 to 10. These parameters must be within the ranges specified in Table I.

TABLE I.        VARIATION RANGE OF THE PARAMETRER

| Name of the parameters | Rating values |
|---|---|
| a (mm) | $45 \leq a \leq 75$ |
| $n_2$ | $2000 \leq n_2 \leq 2800$ |
| $n_3$ | $10 \leq n_3 \leq 18$ |
| e (mm) | $0,45 \leq e \leq 1,05$ |



Fig 6. Magnetron Current Simulation Results as a Function of «a».



Fig 7. Magnetron Current Simulation Results as a Function of «$n_2$».

Fig 8.    Magnetron Current Simulation Results as a Function of «$n_3$».



Fig 9.    Magnetron Current Simulation Results as a Function of «e».

From Fig. 6 to10, we can see that the maximum current and the average magnetron current decreases when the three parameters (a, $n_2$,e) decrease. On the other hand, for the case of the $n_3$ parameter, we notice that the increase in the number of plates constituting each shunt causes a decrease in the maximum value of the current in the magnetrons. We also note that the variation of $n_3$ does not exceed the acceptable limits of the magnetron current mean value.

These observations make us think of reducing the magnetic circuit volume of the transformer, without risk of damaging the magnetron tubes.

### B. Optimization Algorithm

The analysis of the results obtained previously confirms the magnetron current sensitivity with respect to the variation of each geometrical parameter of the transformer. These results prompted us to define an algorithm to study the influence of more than one geometrical parameter on the electrical operation of the power supply. We seek to simultaneously minimize all the parameters (a, $n_2$, $n_3$, e) thus the total volume of the tetrahedral transformer. The different optimal solutions found allow us to make a better choice between the transformer volume and the average power output by the magnetron.

The algorithm is used to minimize an objective function, which is the volume of the transformer in (cm$^3$) while respecting the operating conditions of the three phase magnetron power supply. The objective volume function is defined in the following form:

$$V_{transformer} = V_{core} + V_{copper} \tag{4}$$

From Fig. 2, $V_{core}$ can be calculated as follows:

$$V_{core} = 3 \times [( \text{ Total volume of the magnetic circuit)} \\ + \text{(Volume of the stack of shunts)} \\ - \text{(Total volume of the winding window)}]$$

$$V_{core} = 3 \times [(6a \times 5a \times b) + (2 \times (a \times b \times h)) - (2 \times \\ a \times b \times (3 \times a))] \tag{5}$$

The volume of the copper is defined by :

$$V_{copper} = 3 \times \left\{ S \left[ E_{coil/colum} \left[ \sum_{i=0}^{E_{coil/colum}} (b + a + 4d * i) \right] + \\ \left( N_{coil/colum} - E_{coil/colum} \right)(b + a + 4d(i + 1)) \right] \right\} \tag{6}$$

- d : presents the diameter of the copper cable on the primary or secondary side

- $N_{coll/colum}$ : presents the  number of turns per column

- $E(N_{coll/colum})$ :presents the entire part of  $N_{coll/colum}$

- $N_{coll/colum} - E(N_{coll/colum})$ :  presents the fractional part of $N_{coll/colum}$

By using the different intervals of variation of the geometrical parameters defined previously. The vector X used in our algorithm (Fig. 10) will take all the possible combinations between the different parameters X = [$x_1$, $x_2$, $x_3$, $x_4$] = [a, $n_2$, $n_3$, e]. Table II gives the step and the variation margin of each combined and defined parameter in the vector X.

TABLE II.    STEP AND VARIATION MARGIN OF EACH PARAMETER

| Name of the parameters | Start | Step | End |
|---|---|---|---|
| a(mm) | 45 | 5 | 75 |
| $n_2$ | 2000 | 100 | 2800 |
| $n_3$ | 10 | 2 | 18 |
| e(mm) | 0,45 | 0,15 | 1,05 |



Fig 10.   Different Steps of Optimization Algorithm.

TABLE III.     BEST SOLUTION OBTAINED BY THE OPTIMIZATION ALGORITHM

|  | a (mm) | $n_2$ | $n_3$ | e (mm) | Imax (A) | Iavg (mA) | Volume (cm³) | Pavg (w) |
|------|------|------|------|------|------|------|------|------|
| ref | 75 | 2400 | 18 | 0,75 | 0,95 | 271,17 | 9988,36 | 1100,47 |
| $S_1$ | 75 | 2400 | 18 | 0,45 | 0,95 | 271,17 | 9989,01 | 1100,47 |
| $S_2$ | 65 | 2600 | 10 | 0,75 | 0,91 | 264,34 | 7788,87 | 1074,60 |
| $S_3$ | 70 | 2400 | 10 | 0,75 | 0,90 | 259,73 | 8834,97 | 1054,61 |
| $S_4$ | 70 | 2400 | 16 | 1,05 | 0,84 | 256,53 | 8859,06 | 1036,58 |
| $S_5$ | 55 | 2800 | 14 | 0,45 | 0,79 | 244,85 | 5925,99 | 983,47 |

At each X iteration, we perform a model simulation on Simulink using the "sim" function in Matlab. The results obtained from each simulation will be checked in order to take the one that respects the operating constraints of the power supply. By displaying to each solution found the values of vector X, the max and average current as well as the average power of the magnetron.

After simulating the model with the different combinations of the geometrical parameters, Table III presents the five best solutions selected that meet the current imposed by the manufacturer also that it allows to operate the magnetrons in nominal power.

From Table III, we find that solution $S_2$ presents the best compromise between the volume of the transformer and the operation of the magnetron (Pavg = 1074,6W, Imax= 0.91(A) and Iavg =264,34 mA, Volume= 7788,87cm²). For the solution $S_5$, it has a minimum volume, but it does not allow nominal operation at full power of the magnetrons.

We simulate the model under Matlab/Simulink with the new geometrical parameters of the transformation optimization solution. The waveforms of the voltages and currents obtained (Fig. 11 and 12) are almost identical to those obtained in the case of reference, while respecting the operating constraints of the magnetrons. So we can say that 22 % of the power supply volume is optimized without having a large magnetron power lost after optimization.

Transformer optimization rate ($\tau$).

$$\tau = \frac{Vref - Vopt}{Vref} = \frac{9988,36 - 7788,87}{9988,36} = 0,22 \qquad (7)$$

Magnetron power lost after optimization ($P_{lost}$).

$$P_{lost} = P_{mref} - P_{mopt} = 1100,47 - 1074,60 = 25,87W \qquad (8)$$

$P_{mref}$ : Average magnetron power given by the reference transformer.

$P_{mopt}$ : Average magnetron power given by the optimized transformer.



Fig 11.   Results Obtained by Simulation of the Model Optimized. Currents of Magnetrons, Diodes and Secondary.

Fig 12. Results obtained by Simulation of the Model Optimized. Voltages of Magnetrons, Capacitors and Secondary.

## IV. CONCLUSION

In this work, we have succeeded after a study of optimization, to define a proper algorithm that aims to find an optimal solution of the three-phase tetrahedral type HV transformer used to feed three magnetrons per phase. This study allowed us to reduce the volume, the congestion as well as the cost. The optimized solution obtained is compared to that of reference; it gave a transformer that meets the operating criteria of the entire power supply.

As a perspective, this work can be used as a reference to optimize another type of transformer employed for microwave application at N = 1,2,3 magnetrons per phase while seeking to do a thermal study on the optimized transformer.

REFERENCES

[1] Elghazal, M. Ould ahmedou, M. Chraygane, M. Ferfra, A. Belhaiba. Optimization of high voltage power supply for industrial microwave generators for one magnetron. Journal of Theoretical and Applied Information Technology, 2012, vol. 46, no 1, pp. 001-010.

[2] H.Outzguinrimt, A.Bouzit, M.Chrayagne, M.Lahame, R.Oumghar, M.Ferfra. Design and Modeling of New Configuration of Three Phase Transformer For High Voltage Operation Using in Microwave Industrial. In : 2018 International Conference on Electronics, Control, Optimization and Computer Science (ICECOCS). IEEE, p.1-6.

[3] Naama El Ghazal, M. Chraygane, M. Ferfra, A. Belhaiba, M. Fadel, B. Bahani," Modeling of New Three-phase High Voltage Power Supply for Industrial Microwave Generators with One Magnetron Per Phase" International Journal of Electrical and Computer Engineering (IJECE), Vol. 3, No. 2 , pp. 270-278, April 2013.

[4] M.Lahame, M.Chraygene, H.Outzguinrimt, R.Batit, R.Oumghar, M.ferfra. Modeling Under MATLAB by ANFIS of Three-Phase Tetrahedral Transformer Using in Microwave Generator for Three Magnetrons Per Phase. Telkomnika, vol. 16, no.5, pp. 2406-2414, 2018.

[5] B.Bahani, A. Bouzit, M.Chraygane, M.Ferfra, N. El Ghazal, A. Belhaiba .Modeling of a New High Voltage Power Supply for Microwave Generators with Three Magnetrons. International Journal of Electrical & Computer Engineering (IJECE). vol. 3, no.2, pp.2088-8708, 2013.

[6] Guanghao Liu, Xiao-Bang Xu. "Improved Modeling of the Nonlinear B–H Curve and Its Application in Power Cable Analysis". IEEE Transaction on Magnetics, ; 38(4), 2002.

[7] Z.Zhao, F.Liu, Z.Cheng, W.Yan, L.Liu ,J.Zhang, Y.Fan. Measurements and Calculation of CoreBased B-H Curve and Magnetizing Current in DC-Biased Transformers. IEEE Transactions on Applied Superconductivity. vol .20 ,no.3, pp. 1131-1134, 2010.

[8] M.Bassoui, M.Ferfra, M.Chraygane. Improved modeling of a single-phase high voltage power supply for microwave generators for one magnetron. J Electr Eng, vol.16, no.4, 2016.

[9] M.Buragohain, C.Mahanta . A novel approach for ANFIS modelling based on full factorial design. Applied soft computing, vol.8, no.1, pp. 609-625, 2008.

[10] O.Kisi, Suspended sediment estimation using neuro-fuzzy and neural network approaches. Hydrological Sciences Journal, vol.50, no.4, pp. 683-696, 2005.

[11] M.Ould.Ahmedou, M.Ferfra, N.EL Ghazal, M.Chraygane,M. Maaroufi "Implementation and Optimization Under Matlab Code of a HV Power Transformer For Micowave Generators Supplying Tow Magnetrons", Journal of Theorical and Applied Information Technology, vol. 33 ,no.2, 30 November 2011.

[12] A. Belhaiba, N. Elghazal, M. Chraygane, M. Ferfra, B. Bahani, M.ould Ahmedou "Improved optimization of the Nominal Functioning of a High Voltage power supply N=2 mangnetrons for microwaves generators". International Journal of Electrical and Computer Engineering(IJECE),vol.2, no.5,pp.708-716,October 2012.

[13] H.Satoru, H.Tomoyuki, M.Mitsunori. Hybrid optimization using DIRECT, GA, and SQP for global exploration. IEEE Congress on Evolutionary Computation. pp.1709-1716 , 2007.

# LBPH-based Enhanced Real-Time Face Recognition

Farah Deeba[1], Aftab Ahmed[4]
School of Information and Software Engineering
University of Electronic Science and Technology of China
Chengdu, Sichuan, China

Fayaz Ali Dharejo[3]
Computer Network Information Center
Chinese Academy of Sciences, University of Chinese
Academy of Sciences, Beijing, China

Hira Memon[2]
Department of Computer System Engineering
Quaid e Awam University of Engineering Science and
Technology
Nawabshah, Pakistan

Abddul Ghaffar[5]
Key Laboratory of Instrumentation Science & Dynamic
Measurement, Ministry of Education
North University of China
Taiyuan, China

*Abstract*—**Facial recognition has always gone through a consistent research area due to its non-modelling nature and its diverse applications. As a result, day-to-day activities are increasingly being carried out electronically rather than in pencil and paper. Today, computer vision is a comprehensive field that deals with a high level of programming by feeding the input images/videos to automatically perform tasks such as detection, recognition and classification. Even with deep learning techniques, they are better than the normal human visual system. In this article, we developed a facial recognition system based on the Local Binary Pattern Histogram (LBPH) method to treat the real-time recognition of the human face in the low and high-level images. We aspire to maximize the variation that is relevant to facial expression and open edges so to sort of encode edges in a very cheap way. These highly successful features are called the Local Binary Pattern Histogram (LBPH).**

*Keywords—Face recognition; feature extraction; Local Binary Pattern Histogram (LBPH)*

## I. INTRODUCTION

Among other biometric methods, face recognition is also one of the ways to identify any individual subject. Face recognition identifies anyone by comparing the physical characteristics of the item. There are two face recognition modes, still images and live video. The first step in face recognition is face detection. Therefore, to perform facial recognition, the system must position the face earlier in the input image or video stream. This step is called face acquisition or detection. In this section, the main face recognition methods are described. Such as neural networks, geometric feature matching, feature-based methods and graph matching [1]. The following techniques were studied in the context of facial illustrations. The geometric feature matching method depends on the calculation of a set of geometric features of the facial image. The general specification is defined by a vector showing the position and size of the main facial features, for example, the eyebrows, nose, mouth and facial contours. T. Kanade [2] completed an innovative study on automatic face recognition by implementing conventional features. His system achieved higher results with 75% accuracy on a data set of 20 subjects, and each subject took two photos, one for the model and one for the test. Later, I.J.

Cox et al. [3] proposed a hybrid distance method for automatic face recognition. He scored 95% peak recognition accuracy on a dataset of 685 subjects. In the hybrid distance method, each face image shows 30 manual induced distances. B.S. Manjunath et al. [4] proposed the fragmentation process for the detection of feature units of every single facial image, which decreased the storing capacity of the database, and it generated 35-45 feature units per facial image. Face detection based on geometrical approaches based on the geometrical proportion between regular features by employing some statistical models. These models calculate the distances between features; they could be more beneficial for identifying expected matches in a large dataset. Geometric Feature-based algorithms have some advantages over other methods like rotation independently, faster with execution time, scaling [5].

Like geometrical feature matching and feature-based face recognition approach, graph matching face recognition approach is one of them [6]. Introduced dynamic link architecture to falsification invariable object identification, in this approach researcher utilize elastic graph matching approach for calculating the nearest saved graph. This approach called Dynamic link structure is an addition to standard ANN. Sparse graphs represented memorize objects. The vertices of sparse graphs are tagged through a multi-resolution statement in the context of a local power spectrum, and their borders are tagged with geometric distance vectors.

Recognizing the object from multimedia (Video, image) is called object recognition, can be identifying by employing any efficient method like elastic graph matching. Elastic graph matching process by matching the cost function randomly modified at every node. The better testing outcomes were obtained on the dataset of 87 subjects and a group of office objects containing various expressions with an alternation of 15 degrees. The matching procedure takes more computation time; it takes 25 secs to match with 87 saved objects on the symmetric device by 23 carriers. After that L. Wiskott [6] modified this approach and compared individual front view faces of 112 images. Probe pictures were deformed because of the rotation in depth and variation in facial appearance. Functional outcomes were achieved of facial images on big

rotated angles. Resulting obtained 86.5% recognition percentage on testing of 111 face images at 15 degrees' rotation, and 66.4% recognition percentage on testing of 110 face images at 30 degrees' rotation of a 112 neutral frontal views [7]. Generally, dynamic link structure is dominant on various facial recognition approaches in the context of rotation stability; though, the matching procedure is extensive in terms of computation.

The neural network is so much simplified face recognition approach because of its non-linear architecture in the net system. Therefore, the features extraction phase is more effective than the linear technique, it selects a dimensionality reducing linear projection that increases the scatter of all expected models [8]. ORL database contains 40 objects, with 400 images of each object, recognition accuracy was 96.2% obtained on this dataset. It takes 4 hours of training time and less than 0.5 seconds for classification and provides limited invariance to transformation, variation, scale, and distortion. Though, the number of the individual is proportional to computing time. Increase in the number of people, the computation time also increases. Generally, neural networks approach gets difficulties when the number of people increases. Furthermore, a neural networks approach is not appropriate for only a single model image recognition experiment, since various model pictures per individual subject are required to train the system at "optimum" parameters setting.

SIFT is the most well-known and widely adopted technique for feature extraction. This technique jointly uses the difference of Gaussians (DOG) and image pyramid concepts. By employing this technique, image is processed at different scales by the Gaussian filter. This technique gives excellent achievement if there are any illumination or viewpoint variations presents, it is also invariant to the rotation as well as a scaling factor [10]. Under the SIFT feature matching technique, each feature of the test images is compared to the dataset images. Euclidean distance gives best-extracted feature vector.

SIFT algorithm has four primary phases for the feature matching process named as Scale-Space Extreme Detection, KeyPoint Localization, Orientation Assignment and Key Point Descriptor" [10]. SIFT algorithm is straightforward and gives better results, but one drawback is of its computational complexity and time.

SURF is a features indicator in an input image. SURF detects the local features of the face image. The SURF features detector is more efficient and robust than the SIFT features detector. Comparing to the SIFT, SURF provides excellent results. SURF indicators determine the interest points in the face image. For the detection of interest points, SURF uses decimal number estimation of the determinant of Hessian blob indicator that could be calculated by three decimal numbers actions by utilizing a pre-calculated integral picture [9].

CNN is a very excellent and mathematical tool used for complex computation for multimedia object images videos, for many tasks like recognition, segmentation, etc. CNN is also handy for the study of 2D variability shapes. In CNN

Local features with some share, weights are combined [11] which is also used for sub-sampling purposes such as for shifting the level, scale-invariance, and deformation.

Gabor wavelet frequency and alignment demonstrations based on human graphical method, it is more suitable for representing texture learning. Extracting the feature from particular positions are also supported by Gabor wavelet, also perform image analysis at different scale and orientations. Frequency and rotations variations are also handled in these techniques [12]. By using the Gaussian envelope, the Gabor wavelet is an accentuated [13].

The principal component analysis is a traditional algorithm, broadly employed in machine vision and pattern recognition technology as well as used in feature extraction purpose [14]. In this algorithm, it is stated that "any face image can be reconstructed nearly as a weighted sum of a small set of images which define a facial base (Eigen images), and an average image of the face." In 1991, Eigenfaces technique was suggested by Turk and Pentland [14]. Eigenfaces were also proposed for facial images recognition. Meanwhile, PCA became best effective algorithm for facial recognition. The PCA was frequently used for features extraction and dimension reduction.

So far, the literature is well defined, the rest of the work contains an overview of the facial recognition system, Section III deals with the methodology and Section IV shows the result and the discussion, and finally conclusion and the future meaning is given.

## II. OVERVIEW OF THE FACE RECOGNITION SYSTEM

Face detection and identification process is a machine learning technique, by learning and extracting the physical characteristics of the human. Matching these features with the tested images can identify the person or deny those people to recognize. There are several challenging and varying parameters in face detection and identification like illumination, different poses, change expressions, low-quality input images, etc.

There are several different perspectives about face detection and recognition system; some of the projects only focus on images with high resolution; some of them focus on low resolutions. Recently researchers focus on the different frontal view of images, from different angles, different lighting illuminations, etc.

Traditionally, Face recognition system follows four primary phases, listed follows; also the basic face recognition diagram is shown in Fig. 1.

*a)* Face Detection

*b)* Preprocessing

*c)* Feature Extraction

*d)* Feature Matching

### A. Dataset Preparation

We create our own dataset; the dataset contains a total 1000 images, 333 face images of each person with $60 \times 60$ resolution of each image. It is created based on face detection.

Make different facial expressions and postures to a scene and detect faces. The saved pictures are stored in the same folder to form the generated face dataset. At this stage, the dataset is preprocessed for the feature extraction process. The dataset images have been converted into grayscale images for features extraction, and then normalized those images for good recognition results. For features detection, Haar modules have been used to detect these local features in a given an input image. Here, the input image refers to the digital image captured by the camera. After detecting features, the classifier will classify the input image as a face image as shown in Fig. 2.

In this project, face detection algorithms are developed based on Local Binary Patterns Histogram (LBPH). The LBPH-based algorithm, the first step is to extract the image pattern with the LBPH algorithm. Then, two thresholds are set to calculate the probability of face in the image pattern. After that, the sliding window applied to identify the faces in given images and recognize those faces. From Fig. 3 we can understand well.



Fig. 1. Basic Phases in Face Recognition.



Fig. 2. Image Preprocessing and Feature Selection.



Fig. 3. Face Detection and Recognition Methodology.

## III. METHODOLOGY

### A. Local Binary Patterns Histogram

The goal of face detection is to detect [15] and locate faces in the image, to extract human face to use in other areas. Nowadays, there are many different algorithms to accomplish face detection or recognition, such as Fisher faces, Eigenfaces, Scale-invariant Feature Transform (SIFT) and Speed-Up Robust Features (SURF). In this section, LBPH-based face detection algorithm is introduced. LBPH algorithm is the combination of Local Binary Patterns (LBP) and Histograms of Oriented Gradients (HOG) descriptor. LBP is an easy but powerful way to extract and label the pixels of an image. Using the LBPH, we can easily represent face images with just a straightforward vector.

Ojala et al. first introduced local Binary Patterns (LBP), and it is designed to be a texture analysis for the gray-scale image [16] [17]. To detect faces in an RBG (colored) photo, we have to convert the image into a grayscale image at first. For each pixel $P_{i,j}$, it is a vector that contains three values, which is to represent the degree of red, blue and green. We convert the RBG image into graa y-scale image by:

$$G_{i,j} = (0.2989, 0.5870, 0.1140)^T . P_{i,j} \qquad (1)$$

Where $G_{i,j}$ represented the corresponding pixel in grathe the y-scale image. The LBP operator is going to compare the center value with its P neighbor values on the circle with radius R and assigns neighbor value as 1 if center value is bigger than the neighbor, assigns 0 on the contrary. In this case, we set P = 8 and R = 1 which means that we consider a 3×3 check. The LBP operator labels the center pixel by thresholding the 8 neighbors. For each center pixel, the LBP operator outputs an 8 bits' binary number and we convert it into a decimal between 0 and 255 as a result.

We utilize the following notation to describe the LBP operator:

$$LBP_{P,R}(x,y) = \sum_{i=1}^{8} \sin(G_i - G_c)2^i \qquad (2)$$

Where $LBP_{P,R}(x,y)$ is the results of the LBP operator, (x, y) is the coordinate of center pixel in the 3_3 check. $G_c$ is the center value in gray-scalthe e image, $G_i$ is the neighbor value in gray-scale the image?

The LBP operator reduces the influence of illumination and returns the texture by considering every pixel in an image that excludes the boundary pixels. The figure shows the demo of the LBP operator for one pixel.

### B. Extracting Histograms with LBP

Using the LBP result, we can generate a histogram for this image and formed a data vector to describe the patterns of the original image. The Histogram is about the frequency of the occurrences of LBP result for each pixel. From the last part, after doing the LBP operator, the value of each pixel is between 0 and 255. Thus, the histogram contains only 256 positions.

However, we are not considering the whole image directly. We divide the image into several image pieces. We calculate

the values of 256 positions for each image piece given Nx and Ny, representing the number of cells in the horizontal direction and vertical direction respectively.

Then, we concatenate all histograms of all image pieces to form a more significant histogram. For instance, if we set Nx = 4 and Ny = 4, we will have 4 ×4 × 256 = 4.096 positions in the total histogram. Finally, we can use this total histogram to represent the image by just one data vector.

### C. LBPH Algorithm

The Local Binary Patterns Histograms (LBPH), was introduced in the year 2006 [18]. LBPH algorithm was commonly used for facial recognition. This algorithm is based on the local-binary-operator [19], broadly implemented in face recognition, due to its discriminating strength and calculation easiness [20].

Face recognition is performed by employing the Local Binary Pattern Algorithm. The LP operator is applied for local binary features by considering the Local Binary patterns [21] which helps to shorten the local special features of the face image. The LBP is the binary ratio of pixels intensities within the center pixel. And it's around eight pixels. The mathematical description is described in the below equation.

$$LBP(p_c - q_c) = \sum_{m=0}^{7} S(t_m - t_c)^{2^t} \qquad (3)$$

Center pixel is shown by $t_c$ and (pc, qc) represents the surrounded eight pixels, it is very though very useful to determine the face feature. In face matrix feature extracted from the image to compare the values with center pixel values to finally generate binary code.

### D. Feature Vectors

Images divided into the region to show the faces efficiently then these images have been subdivided into $A^2$ regions, i-e 82= 64 regions. Histogram of each image is composed by each potential label, where each bin in histogram tells the information about the pattern to get the feature vectors from histograms. The each regional histogram V (V−1) + 3 bins: V ($V_1$).

To perform the specific area with the help of the LBP operator from the edges of the image if not exist this means some section of the border is not related. For the image (CxD), the feature vector is designed with the help of calculating the LBP code for all pixels $(P_c, Q_c)$ with xcϵ {U + 1,. .., C − U} and $q_c$ {U + 1, . . . , D − U}. If an image is divided into a × a regions, then the histogram for region$(a_p, a_q)$, with $a_p$ϵ {1, . . ., a} and $a_q$ϵ {1, . . . , a}, Mathematically,

$$J_k(A_p - A_q) = \sum p, q\{LBP_{v,u(p,q)=M(K),k=1,\ldots\ldots V(V-1)+3}\} \quad (4)$$

$$P \in \left\{U + 1, \ldots \ldots \ldots \frac{C}{A}\right\} A_p = 1$$

$$\{(A_p - 1)(\frac{C}{A})\} 1, \ldots \ldots \ldots C - U \quad A_p = A$$

$$\{(A_p - 1)(\frac{C}{A})\} 1, \ldots \ldots \ldots A_p((\frac{C}{A})\text{else}$$

$$\in \left\{U + 1, \ldots \ldots \ldots \frac{D}{A}\right\} A_p = 1$$

$$\{(A_q - 1)(\frac{D}{A})\} 1, \ldots \ldots \ldots D - U \quad A_q = A$$

$$\{(A_q - 1)(\frac{D}{A})\} 1, \ldots \ldots \ldots A_q((\frac{D}{A})\text{else}$$

In which M is the label of binary k and B(Z) = {1, Z is True 0, Z is False

From the feature vector, we can get three districts levels of the locality of the face; these labels combine information of the little background level and architecture histogram which provide the knowledge about the face.

### E. Comparing the Feature Vectors

To measure the feature of images, the sample (H) and a model (I) are used as so that the difference values between feature vectors can be measured. Here with the help of histograms can measure the difference between two images. - Histogram Intersection

$$F(H, I) = \sum_{v=1}^{j^2} \left(\sum_{e=1}^{Q(Q-1)+3} \min(H_{e,v}, I_{e,v})\right) \qquad (5)$$

Log-likelihood statistics

$$M(H, I) = \sum_{v=1}^{j^2} \left(-\sum_{e=1}^{Q(Q-1)+3} H_{e,v} \log I_{e,v}\right) \qquad (6)$$

Chi-square statistics

$$Y^2(H, I) = \sum_{v=1}^{j^2} \left(\sum_{e=1}^{Q(Q-1)+3} \frac{(H_{e,v} - I_{e,v})^2}{H_{e,v} + I_{e,v}}\right) \qquad (7)$$

The $x^2$ wthe eight of face rank the similarity of images by the computing the histograms. The deeper the value of the x2, the greater is the similarity.

## IV. RESULTS AND DISCUSSIONS

The proposed primary face recognition application has been implemented in Python, Open CV image processing library and LBPH algorithm over HD camera. In this application, the algorithm applied for face recognition is distributed in three dissimilar and independent sections.

1) Face image acquisition Module, 2) Dataset training module, 3) Face recognition module. In the image acquisition module, the user needs to run "Dataset Creator. Py" file from Python IDLE shell and enter the subject ID (see Fig. 4). This module will open the "Image Acquisition" window, which will detect and capture the face images.

Application opens the externally connected camera with PC, apply the Haar classifier to detect the face and capture face images (see Fig. 5). The algorithm permits to store face images in a folder with the subject ID and sample number.

After the preprocessing, all the 2000 captured images of 4 subjects will be stored in the same folder. Each image will be assigned subject ID and sample number. Sample number is the number of images per face image. So the sample number will be different, while the subject ID will be the same for a single face image (see Fig. 6).

Fig. 4.    Image Acquisition.



Fig. 5.    Face Detection.



Fig. 6.    Subject ID with Sample Number.

### A.  Experiments

This part states that the experiments executed to match the performance of the facial recognition algorithm applied in this system. It is required to describe the experimental conditions that must be acceptable for the proposed system. So, the proposed experiments must be executed by changing certain factors that are significant in the learning and recognition method of the system. The following test was implemented to examine the efficiency of the operating system—identification percentage based on a threshold. To perform the experiments, besides my dataset images, I also have taken dataset images of

my three lab mates from the school of Information & Software Engineering. The tests were executed, and the efficiency ratios were examined in every situation. Fig. 7 displays the subjects with their corresponding IDs. After the execution of the training of the experiment subjects, facial recognition is executed. The expected results generated by the system are defined below.

### B.  True Positive

The real positive condition happens when the observed individual's data stored in the dataset folder and the recognized subject matches to the one that is available in the training dataset. Fig. 8 displays the result after applying facial recognition; in this situation, the result is accurate.

### C.  True Negative

The true adverse condition occurs when the tested subject's data is not stored in the database, and the system could not recognize that subject. Fig. 9 displays the result after applying face recognition; in this condition, the recognition result will be "Unknown," and hence the effect will be considered as correct.

### D.  True Occlusion

True occlusion condition happens when the subject's data saved in the dataset folder without occlusion condition and the recognized subject matches to the one that is available in the training dataset. Fig. 10 displays the result after applying facial recognition; in this situation, the result is accurate.



Fig. 7.    Test Subjects.



Fig. 8.    True Positive.



Fig. 9.    True Negative.

Fig. 10. True Occlusions.

### E. True Pose Variation

True pose variation occurs when the subject's position is at a different angle concerning the camera, and the system can recognize that subject even if in the dataset, only frontal view face images are stored as shown in Fig. 11.

We evaluate the parameters and different values of our model during the testing experiment which is given in the Table I.

The accuracy of model is shown in Fig. 12 the accuracy is taken against the Training and Validation dataset. We divided our dataset, such as 50% training dataset, 30% testing and rest 20% validation dataset.

### F. Setup

To implement the above jobs subsequent Hardware and Software are required to grow the proposed the scheme. We used software tools: Windows/UNIX operating system, Python 2.7/3.6, OpenCV library, Numpy library, Matplotlib library, Pillow library.



Fig. 11. True Occlusions.

TABLE I.        PARAMETERS TO TEST THE EXPERIMENT

| Parameter | Value |
|---|---|
| **Training** | |
| Number of Subjects | 4 |
| Cam-Subject Distance(cm) | 200 |
| **Recognition** | |
| Threshold | <250;<500;<750;<100 |
| times | 80[5 each subject] |



Fig. 12. Training v/s Testing Accuracy Curve.

## V. CONCLUSION AND FUTURE WORK

We purposed LBPH for image recognition and face detection in the surveillance camera in a specific area. Having obtained good results from various experimental analyzes of this technique, they also provide valid results for occlusion, pose variation, and illumination. Therefore, the proposed system allows recognition and recognition of faces in a controlled environment. As machine learning is very important nowadays, there are many areas where this work can be expanded. In implementing this project, we have identified some areas for improvement, such as Limitations of distance, the maturity of algorithms and camera qualities, even using DNN techniques. Accuracy can be improved in the future, which is more directly related to our work.

### REFERENCES

[1] Elham Bagherian, Rahmita Wirza O.K. Rahmat, "Facial feature extraction for face recognition: a review", IEEE 2008 International Symposium on Information Technology.

[2] T. Kanade, "Picture processing by computer complex and recognition of human faces, "technical report, Dept. Information Science, Kyoto Univ., 1973

[3] I.J. Cox, J. Ghosn, and P.N. Yianios, "Feature Based face recognition using mixture distance," Computer Vision and Pattern Recognition, 1996

[4] B.S. Manjunath, R. Chellappa, and C. von der Malsburg, "A Feature based approach to face recognition," Proc. IEEE CS Conf. Computer Vision and Pattern Recognition, pp. 373-378,1992

[5] P. Suja and S. Tripathi, "Analysis of emotion recognition from facial expressions using spatial and transform domain methods," International Journal of Advanced Intelligence Paradigms, vol. 7, pp. 57–73, 2015.

[6] M. Lades, J.C. Vorbruggen, J. Buhmann, J.Lange, C. Von Der Malsburg, R.P. Wurtz, and M. Konen, "Distortion Invariant object recognition in the dynamic link architecture," IEEE Trans. Computers, vol. 42, pp. 300-311,1993.

[7] Dharejo FA, Jatoi MA, Hao Z, Tunio MA. PCA based improved face recognition system. Frontiers in Artificial Intelligence and Applications. 2017. https://doi.org/10.3233/978-1-61499-785-6-429

[8] KIRBY, M. AND SIROVICH, L.. Application of the Karhunen-Loeve procedure for the characterization of human faces. IEEE Trans. Patt. Anal. Mach. Intell. 12 ,1990.

[9] E. Paul and A S Ajeena Beegom, "Mining images for image annotation using SURF detection technique," IEEE International Conference on Control Communication & Computing India, Trivandrum, 2015, pp.724-728.

[10] V. Purandare and K. T. Talele, "Efficient heterogeneous face recognition using Scale Invariant Feature Transform," IEEE International Conference on Circuits, Systems, Communication and Information Technology Applications, Mumbai, 2014, pp. 305-310

[11] C. Yan, C. Lang, T. Wang, X. Du and C. Zhang, 2014, "Age estimation based on convolutional Master Thesis of University of Electronic Science and Technology of China 42 neural network" Pacific Rim Conference on Multimedia, springer, Malaysia, 2014, pp. 211-220

[12] J. K. Kamarainen, "Gabor features in image analysis," IEEE International Conference on Image Processing Theory, Tools and Applications, Istanbul, 2012, pp. 13-14.

[13] S. Murala, A. B. Gonde and R. P. Maheshwari, "Color and Texture Features for Image Indexing and Retrieval," IEEE International Advance Computing Conference, Patiala, 2009, pp. 1411-1416

[14] M. Kirby and L. Sirovich, "Application of the KL Procedure for the Characterization of Human Faces," IEEE Trans. Pattern Analysis and Machine Intelligence, vol. 12,1no. 1, pp. 103-108,Jan. 1990

[15] Dr C.Sunil Kumar1 , C.N Ravi2 and J.Dinesh3. Human Face Recognition and Detection System with Genetic and Ant Colony Optimization Algorithm IOSR Journal of Computer Engineering (IOSR-JCE) e-ISSN: 2278-0661,p-ISSN: 2278-8727, Volume 16, Issue 4, Ver. VII, Jul – Aug. 2014

[16] A. Ahmed, J. Guo, F. Ali, F. Deeba and A. Ahmed, "LBPH based improved face recognition at low resolution," 2018 International Conference on Artificial Intelligence and Big Data (ICAIBD), Chengdu, 2018, pp. 144-147.doi: 10.1109/ICAIBD.2018.8396183

[17] C. Darwin. The Expression of the Emotions in Man and Animals. London:John Murray, 1872.

[18] H. A. Ahonen, T. and M. Pietikinen, "Face description with local binary patterns," IEEE Trans. Pattern Analysis and Machine Intelligence, vol. 28, no. 12, pp. 2037–2041, 2006.

[19] P. M. Ojala, T. and D. Harwood, "A comparative study of texture measures with classification based on feature distributions," Pattern Recognition, vol. 19, no. 3, pp. 51–59, 1996

[20] M. Pietikinen, "Local Binary Patterns," vol. 5, no. 3, p. 9775, 2010

[21] T. Chen, Y. Wotao, S. Z. Xiang, D. Comaniciu, and T. S. Huang, "Total variation models for variable lighting face recognition" IEEE Transactions on Pattern Analysis and Machine Intelligence, 28(9):1519{1524,2006

# Visualizing Code Bad Smells

Maen Hammad[1], Sabah Alsofriya[2]

Department of Software Engineering, The Hashemite University

Zarqa, Jordan

*Abstract*—Software visualization is an effective way to support human comprehension to large software systems. In software maintenance, most of the time is spent on understanding code in order to change it. This paper presents a visualization approach to help maintainers to locate and understand code bad smells. Software maintainers need to locate and understand these bad smells in order to remove them via code refactoring. Object oriented code elements are visualized as well as their bad smells if they exist. The proposed visualization shows classes as building and bad smell as letter avatars based on the initials of the names of bad smells. These avatars are shown as warning signs on the buildings. A framework is proposed to automatically analyze code to identify bad smells and to generate the proposed visualizations. The evaluation of the proposed visualizations showed they reduce the comprehension time needed to understand bad smells.

*Keywords*—*Software visualization; program comprehension; data modeling; bad smells*

## I. INTRODUCTION

Code bad smells are symptoms of poor design and implementation choices [1]. These bad smells have negative impact on the maintainability of the code. Badly written code is hard to understand, test and change. As a result, the changes of bugs increase. So, maintainers have to locate these smells in the code in order to remove them. Code smells are removed by a process called refactoring [1]. It is the process of rewriting the code to improve its internal structure without changing its external behavior.

The problem is how to identify these bad smells and locating code elements affected by these smells. Most of bad smells detecting tools reports results as formatted text. Developers have to go back to the source code and check the identified smell. They need to understand the cause of the smell in order to remove it. Understanding the smell with its causes in the code is essential to the refactoring process. The research question that we are trying to address is; how to represent or model code smells within its static code environment?

Program comprehension is essential to software maintenance activities. Most maintenance cost is spent on understanding the current status of the code and the system in general. Maintainers consume time and effort when interacting with large scale projects in order to understand them. Visualization is an effective way to ease the interaction process with code and hence support comprehension tasks. Our premise is that visualizing code smells with structural code elements supports maintenance activities by reducing comprehension time.

In this paper, we propose a visualization technique to model bad smells as well as their locations in the structural code environment. Classes are modeled as buildings. Each building consists of a number of floors that match the number of methods in the class. The number of class attributes, LOC for each method and its parameters are also visualized in the buildings. Bad smells are visualized as signs of letter avatars on the buildings. Each bad smell is modeled by a different avatar based on the initials of the smell's name. A framework is proposed to automatically analyze object oriented source code and to generate the proposed visualizations.

This paper presents our early work towards realizing the proposed framework as a complete visualization tool. We also illustrate the proposed visualizations and we show how they can be useful for program comprehension tasks. The main two research contributions of this paper are:

- Easy to understand visualizations to locate and identify bad smells in code.

- A framework to automatically analyze source code to generate the proposed visualizations.

This paper is organized as follows. Section 2 summarizes the main related research in the area. Section 3 presents the proposed visualizations. The proposed framework is detailed in Section 4. The evaluation of the proposed visualizations is presented in Section 5 followed by our conclusions and future work.

## II. RELATED WORK

There are many researches in the software visualization area that hard to cover in this paper. We focus on the most related work to ours that can be categorized into visualizing static code structure and visualizing bad smells.

### A. Visualizing Static Code Structure

Ducasse and Lanza [2][3] presented a novel visualization for classes named class blueprint. It visualizes the internal structure of classes to support class understanding. A well-known 3D visualization approach, that model software as cities, is presented in [4][5][6]. The visualization maps the information about the source code in meaningful ways related to real cities. Another visualization approach for architecture and metrics of software systems as 3D software cities is presented in [7]. Panas et al. [8] proposes a 3D visualization metaphor to model software production cost as real cities. A reverse engineering environment called Rigi is presented in [9] to analyze, interactively explore, summarize, and document large projects. Marcus et al. [10] presented the sv3D framework for software visualization. The visualization is

focused on source code and testing levels. Langelier et al. [11] proposed a visualization framework to supports and visualizes quality analysis of large software systems. Fittkau et al. [12] presented a live visualization approach to monitor traces for large software landscapes. In [13], Fittkau et al. presented ExplorViz to visualize the hierarchical abstractions of large software. The goal is supporting programming comprehension tasks. Merino et al. [14] introduced an interactive software visualization tool called CityVR. The tool implements the city metaphor technique using virtual reality to support comprehension tasks.

### B. Visualizing Code Smells

The focus of this discussion is on visualizing bad smells not detecting them. Parnin and Goorg [15] presented visualizations to inspect bad coding patterns to assist developer finding relevant methods to inspect. Parnin et al. [16] proposed a catalogue of visualizations to assist reviewers to identify bad smells. They implemented a visualization tool called NOSEPRINTS. Murphy-Hill and Black [17] presented a bad smell detector called Stench Blossom. It detects and visualizes bad smells using the ambient view. The smell is shown with the source code. Mumtaz et al. [18] analyzed multivariate software metrics that link two visualization techniques, Parallel coordinates' plots and RadViz, for detecting outliers that may indicate for bad smells. Steinbeck [19] presented a visualization technique that consists of several Treemaps as a circle in order to integrate more bad smells visualizations. Carneiro et al. [20] presented a multiple views for code concern properties. These showed how these views support code smell detection.

In Summary, most of the related works in this area visualize either code elements or bad smells. We distinguished by modeling both; code and bad smells in the code with meaningful.

### III. THE PROPOSED VISUALIZATIONS

The proposed visualization models classes as buildings with the following characteristics:

- Each building consists of a number of floors equals to the number of methods in the class. The first floor is not counted. The doors of the building are shown in this floor.

- The height of each floor represents the number of lines of code (LOC) in that method.

- The number of windows in each floor equals to the number of parameters for that method.

- The number of doors, shown in the first floor, of each building equals to the number of the data fields in that class.

The bad smells are visualized as red signs on the buildings with letter avatars based on the initials with the following characteristics:

- Each avatar represents one code bad smell.

- Method related bad smells appears as avatars in the corresponding floor of that method.

- Class related bad smells appears as avatars on the roof of the building for that class

For example, Fig. 1 shows the source code for two classes; Phone and Customer from (from https://elearning.industriallogic.com). Phone class has four methods and one data field. The Customer class has one data field and one method. The two classes are modeled and visualized in Fig. 2. To illustrate the proposed visualizations, we generated the visualizations using the SketchUp (www.sketchup.com) tool. We used it to generate the visualizations in this paper based on the descriptions that we proposed.

The visualized building for the Phone class in Fig. 2 has four floors that correspond to the four methods in the class. The name of each method is shown in its corresponding floor. The first floor is not counted. We consider it level zero. It is always shown for all classes even if they have no methods. Buildings' doors are shown in level zero with the name of the class. The number of doors corresponds to the number of class's attributes. For the Phone class, one door is shown which models the single attribute of the class; unformattedNumber. The same applies for the visualized building for the Customer class. The class has one method (getMobilePhoneNumber) and one attribute (mobilePhone). It is modeled as a building with one floor and one door. The height of the first floor is five units with no windows. The floor models the getMobilePhoneNumber method that has five LOC and zero parameters.

The heights of all four floors are equal since all methods have two LOCs. The first floor has one window which means the first method has one parameter. Other methods have no parameter modeled by zero windows.

```
class Phone {
    private String unformattedNumber;

    public Phone(String unformattedNumber) {
        this.unformattedNumber = unformattedNumber;
    }
    public String getAreaCode() {
        return unformattedNumber.substring(0,3);
    }
    public String getPrefix() {
        return unformattedNumber.substring(3,6);
    }
    public String getNumber() {
        return unformattedNumber.substring(6,10);
    }
}

public class Customer…
    private Phone mobilePhone;

    public String getMobilePhoneNumber() {
        return "(" +
            mobilePhone.getAreaCode() + ") " +
            mobilePhone.getPrefix() + "-" +
            mobilePhone.getNumber();
    }
}
```

Fig. 1.   Phone and Customer Classes.

Fig. 2.    The Visualization of the Phone and Customer Classes in Fig. 1.

Data class smell results when a class has data with only setters and getters. When a method is more interested in the data of other class, this method suffers from the feature envy smell. A method that has large number of parameter has a long parameters smell. Finally, long method bad smell results when the method has too many responsibilities and performs different tasks. Shotgun surgery results when a small change occur, many other changes have to be made in many classes and methods. The avatars that model these five bad smells are shown in Fig. 3. We tried to make the design of the avatars reflects the meaning of each bad smell. We used the initials of the bad smells names. The background is red to reflect the warning status of the smell. The letter avatars was designed using a tool from (http://google-avatar.herokuapp.com/)

Using initials has two main advantages. The first one is the readability and clarity. They are easy to read and distinguished among each other. This is essential in case large number of buildings with many smells is visualized. The second advantage is the supporting of adding more smells to be visualized. Letter avatars are easy to design and visualize. So, they support the extensibility of the approach to include more smells.



Fig. 3.    The Proposed Avatars that used to Model the Bad Smells.

Method related avatars are visualized on the floors while class related smells are placed on the roofs. In Fig. 2, the Feature Envy avatar is shown on the floor that models the getMobilePhoneNumber method in the customer class. This is because this method is more interested in the data field of the Phone class and hence has the Feature Envy smell.

Fig. 4 visualizes the code shown in Fig. 5. The code in Fig. 5 shows a Java code example for a shotgun surgery smell (from http://javaonfly.blogspot.com). The visualized building shown in Fig. 4 models the class Account. It consists of four floors and three doors. The number of windows in each floor models the number of parameters for each method. The shotgun surgery avatar is shown on the roof of the building. The smell resulted from the account balance validation condition in the methods.  In case this validation is updated, all methods have to be updated.

Fig. 6 shows two more visualizations examples for two different classes.  The data class smell avatar is shown on the roof of the first building since it is a class related bad smell. The avatars of long parameter and long method smells are visualized on the second floor of the second building. This building has a very tall floor with many windows. This floor corresponds to the method that has the identified method related bad smells. The two avatars of these two smells are shown together.



Fig. 4.    The Visualization of the Account Class Shown in Fig. 5.

```
public class Account {
      private String type;
      private String accountNumber;
      private int amount;

public Account(String type, String accountNumber,
int amount){
            this.amount=amount;
            this.type=type;
            this.accountNumber=accountNumber;
      }
public void debit(int debit) throws Exception{
      if(amount <= 500){
        throw new Exception("Mininum balance
shuold be over 500");
      }
      amount = amount-debit;
      System.out.println("Now amount is" +
amount);
      }

public void transfer(Account from,Account to,int
cerditAmount) throws Exception{
            if(from.amount <= 500){
              throw new Exception("Mininum balance
shuold be over 500");
            }
            to.amount = amount+cerditAmount;
      }
public void sendWarningMessage(){
      if(amount <= 500){
        System.out.println("amount should be over
500");
            }
      }
}
```

Fig. 5.    A Class with Shotgun Surgery Bad Smell.



Fig. 6.    Two Visualizations for two different Classes with Bad Smells.

It is important to mention that the floors of the building may have different heights because the LOC of the class methods are not necessarily equal. Also, the number of windows in all floors varies because it is based on the number of parameters for each method. The visualization of the LOC and the number of parameters for methods helps in preventing bad smells in advance. By browsing buildings, developers can quickly locate methods with potential long methods and long parameter smells. These methods have tall floors and many windows. These methods need to be carefully changed to avoid smells.

## IV.  PROPOSED FRAMEWORK

In this section, the proposed automated process for generating the views is detailed. This automated process can be realized by the proposed framework shown in Fig. 7. It shows the block diagram for the main components of the proposed framework. The process starts with the source code as an input to the framework. Then, the code is analyzed to identify bad smells. The same code is also parsed to extract the code elements that will be visualized. The identified bad smells with their locations in the extracted code elements are used to generate the data model. In the next step, the generated data model is used by a visualization tool to generate the proposed visualizations. The following subsections detail the components of the proposed framework.

### A.  Code Elements Extractor

This component is responsible for parsing the source code and extracting the needed code elements. The input could be a source file or a package of files. Each source file is transformed to the XML representation srcML [21]. This representation tags each code element with its syntactic information. The srcML representation can be automatically generated by using its tool that is available from (http://www.srcml.org/).

A set of XPath queries are applied on srcML to extract all needed code elements. The first extracted elements are classes. Then for each class all its methods and its data fields are extracted. Finally, the attributes of each method are extracted. Each extracted code element is given a unique label. The given label indicates the location of the code element. For example, the label of a data field in a specific class within a specific package is written as; package (Name) .class (Name) .field (Name). Another example, the label of a method with its (LOC) value is written as: package (Name). class (Name). method (Name). parameters (Names). LOC (number). All extracted labels are sent to the Data Model Generator component.

### B.  Bad Smell Identification

Code bad smells are identified using any specialized tool. There are many tools in the literature that can be utilized. More than one tool can be used in this component to identify a variety of bad smells. JDeodorant [22][23] and JFly [24] are our target tools to be used in this component. JDeodorant is an Eclipse plug-in tool that identifies code smells with refactoring suggestions. The bad smells identified by the tool are; Feature Envy, Type Checking, Long Method, God Class and Duplicated Code. JFly is also an Eclipse plug-in tool that detects bad smells from code changes as well as static code. The tool keeps track on code changes to identify nine bad smells. The identified code smells are; Inappropriate Intimacy, Data Class, Middle Man, Message Chain, Long Parameter, Lazy Class, Brain Method, Speculative Generality and Temporary Field.

The identified bad smells are forwarded to the Data Model Generator. Each bad smell is tagged with unique label. Each label represents the name of the bad smell and the location of the affected code element. For Example, the identified feature envy bad smell in a specific class within a specific package is labeled with the following label; smell (FeatureEnvy). Package (Name). class (Name). In case a long parameters code bad smell is detected in one method, its label is written as follows; smell (LongPara). Package (Name). class (Name). method (Name). parameters (Names).

Fig. 7.    The Proposed Framework for Generating the Proposed Visualizations.

*C.  Data Model Generator*

Two types of labels are sent to the Data Model Generator component. The first type represents the extracted code elements and the other type represents the identified bad smells in code elements. The labels of code elements are used to generate the data model of the buildings. Then, the labels of bad smells are used to determine the avatars and their locations on the building.

The data model is generated automatically and mainly contains the following information:

- Locations and sizes of buildings that will be rendered on the screen. The base size of the buildings varies based on their numbers and the size of the screen.

- Specifications of buildings; the number of floors, the height of each floor, the number of floor's windows and the number of doors.

- Specifications of the avatars; type, number and location on the building.

The specifications are stored as data meta-model in a flexible XML format. The goal is to ease the rendering process using any visualization tool. For example, Fig. 8 shows a snapshot from the XML representation for the specification of one building. Each building has its own tag that corresponds to a single class. Within the building, more information is stored about its contents. This information includes the location on the screen, number of doors and the specification of each floor. Also, the avatar of the identified bad smell is stored as a tag. It is important to note that the width of the building is determined based on the number of class attributes. For example, a building that models a class with ten attributes has longer width to visualize ten adjacent doors than a building with only one door.

*D.  Visualizations Generator*

Finally, the proposed visualizations are ready to be rendered.  This component is responsible for rendering the generated data model on the screen. A tool can be developed using any programming language to draw the specifications stored in the data model. A specialized visualization tool can also be utilized to generate the visualizations.

```
<building>
      <location><x>200</x>  <y>220</y>
      <width>50</width>
      </location>
      <doors> 1 </doors>
      <floors>
      <floor> <height> 25</height>
      <windows> 10 < /windows>
      <avatar>long method</avatar> </floor>
      <floor> <height> 3 </height>
      <windows> 0 < /windows>
      <avatar>NONE</avatar> </floor>
      </floors>
</building>
```

Fig. 8.    A Snapshot for the XML Representation of One Building.

We are working on developing a visualization tool to generate the views from the data model. The tool will be able to generate 3D visualizations for buildings. Zooming, localization and browsing are essential features that are under consideration. Developers will have the ability to search and locate a specific building that corresponds to a specific class. Also, they will be able to zoom in or out the buildings. The browsing feature helps developers to navigate through buildings in 3D environment. These features help developers to understand and handle large number of buildings that model large scale systems with many classes.

V.  Evaluation

We need to evaluate the usefulness of the proposed visualizations in supporting comprehension tasks. So, we performed a controlled pilot experiment on software engineering undergraduate students. The goal of the experiment is to check how the proposed visualizations help in quickly understand and locate bad smells in code. The steps of the experiments were as follows:

*1)* Five java classes were carefully selected and implemented to have intentionally bad smells. The five bad smells that are under consideration in this paper where distributed over the five classes.

*2)* The data model of the classes was generated to model the classes based on the proposed visualizations.

*3)* The visualizations were generated using SketchUp and based on the data model of the classes.

*4)* Four software engineering students who are familiar with code bad smells were divided equally into two groups. All students have very good GPA rating.

*5)* The first group was given the five classes with textual report about the bad smells and their locations in these five classes.

*6)* The second group was given the same five classes and their visualizations with their bad smells.

*7)* Each student was asked to write down why the dedicated bad smell has occurred.

*8)* The answers of the two grouped were compared and the average completion time.

Table I shows the completion time in minutes for all four students on the two groups. The average completion time for the first group is 4.9 minutes. The first student needed five

minutes to complete the task while the second student needed four minutes and 48 second (4.8 minutes). The second group, who used the visualizations, achieved better average time which is three minutes.

The comparison results between the two groups showed also that both groups answered the questions correctly. But the average completion time was different. The group who used the visualizations completed the task with about 40% less average time than the other group.

After completing the experiment, we also, asked the four students if they find the visualizations useful in understanding code smells. Three out the four students found it useful. The fourth student found the design of the avatars is not useful in modeling the bad smells.

TABLE I. THE COMPLETION TIME FOR THE SUBJECTS

| Group | Student | Time in Minutes | Average |
|---|---|---|---|
| Without Visualizations | S1 | 5 | 4.9 |
| | S2 | 4.8 | |
| With Visualizations | S3 | 2.5 | 3 |
| | S4 | 3.5 | |

## VI. CONCLUSIONS AND FUTURE WORK

Software visualization support program comprehension tasks for maintainers. Useful visualizations have been proposed to help developers locate and identify bad smells. These visualizations show object oriented code elements with the types of identified bad smells in these code elements. A framework is presented to automatically analyze source code and generate the proposed visualizations. The framework is automated and can be extended to include more smells and visualize more code elements. The evaluation of the proposed visualizations showed their positive impact on understanding bad smells and their causes.

Our future work aims to completely implement the framework and realize it as a plug-in tool in an IDE as Eclipse. More avatars will be considered to cover more bad smells. We are also working on visualizing more code elements as relationships among classes, data types and methods invocations.

## REFERENCES

[1] M.Fowler, Refactoring improving the design of existing code .Addison-Wesley, 1999.

[2] Ducasse, Stéphane, and Michele Lanza. "The class blueprint: visually supporting the understanding of glasses." IEEE Transactions on Software Engineering 31, no. 1 (2005): 75-90.

[3] Lanza, Michele, and Stéphane Ducasse. "A categorization of classes based on the visualization of their internal structure: the class blueprint." ACM SIGPLAN Notices 36, no. 11 (2001), pp. 300-311.

[4] Wettel, Richard, and Michele Lanza. "Visualizing software systems as cities." In 2007 4th IEEE International Workshop on Visualizing Software for Understanding and Analysis, pp. 92-99. 2007.

[5] Wettel, Richard, and Michele Lanza. "Codecity: 3d visualization of large-scale software." In Companion of the 30th international conference on Software engineering, pp. 921-922. 2008.

[6] Wettel, Richard, Michele Lanza, and Romain Robbes. "Software systems as cities: A controlled experiment." In 2011 33rd International Conference on Software Engineering (ICSE'11), pp. 551-560. 2011.

[7] Alam, Sazzadul, and Philippe Dugerdil. "Evospaces visualization tool: Exploring software architecture in 3d." In 14th Working Conference on Reverse Engineering (WCRE 2007), pp. 269-270. 2007.

[8] Panas, Thomas, Rebecca Berrigan, and John Grundy. "A 3d metaphor for software production visualization." In Proceedings on Seventh International Conference on Information Visualization (IV 2003), pp. 314-319. 2003.

[9] Kienle, Holger M., and Hausi A. Müller. "Rigi—An environment for software reverse engineering, exploration, visualization, and redocumentation." Science of Computer Programming 75, no. 4 (2010), pp. 247-263.

[10] Marcus, Andrian, Louis Feng, and Jonathan I. Maletic. "3D representations for software visualization." In Proceedings of the 2003 ACM symposium on Software visualization, p. 27. 2003.

[11] Langelier, Guillaume, Houari Sahraoui, and Pierre Poulin. "Visualization-based analysis of quality for large-scale software systems." In Proceedings of the 20th IEEE/ACM International Conference on Automated software engineering (ASE'05), pp. 214-223, 2005.

[12] Fittkau, Florian, Jan Waller, Christian Wulf, and Wilhelm Hasselbring. "Live trace visualization for comprehending large software landscapes: The ExplorViz approach." In 2013 First IEEE Working Conference on Software Visualization (VISSOFT), pp. 1-4. 2013.

[13] Fittkau, Florian, Alexander Krause, and Wilhelm Hasselbring. "Software landscape and application visualization for system comprehension with ExplorViz." Information and software technology 87 (2017).pp. 259-277.

[14] Merino, Leonel, Mohammad Ghafari, Craig Anslow, and Oscar Nierstrasz. "CityVR: Gameful software visualization." In 2017 IEEE International Conference on Software Maintenance and Evolution (ICSME), pp. 633-637, 2017.

[15] Parnin, Chris, and Carsten Görg. "Lightweight visualizations for inspecting code smells." In Proceedings of the 2006 ACM symposium on Software visualization, pp. 171-172, 2006.

[16] Parnin, Chris, Carsten Görg, and Ogechi Nnadi. "A catalogue of lightweight visualizations to support code smell inspection." In Proceedings of the 4th ACM symposium on Software visualization, pp. 77-86. 2008.

[17] Murphy-Hill, Emerson, and Andrew P. Black. "An interactive ambient visualization for code smells." In Proceedings of the 5th international symposium on Software visualization, pp. 5-14. 2010.

[18] H. Mumtaz, F. Beck and D. Weiskopf, "Detecting Bad Smells in Software Systems with Linked Multivariate Visualizations," In 2018 IEEE Working Conference on Software Visualization (VISSOFT'18), pp. 12-20, 2018.

[19] Steinbeck, Marcel. "An arc-based approach for visualization of code smells." In 2017 IEEE 24th International Conference on Software Analysis, Evolution and Reengineering (SANER'17), pp. 397-401. 2017.

[20] Carneiro, Glauco de F., Marcos Silva, Leandra Mara, Eduardo Figueiredo, Claudio Sant'Anna, Alessandro Garcia, and Manoel Mendonca. "Identifying code smells with multiple concern views." In 2010 Brazilian Symposium on Software Engineering, pp. 128-137, 2010.

[21] Collard, M. L. , Kagdi H. H., Maletic, J. I., "An XML-based lightweight C++ fact extractor," Proc. of 11th IEEE International Workshop on Program Comprehension (IWPC'03), pp. 134-143, 2003.

[22] Fokaefs, Marios, Nikolaos Tsantalis, Eleni Stroulia, and Alexander Chatzigeorgiou. "JDeodorant: identification and application of extract class refactorings." In 2011 33rd International Conference on Software Engineering (ICSE'11), pp. 1037-1039, 2011.

[23] Mazinanian, Davood, Nikolaos Tsantalis, Raphael Stein, and Zackary Valenta. "JDeodorant: clone refactoring." In 2016 IEEE/ACM 38th International Conference on Software Engineering Companion (ICSE-C), pp. 613-616. 2016.

[24] Maen Hammad, Asma Labadi, "Automatic Detection of Bad Smells from Code Changes", International Review on Computers and Software, Vol. 11, No. 11, pp. 1016-1027, 2016

# Novel Carrier based PWM Techniques Reduce Common Mode Voltage for Six Phase Induction Motor Drives

Ngoc Thuy Pham[1]
Dept. of ET, Industrial University of Ho Chi Minh
Ho Chi Minh City, Viet Nam

Nho Van Nguyen[2]
Dept. of EEE, Ho Chi Minh University of Technology
Ho Chi Minh City, Viet Nam

*Abstract*—This paper proposes a novel pulse width modulation (CBPWM) technique for reducing the common mode voltage for a six-phase induction motor (SPIM) drive. This proposed CBPWM technique relies on setting up offset functions and the phase shift of carrier wares. Common mode voltage occurs under the effect of DC power Vd always in Vd/6 limits. Some ways of designing the offset function are proposed; these proposed strategies permit to reduce either the mean value or the instantaneous value of the common mode voltage. Features of proposal CBPWM solutions have been compared. Simulation and experimental results demonstrate the feasibility of the proposed solution.

*Keywords—Six-phase induction motor; six-phase voltage source inverter; common mode voltage; carrier based pulse width modulation*

## I. INTRODUCTION

In recent decades, multi-phase motors have become increasingly popular, especially in medium to large power applications such as automotive, aerospace, military and nuclear [1,2]. The use of multi-phase drives has been considered an effective approach to achieving high power without increasing the stator currents per phase. Among the numerous possibilities of multiphase ac machines, SPIM is probably the most popular in industrial applications. Nowadays, SPIM are even considered for small power in all applications requiring reliability and fault tolerance [3], [4]. In this way, it is expected that the loss of one or more phases allows the machine to provide a significant electromagnetic torque to run the system. On the other hand, the high performance of modern power converters in terms of switching frequency and control capability can be used to reduce torque oscillations in the case of phase loss [5]. In the SPIM drive, the use of a six-phase inverter (6P_VSI) as a necessary choice because in fact, the six-phase power supply is not available. As we all know, in order to control 6P_VSI, different pulse width modulation (PWM) techniques are employed to achieve the voltage quality criterion. The classic PWM techniques such as space vector modulation (SVPWM), continuous carrier modulation (CPWM) technique, and the sine modulation (SINPWM) technique is commonly used to generate the good voltage quality, voltage distortion caused by harmonics is low. Discontinuous modulation (DPWM) techniques have also been developed, which allows reducing the switching losses, higher tolerance distortion. However, the above PWM modulation methods cause high common mode pulse. This common mode voltage generates the parasitic current components that appear between the metal part of the stator and the rotor, between the stator and the housing. As a consequence of this common mode voltage effect, electromagnetic interference (EMI) current appears to corrode the bearing surface of the motor, on the other hand, EMI current with the interference activation of protective devices and heating the wires [6-9]. Thus, the application of six-phase inverter does not completely satisfy the high- quality requirements of practical electric drives, especially when DC power is high [6], [7], [9]. The goal of the inverter voltage control is to suppress the negative influence caused by common mode voltage. Some practical solutions use serial reactor with the inverter output or use a hardware circuitry contain semiconductor switch to control common mode voltage compensation [10,11]. However, techniques use expensive hardware circuits and even reduce system reliability. Thus, the current trend is more concerned with PWM techniques that use the space vector and carrier based on techniques that reduce or eliminate the common mode voltage [12]. Although the SPIM drive has been studied for a quite long time, common research results for common-mode voltage reduction are less well known than the common mode reduction of three-phase induction motor drive [9], [13, 14]. One of the few common uses of PWM technology is the reduction of common mode for symmetric six phase inverters or 5-phase inverters [1], [15-16]. The 6P_VSI is capable of controlling the common mode voltage reduction/elimination because the space vector schema contains a number of zero common-mode voltages. Compared with space vector PWM technique, carrier based on PWM techniques have available some conveniences as less calculating, apply in cases extending PWM techniques for power conversion systems easily such as three phase multilevel inverters, or multiphase inverter. For that purpose, the paper proposes a common mode voltage reduction technique using CBPWM methods, called RCMV CBPWM (Reduced Common Mode Voltage CBPWM). The phenomenon of switching simultaneously at two inverter branches also leads to a large change in dv/dt between phases. Hence, it creates a load current has large noise peak and causes leakage currents through parasitic capacitors between phases in the cables. On the other hand, the PWM techniques suppress the CMV have limited range of output voltage. Therefore, another possible solution is researched to reduce CMV.

The paper presented the new CBPWM technique for controlling voltage reduction common mode for six phase inverter. The simple implementation method by applying offset function combinate with carrier technique in six phase control voltage. These 4S-CBPWM techniques control and reduce the common mode voltage amplitude in the range Vd/6. Compared to the space vector modulation techniques, the carrier based PWM modulation technique are less computation, easy to apply when expanding PWM technology for power conversion systems such as multi-level inverter or multi-phase inverter.

In this paper, the new CBPWM techniques are presented for controlling and reducing the common mode voltage for six phase inverter. The methods implemented by applying offset function combine with carrier technique. These 4S-CBPWM techniques control and reduce the common mode voltage amplitude in the range Vd/6. Compared to the space vector modulation techniques, the carrier based PWM modulation technique are less computation, easy to apply when expanding PWM technology for power conversion systems such as multi-level inverter or multi-phase inverter. The achieved simulation and experiment results demonstrate the effectiveness of the proposed CBPWM RCMV techniques.

This paper is organized into five sections. In section 2, the basic theory of the model of SPIM drive and SPIM are presented. Section 3 introduces the principle of proposed RCMV CBPWM. Simulation and experiment results are presented in Sections 4 and 5. Finally, the concluding is provided in Section 5

## II. MODEL OF SPIM DRIVE AND SPIM

The system under study consists of SPIM fed by a six-phase VSI (Voltage Source Inverter) and a DC link. A detailed scheme of the drive is provided in Fig. 1. This SPIM is a continuous system that can be described by a set of differential equations. The model of the system can be simplified by means of the vector space decomposition (VSD). By applying this technique, the original six-dimensional space of the machine is transformed into three two-dimensional orthogonal subspaces in the stationary reference frame (α-β), (x-y) and (zl -z2). This transformation is obtained by means of 6 x 6 transformation matrix:

$$T_6 = \frac{1}{3} \begin{bmatrix} 1 & \cos\gamma & -\frac{1}{2} & \cos(\frac{2\pi}{3}+\gamma) & -\frac{1}{2} & \cos(\frac{4\pi}{3}+\gamma) \\ 0 & \sin\gamma & \frac{\sqrt{3}}{2} & \sin(\frac{2\pi}{3}+\gamma) & -\frac{\sqrt{3}}{2} & \sin(\frac{4\pi}{3}+\gamma) \\ 1 & \cos(\pi-\gamma) & -\frac{1}{2} & \cos(\frac{\pi}{3}-\gamma) & -\frac{1}{2} & \cos(\frac{5\pi}{3}-\gamma) \\ 0 & \sin(\pi-\gamma) & -\frac{\sqrt{3}}{2} & \sin(\frac{\pi}{3}-\gamma) & \frac{\sqrt{3}}{2} & \sin(\frac{5\pi}{3}-\gamma) \\ 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 \end{bmatrix}$$

(1)

In that, an amplitude invariant criterion was used. From the motor model obtained by using the VSD approach, the following conclusions should be emphasized:

- The electromechanical energy conversion variables are mapped to the (α-β) subspace. Therefore, the fundamental supply component as well as the supply harmonics of order 12n ± 1 (n = 1,2,3,...), are

represented in this subspace. The non-electromechanical energy conversion variables can be found in other subspaces.

- The current components in the (x-y) subspace do not contribute to the air gap flux and are limited only by the stator resistance and stator leakage inductance, which are usually small. These components represent the supply harmonics of the order 6n ± 1 (n = 1,3,5,...) and only produce losses, so consequently they should be controlled to be as small as possible.

- The voltage vectors in the (zl -z2) are zero due to the separated neutrals configuration of the machine.



Fig 1. A general scheme of SPIM drive.

A 3P_VSI has a discrete nature, actually, it has a total number of $2^6 = 64$ different switching states defined by six switching functions corresponding to the six inverter legs [Sa, Sx, Sb, Sy, Sc, Sz], where Si ϵ {0,1}. The different switching states and the voltage of the DC link define the phase voltages which can turn be mapped to the (α-β) - (x-y) space according to the Vector space decomposition VSD approach. For this reason, the 64 different on/off combinations of the six legs of VSI lead to 64 space vectors in the (α-β) and (x-y) subspaces. Fig. 2 shows the active vectors in the (α-β) and (x-y) subspaces, where each vector switching state is identified using the switching function by two octal numbers corresponding to the binary numbers [SaSbSc] and [SxSySz], respectively.



Fig 2. Voltage space vectors and switching states in the (α-β) and (x-y) subspaces for a six-phase asymmetrical VSI.

Fig. 2 the 64 possibilities lead to only 49 different vectors in the (α-β) - (x-y) subspace. On the other hand, a transformation matrix must be used to represent the stationary reference frame (α-β) in the dynamic reference (d - q). This matrix is given:

$$T_{dq} = \begin{bmatrix} \cos(\delta_r) & -\sin(\delta_r) \\ \sin(\delta_r) & \cos(\delta_r) \end{bmatrix} \qquad (2)$$

Where δ is the rotor angular position referred to the stator as shown in Fig. 1.

A SPIM which contains two sets of three-phase winding spatially are shifted 30 electrical degrees with isolated neutral points or double neutral point, as depicted in Fig. 3, is modeled. Stator and rotor voltage equation for this model is as follows:

$$[V_s] = [R_s][I_s] + P([L_s][I_s] + [M][I_r]) \qquad (3)$$

$$[V_r] = [R_r][I_r] + P([L_r][I_r] + [M][I_s]) \qquad (4)$$

Where: [V], [I], [R], [L] and [M] are voltage, current, resistant, self, and mutual inductance vectors, respectively. P is a differential operand. Subscript r, s related to the rotor and stator vectors respectively. Since the rotor is squirrel cage, [Vr] is equal to zero.



Fig 3. Distribution of coils in six-phase motors.

Applying the transformation matrix [T6], where γ = π / 6, the six-phase motor can be represented in three two-dimensional space coordinates: (D-Q), (x, y) and ( z1, z2):

(D-Q) subspaces :

$$\begin{bmatrix} V_{sD} \\ V_{sQ} \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} R_s + PL_s & 0 & PM & 0 \\ 0 & R_s + PL_s & 0 & PM \\ PM & \omega_r M & R_r + PL_r & \omega_r L_r \\ -\omega_r M & PM & -\omega_r L_r & R_r + PL_r \end{bmatrix} \begin{bmatrix} I_{sD} \\ I_{sQ} \\ I_{rD} \\ I_{rQ} \end{bmatrix} \qquad (5)$$

(x,y) subspaces :

$$\begin{bmatrix} V_{sx} \\ V_{sy} \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} R_s + PL_s & 0 & 0 & 0 \\ 0 & R_s + PL_s & 0 & 0 \\ 0 & 0 & R_r + PL_r & 0 \\ 0 & 0 & 0 & R_r + PL_r \end{bmatrix} \begin{bmatrix} I_{sx} \\ I_{sy} \\ I_{rx} \\ I_{ry} \end{bmatrix} \qquad (6)$$

(z1,z2) subspaces:

$$\begin{bmatrix} V_{sz1} \\ V_{sz2} \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} R_s + PL_s & 0 & 0 & 0 \\ 0 & R_s + PL_s & 0 & 0 \\ 0 & 0 & R_r + PL_r & 0 \\ 0 & 0 & 0 & R_r + PL_r \end{bmatrix} \begin{bmatrix} I_{sz1} \\ I_{sz2} \\ I_{rz1} \\ I_{rz2} \end{bmatrix} \qquad (7)$$

where: Ls=Lls+M, Lr= Llr+M, M=3.Lms, P=d/dt. As these equations imply, the electromechanical conversion only takes place in the α-β subspace (DQ subspace) and the other subspaces just produce losses.

The torque equation can be written as follows:

$$T_e = 3P(\Psi_{\beta r} i_{\alpha r} - \Psi_{\alpha r} i_{\beta r}) \qquad (8)$$

$$J_i \frac{d}{dt}\omega_r + B_i \omega_r = P(T_e - T_L) \qquad (9)$$

where Ji, ωr, Bi, Tm, TL, P: are inertia coefficient, angular speed, friction factor, the electromagnetic torque that generated by the motor, load torque, number of poles and stator flux linkage at the related subspace.

### III. PRINCIPLE OF RCMV PWM

#### A. CBPWM Technique for 6P_VSI:

*1) SVPWM Technique for 6P_VSI*: Fig. 4 describes the principle of implementing SVPWM for 6P_VSI. A 6P_VSI contains two the three phase VSI (3P_VSI I and 3P_VSI II) and they are independently PWM controlled. The reference voltage vector is carried out in SVPWM of 3P_ VSI I, which generates the control pulse ( $S_a S_b S_c$ ), is shifted by 30 degrees compared to the reference voltage vector of VSI II. This reference voltage vector is carried out SVPWM of 3P_ VSI II generate the control pulse ( $S_A S_B S_C$ ).



Fig 4. SVPWM technique for SPIM drive.

*2) CBPWM Technique for 6P_VSI*: Fig. 5 describes the principle of implementing CBPWM for 6P_VSI. Typically, each 3P_VSI I and II can be independently controlled PWM. The control voltage signal of 3P_VSI I is generated by synthesizing the three-phase basic voltage and the offset voltage, This is then compared to the carrier carried out in the CBPWM I block (Fig. 6) that generated the control pulse (s_a, s_b, s_c).



Fig 5. CBPWM technique for SPIM drive.

Fig 6.   Detail CBPWM block for 3P_VSI I.

Similarly, the basic control voltage signal of 3P_VSI II has a phase angle shift 30 degrees compared to VSI I, This voltage compared to the carrier carried out in the CBPWM II block that generated the control pulses for three phase VSI II ($s_A$, $s_B$, $s_C$). The CB PWM block generates the excitation pulse $s_a s_b s_c$ by comparing the control voltage $v_{dka}, v_{dkb}, v_{dkc}$ to the carrier waves V carrier. The relationships are described by:

$$v_{dka} V_d = V_{refa} + V_{comI} + \frac{V_d}{2}$$
$$v_{dkb} V_d = V_{refb} + V_{comI} + \frac{V_d}{2} \qquad (10)$$
$$v_{dkc} V_d = V_{refc} + V_{comI} + \frac{V_d}{2}$$

Similar, we also can set-up relationship between the control functions for 3P_VSI II:

$$v_{dkA} V_d = V_{refA} + V_{comII} + \frac{V_d}{2}$$
$$v_{dkB} V_d = V_{refB} + V_{comII} + \frac{V_d}{2} \qquad (11)$$
$$v_{dkC} V_d = V_{refC} + V_{comII} + \frac{V_d}{2}$$

The reference voltage for the phases of VSI I and II are described as follows:

$$V_{refa} = V_m \cos\theta$$
$$V_{refb} = V_m \cos(\theta - \frac{2\pi}{3}) \qquad (12)$$
$$V_{refc} = V_m \cos(\theta - \frac{4\pi}{3})$$

$$V_{refA} = V_m \cos(q - \frac{p}{6})$$
$$V_{refB} = V_m \cos(q - \frac{2p}{3} - \frac{p}{6}) \qquad (13)$$
$$V_{refC} = V_m \cos(q - \frac{4p}{3} - \frac{p}{6})$$

The configuration of 6P_VSI is not symmetrical with two neutral points of two 3P_VSI I and II isolation. The voltage of phase a, b, c called sequentially $v_{a0}, v_{b0}, v_{c0}$ and $v_{A0}, v_{B0}, v_{C0}$. $V_{comI}$ and $V_{comII}$ are common mode voltage of 3P_VSI I and II. Calling the switching status sequentially: $s_a s_b s_c$ and Vcom for 6P_VSI is defined:

$$v_{comI} = \frac{v_{a0} + v_{b0} + v_{c0}}{3} - \frac{V_d}{2} \qquad (14)$$
$$v_{comII} = \frac{v_{A0} + v_{B0} + v_{C0}}{3} - \frac{V_d}{2}$$

Assuming that the three-phase stator coils are arranged symmetrically, the common mode voltage can be calculated as follows [1]:

$$v_{com} = \frac{v_{comI} + v_{comII}}{2} \qquad (15)$$

$$v_{com} = \frac{v_{a0} + v_{b0} + v_{c0} + v_{A0} + v_{B0} + v_{C0}}{6} - \frac{V_d}{2} \qquad (16)$$

Description of common mode voltage to the switching states of the semiconductor switch, we have:

$$v_{comI} = \left( \frac{s_a + s_b + s_c}{3} - \frac{1}{2} \right) V_d \qquad (17)$$

$$v_{comII} = \left( \frac{s_A + s_B + s_C}{3} - \frac{1}{2} \right) V_d \qquad (18)$$

$$v_{com} = \left( \frac{s_a + s_b + s_c + s_A + s_B + s_C}{6} - \frac{1}{2} \right) V_d \qquad (19)$$

In order to reduce Vcom, the switching state of the switch are designed to Vcom reaches the smallest values. One of the possibilities, that is each of the 3P_VSI I and II will be controlled independently to minimize Vcom.

*B. RCMV 4S-PWM Technique for 6P_VSI:*

*1) RCMV 4S- SVPWM technique*: A 6P_VSI contains two the three phase VSI (3P_VSI I and 3P_VSI II) and they are independently PWM controlled. To analyze RCMV 4S-SVPWM technique we have used a 3P__VSI scheme as Fig. 7. Vcom can be defined as follows (see Fig. 7).



Fig 7.   Diagram 3__VSI model

Fig 8.    Space voltage vector schema in RCMV 4S-SVPWM.



Fig 9.    The average voltage model applied to one phase of VSI (VA=VrefA, Voff=voff*Vd).

The state Sj of the semiconductor switch can be set to get the 0 or 1 values, respectively with the sates of semiconductor switch conduct or dis-conduct. According to formula (1), Vcom can reach Vd/6, Vd/3, and Vd/2 values. The voltage vector states (000) and (111) will produce the maximum common mode voltage (Vd/2); Active voltage vectors produce low common mode voltages (Vd/6, Vd/3). In order to achieve the purpose reducing common mode voltage (RCMV), the PWM technique eliminates the (000) and (111) vector states in the tripping state. Therefore, we call the PWM technique called RCMV-PWM. Obviously, the RCMV PWM method has a state schema only containing positive voltage vectors.

The PWM technique reduces CMV can be implemented by the space vector modulation technique, called RCMV SVPWM. The vectors state zero will be replaced by the sum of the two positive vectors that are symmetric across the center of the hexagon, as shown in Fig. 8. For example, when considering the Vref vector in the first sector of a hex, the zero state vector can be replaced by one of the sums of the vectors (V100 + V011), (V110 + V010) and (V010 + V101).

In these three alterations, the combine two vectors (V101 + V010) are closest to the area that contains Vref vector, the RCMV SVPWM technique uses the four closest vectors to the Vref vector, Helps to reduce the common mode voltage. Because the SVPWM technique uses Four State vector (4S) with reduced common-mode voltage, we call this the RCMV 4S SVPWM.

*2) RCMV 4S-CBPWM technique*: Fig. 5 describes the principle of implementing CBPWM for 6P_VSI. The PWM control principle reduces the independent common mode voltage will be illustrated for the VSI I (Fig. 6). The RCMV CBPWM technique base on the way generates excitation pulse sj by comparing the three phase control voltage vdkA, vdkB and vdkC with the triangular carrier. The RCMV technical principle is based on two main characteristics:

    _  Determines the offset functions and control voltages;

    _  Carrier technology

    *a) Determine the offset functions and control voltages*

The three-phase control voltage vdkA, vdkB and vdkC can be deduced from the three-phase average voltage model ( Fig. 9 illustrates the average voltage model applied to a phase of VSI). We have the output voltage of VSI are: vA0, vB0, vC0:

$$V_{A0} = v_{dkA}V_d = V_{refA} + V_{com} + \frac{V_d}{2}$$

$$V_{B0} = v_{dkB}V_d = V_{refB} + V_{com} + \frac{V_d}{2} \qquad (20)$$

$$V_{C0} = v_{dkC}V_d = V_{refC} + V_{com} + \frac{V_d}{2}$$

The offset function in the average voltage model can be determined by the formula:

$$v_{off} = \frac{v_{dkA} + v_{dkB} + v_{dkC}}{3} \qquad (21)$$

$$v_{off}.V_d - \frac{V_d}{2} = V_{com} \qquad (22)$$

The maximum, minimum, medium values of the three-phase reference base voltage components are called Max, Min and Mid, these components required to be normalized according to the source Vd:

$$Max = \frac{Maximum(v_{refA}, v_{refB}, v_{refC})}{V_d} \qquad (23)$$

$$Min = \frac{Minimum(v_{refA}, v_{refB}, v_{refC})}{V_d} \qquad (24)$$

$$Mid = M_{max} - M_{min} \qquad (25)$$

The offset voltage can be inferred in the range between the voffmax and voffmin, it determined as follows:

$$v_{offMax} = 1 - Max \qquad (26)$$

$$v_{offMin} = -Min \qquad (27)$$

Two maximum and minimum offset function values correspond to the average Vcom values as follows:

$$V_{comMax} = (0.5 - Max)V_d \qquad (28)$$

$$V_{comMin} = (-0.5 - Min)V_d \qquad (29)$$

Max and Min are the maximum and minimum values of the three-phase control voltages. This Vcom indicates limits of instantaneous Vcom may occur at the working point of 3P_VSI.

*3) The conditions set up the offset function to reduced Vcom*: In order, states string ensure limit Vcom as shown in Fig. 10, it is necessary to set up the constraint conditions of the offset function (22). Besides the general limitations of carrier technology (23), (24); The offset function needs to be more strictly limited under the condition that it generates the state order in Fig. 10. Calling these limits are eoMin and eoMax, which correspond to common mode ecomMin and ecomMax. We have:

$$v_{dkMax} + v_{dkMid} > 1 \tag{30}$$

$$v_{dkMid} + v_{dkMin} < 1 \tag{31}$$

where:

$$v_{dkMax} = Max\ (v_{dkA}, v_{dkB}, v_{dkC}) \tag{32}$$

$$v_{dkMin} = Min\ (v_{dkA}, v_{dkB}, v_{dkC}) \tag{33}$$

$$v_{dkMid} = v_{dkA} + v_{dkB} + v_{dkC} - v_{dkMax} - v_{dkMin} \tag{34}$$

The ecomMin and ecomMax values are the determine functions that depend on the working position of the Vref space vector. ecomMin and ecomMax are related to the common mode extremes of CBPWM (28), (29) as follows:

$$v_{comMin} \le e_{comMin} \le v_{com} \le e_{comMax} \le v_{comMax} \tag{35}$$

Similar, we also can set-up relationship between the control functions for 3P_VSI II.

To perform the state sequences (101), (100), (110), (010) in the sampling cycle Ts, it is necessary to establish the difference between the carriers of the phases. For example, in the first sector, the carrier of phase B is shifted 180 degrees to the carriers of phase A and C, see Fig. 10.



Fig 10. Excitation pulse scheme of RCMV 4S CBPWM technical for 3P_VSI.

The problem decides to the implement of 4S-CBPWM technique to reduce Vcom of the 6P_VSI is set up the average common mode voltage $V_{comI}$, $V_{comII}$ in formulas (16).

It can be said that to implement RCMV PWM reduces common mode voltage using four close state vectors (4S-CBPWM). Common mode function values setup for two 3P_VSI I and II should meet the following conditions:

$$e_{comMinI} \le v_{comI} \le e_{comMaxI} \tag{36}$$

$$e_{comMinII} \le v_{comII} \le e_{comMaxII} \tag{37}$$

Set the limits for common mode voltage of 6P_VSI when applying 4S-CBPWM:

$$e_{comMin} = \frac{e_{comMinI} + e_{comMinII}}{2} \tag{38}$$

$$e_{comMax} = \frac{e_{comMaxI} + e_{comMaxII}}{2} \tag{39}$$

CMV of the 6P_VSI can be set up any value within the following limits:

$$e_{comMin} \le V_{com} \le e_{comMax} \tag{40}$$

The proposed method satisfies if the common mode voltage value CMV satisfies (40) to implement the RCMV CBPWM technique. The values will need to be determined to implement the PWM technique for the VSI I and II.

*C. Proposal RCMV CBPWM Techniques*

*1) RCMV CBPWM technique with average common mode voltage VcomMid*: All common mode function values are set up individually and they allow reducing $V_{com}$. This paper introduces two proposals of the choice of common mode function and its effects on the quality of six phase drive systems.

$$Max_I = \frac{Maximum(v_{refa}, v_{refb}, v_{refc})}{V_d} \tag{41}$$

$$Min_I = \frac{Minimum(v_{refa}, v_{refb}, v_{refc})}{V_d} \tag{42}$$

$$Mid_I = -Max_I - Min_I \tag{43}$$

Similarly, for VSI II, we have:

$$Max_{II} = \frac{Maximum(v_{refA}, v_{refB}, v_{refC})}{V_d} \tag{44}$$

$$Min_{II} = \frac{Minimum(v_{refA}, v_{refB}, v_{refC})}{V_d} \tag{45}$$

$$Mid_{II} = -Max_{II} - Min_{II} \tag{46}$$

The common mode function for two VSI I and II performs to reduce common mode voltage by the formula:

$$v_{comI} = v_{comMidI} = Mid_I \cdot \frac{V_d}{2} \tag{47}$$

$$v_{comII} = v_{comMidII} = Mid_{II} \cdot \frac{V_d}{2} \tag{48}$$

Average value of common mode voltage of 6P_VSI:

$$v_{com} = v_{comMid} = (Mid_I + Mid_{II}) \cdot \frac{V_d}{4} \tag{49}$$

It can be seen that the functions (47), (48), (49) satisfy conditions (36), (37) and (40), and the CBPWM technique with generating impact time of two farther vectors are the same. Two further vectors have effect as two vectors zero; therefore, harmonics will be less distorted. The relationship between $V_{commax}$, $V_{commin}$ and $V_{comMid}$ is described in Fig. 11.(a), (b). The 6P_VSI use the RCMV PWM technique with $V_{comMid}$ always contains the average third harmonic mode common component with increasing magnitude according to the modulation index.



(a)

(b)

Fig 11. (a):Correlations diagram the average values of common mode voltage $e_{comMax}$ (green), $VcomMid$ (red) and $e_{comMin}$ (blue) when changing modulation index m=0.2;0.4;0.6;0.8 and 1. (b):Correlations diagram the average values of common mode voltage $e_{comMax}$ (green), $VcomMid$ (red) and $e_{comMin}$ (blue) during m=1.

In formula (47), (48 (49) the Mid common mode voltage function is proportional to load phase voltage and is between the largest and smallest phase. Vm be called the amplitude of the load phase, the amplitude of the vcom varies depending on the magnitude of the load voltage and this voltage will greater when larger load voltage, as explained below:

$$V_{comMid} = V_m \sin(30^0) = 0.5V_m \tag{50}$$

The peak amplitude of the average common mode voltage is equal to half the amplitude of the base component of the output voltage. Therefore, as the load voltage increases, this method will make increase the triple harmonic of the common mode voltage component for the drive system.

Fig. 12 describes FFT analysis of the triple harmonic of VcomMid for m = 0.4.



Fig 12. Correlations diagram the average values of common mode voltage $e_{comMax}$ (green), $V_{comMid}$ (red) and $e_{comMin}$ (blue) when changing modulation index m=0.2; 0.4;0.6;0.8 and 1.

*2) 4S-CBPWM technique with the minimum common mode voltage $V_{comOpt}$*: Recent, some published results showed that the existence of the triple harmonic component of Vcom would cause a high voltage stress on the motor particular when the inverter connected to the motor with a relatively long cable [14]. Therefore, in order to limit the triple harmonic component of CMV, we can choose the optimal and minimum common mode voltage in the equation (22) as follows:

$$V_{comOpt} = Min(V_{com}) \tag{51}$$

In the range can control SIN PWM technique, mean the average value $V_{com} = 0$. Thus, the output voltages don't contain the common mode triple harmonic component. Outside of this range, the $V_{com}$ values are chosen minimum from (22), (23) and (24). Detail, we will define the voltage $\mathbf{V}_{comOpt}$ as the case of CBPWM having the function $V_{com}$ reach the least absolute value. Obviously, the special case of the CBPWM technique with low m modulation index is the SIN PWM technique. We have $v_{comI} = v_{comII} = 0$. In the larger modulation index area, For example, when controlling VSI I independently with modulation index m=1 conditions $e_{comMinI} \le v_{comI} = 0 \le e_{comMaxI}$ cannot occur. Similarly, apply to VSI II.

The optimum mode common mode function has minimal common mode value that can be set as follows:

$$v_{comOpt} = \begin{cases} e_{comMin} & \text{if} & e_{comMin} > 0 \\ e_{comMax} & \text{if} & e_{comMax} < 0 \\ 0 & \text{if} & e_{comMin} < 0 < e_{comMax} \end{cases} \tag{52}$$

It is easy to see that defining of CMV function according to two conditions of the (53) relation is easy to solve:

$$\left. \begin{array}{l} v_{comI} = e_{comMinI} \\ v_{comII} = e_{comMinII} \end{array} \right\} if \quad v_{comOpt} = e_{comMin} \tag{53}$$

And:

$$\left. \begin{array}{l} v_{comI} = e_{comMaxI} \\ v_{comII} = e_{comMaxII} \end{array} \right\} if \quad v_{comOpt} = e_{comMax} \tag{54}$$

Fig 13.  Correlations diagram the average values of common mode voltage $e_{comMax}$ (green), $V_{comMid}$ (red) and $e_{comMin}$ (blue) when changing modulation index m=0.2; 0.4;0.6;0.8 and 1.



Fig 14.  Correlations diagram the average values of common mode voltage $e_{comMax}$ (green), $V_{comMid}$ (red) and $e_{comMin}$ (blue) during m=1.

Find the common mode function for VSI I and II when the function  =0 can be implemented by giving the parameter k, 0 <k <1, as follows:

$$v_{comI} = (1-k).e_{comMinI} + k.e_{comMaxI}$$
$$v_{comII} = (1-k).e_{comMinII} + k.e_{comMaxII}$$

(55)

The value of k is determined by the condition $v_{comOpt} = 0$:

$$k = - \frac{e_{comMin1} + e_{comMin2}}{e_{comMax1} + e_{comMax2} - e_{comMin1} - e_{comMin2}}$$

(56)

Or short form:

$$k = - \frac{e_{comMin}}{e_{comMax} - e_{comMin}}$$

(57)

The correlation between the maximum and minimum common mode voltages with common mode voltage VcomOpt is shown in Fig. 13 and Fig. 14.

The 6P_VSI uses the RCMV PWM technique using VcomOpt to reach the minimum average common mode value. When the modulation index m <0.866, VcomOpt = 0. It can be said that the RCMV PWM technique with VcomOpt acts as the SIN PWM technique extends to the maximum voltage value of the VSI, where the modulation index is equal to 1

*3) RCMV POD-CBPWM technique:* The POD_CBPWM technique for 6P_VSI can be defined as a conventional SIN PD-CBPWM technique. The difference is that the CBPWM block in Fig. 4, when applied to VSI II, will use a 180 degree phase shift carrier compared to the carrier used for VSI I CBPWM block. As described in Fig. 4, the Vcom function is determined by zero (Vcom = 0).

## IV. SIMULATION RESULTS AND DISCUSSION

Simulation analysis has been carried out by MATLAB/ Simulink to verify the proposed methods. Comparisons with the conventional SIN CBPW technique are also presented in this section.

*1) The conventional SIN CBPW technique*: SIN PD CBPWM technique uses the same phase carrier for both 3P_VSI I and II. The results of Fig. 15 and 16 show that the Vcom common mode voltage changes to the highest peak values $\pm$ Vd/2. The limited control range of the SIN CBPWM technique is the largest modulation index m = 0.866.



Fig 15.  SIN PD PWM technique, m=0.8- Diagram of the line voltage and phase voltage.



Fig 16.  SIN PD PWM technique, m=0.8- Diagram of  common mode voltage $v_{comI}$, $v_{comII}$ and $v_{com}$ of 6P_VSI.

*2) SIN POD CBPWM technique*: SIN POD CBPWM technique utilizes the carriers for VSI II is shifted phase 180 degrees to the carriers for VSI I. The results show that, in Fig. 17 and 18, the SIN POD CBPWM instantaneous-mode voltage value decreases within the limits $V_{com(peak)=\pm}V_d/6$, although the common mode voltage components have peak values $V_{com(peak)= \pm}V_d/2$. It is easy to see that, because the voltage of each 3P_VSI is controlled  PWM independently, the output voltage quality of the SIN POD CBPWM technique is also equal to the harmonic quality of the SIN PD CBPWM technique. The SIN POD CBPWM technique has a control range limited to the largest modulation index m = 0.866.

Fig 17.   SIN POD PWM technique, m=0.8- Diagram of the line voltage and phase voltage, m=0.8.



Fig 18.   SIN POD PWM technique, m=0.8- Diagram of  common mode voltage $v_{comI}$ , $v_{comII}$  and  $v_{com}$ of 6P_VSI.

To extend RCMV PWM to m = 1, it is possible to use (45) to create the VcomOptI, VcomOptII optimizer for each of VSI I, II.  Vcom voltage of 6P_VSI (Fig. 13a) will be larger than the coordinate control case (Fig. 13b).

*3) RCMV 4S-PWM technique with VcomMid*: The results of Fig. 19 and 20 show that the instantaneous common-mode voltage value of the RCMV 4S-CBPWM technique decreased within the limits $V_{com(peak)=\pm}V_d/6$. Unlike the RCMV POD CBPWM, this method simultaneously controls the common mode voltage components within the above limits  $V_{comI,II}$ $_{(peak)=\pm}V_d/6$,. The voltage of each 3P_VSI is controlled PWM independently to m = 1.





Fig 19.   RCMV 4S-PWM technique with VcomMid, m=0.8- Diagram of the line voltage and phase voltage, m=0.8.



Fig 20.   RCMV 4S-PWM technique with VcomMid, m=0.8- Diagram of  common mode voltage $v_{comI}$ , $v_{comII}$  and  $v_{com}$ of 6P_VSI.

*4) 4S-PWM technique with VcomOpt*: The results in Fig. 21 and 22 show that the instantaneous common-mode voltage value of the RCMV 4S-PWM technique with VcomOpt decreases within the limits, while the common mode voltage component also decreases within the above limit,  $v_{comI,II}(peak)=\pm V_d / 6$ . The voltage of each 3-P_VSI is controlled independently of the PWM to the limit of m = 0.866. In range (m> 0.866) from the modulation index m= 0.866 to 1, the offset voltages of the two VSI I, II will be constrained under the extreme condition of the V$_{com}$ function.

From the principle of the RCMV methods performed and the load current diagram obtained through Simulink, it can be observed that the RCMV POD CBPWM technique will give better output quality than the other two methods.  The RCMV method 4S CBPWM has an operational range up to modulation index m = 1. In that, due to the symmetry time distribution of the two far vectors, the PWM method with V$_{comMid}$ will be of better quality. Evaluation of the average common mode voltage, the RCMV-4S CBPWM method will provide a minimum common mode voltage, which can help to limit the common mode amplitude when the drive system is connected to SPIM with a long cable.

Fig 21. RCMV 4S-PWM technique with $V_{comOpt}$, m=0.8- Diagram of the line voltage and phase voltage, m=0.8.



Fig 22. RCMV 4S-PWM technique with VcomOpt, m=0.8- Diagram of common mode voltage $V_{comI}$, $V_{comII}$ and $V_{com}$ of 6P_VSI.

## V. EXPERIMENT RESULTS

The characteristics of the proposal CBPWM methods were investigated experimentally; the experimental set-up is illustrated in Fig. 23. The proposal RCMV technologies, described in Section 4, is implemented on TMS320F28335 digital signal processor (DSP), which is used to control six-phase VSI of SPIM drive with the parameters of SPIM: 1HP, phase voltage 240 V, 50 Hz, 4 pole , 1450 rpm. Rs = 10.1Ω, Rr = 9.8546Ω, Ls = 0.833457 H, Lr = 0.830811 H, Lm = 0.783106H, Ji = 0.0088 kg.m2.

The theoretical analysis and simulation results have been verified with experiments carried out the SPIM drive with the same ratings, parameters, and operating conditions as those of the simulations. As discussed and simulink in part 3 and part 4,

the conventional Sin_PD PWM technique has been implemented to compare to three proposal methods: RCMV_POD_CBPWM technique, RCMV-4S CBPWM techniques with $V_{mid}$ and $_{Vopt}$.

Fig. 24 and 25 are the experimental results corresponding to the simulation waveforms of Fig. 15 and 16, respectively; similar, Fig. 26 and 27 are the experimental results corresponding to the simulation waveforms of Fig. 17 and 18, respectively; Fig. 28 and 29 are the experimental results corresponding to the simulation waveforms of Figs. 19 and 20, respectively, Fig. 30 and 31 are the experimental results corresponding to the simulation waveforms of Fig. 21 and 22, respectively. There is a strong correlation between the waveforms of experiments and simulations. In the experimental waveforms, slightly larger ripple than the simulation waveforms is observed.

The CMV comparison indicates that both RCMV-4S CBPWM techniques with $V_{mid}$ and $V_{opt}$ have low CMV components (Vcom I, VcomII) compared to other methods. With conventional SinPD-PWM and RCMV_POD_CBPWM technique, the peak component CMV are ±Vd/2, RCMV_4S Vmid_PWM, RCMV_4S_Vopt_PWM methods are $\pm$ Vd/6 reduce approximately 77% to the conventional methods. However, general CMV(Vcom) of the RCMV_POD_CBPWM technique is equal the general CMV(Vcom) of RCMV_4S_Vmid_PWM and RCMV_4S_V$_{Opt}$_PWM methods by ±Vd/6 reduce approximately 77% to the conventional methods.



Fig 23. Six-phase voltage source inverter.



Fig 24. SIN PD PWM technique, m=0.8- Diagram of the line voltage and phase voltage.

Fig 25.   SIN PD PWM technique, m=0.8- Diagram of  common mode voltage V$_{comI}$ , V$_{comII}$ and  V$_{com}$ of 6P_VSI.



Fig 26.   SIN POD PWM technique, m=0.8- Diagram of the line voltage and phase voltage.



Fig 27.   SIN POD PWM technique, m=0.8- Diagram of  common mode voltage  V$_{comI}$ , V$_{comII}$ and  V$_{com}$ of 6P_VSI.



Fig 28.   RCMV 4S-PWM technique with Vcommid, m=0.8- Diagram of the line voltage and phase voltage.



Fig 29.   RCMV 4S-PWM technique with Vcommid,  m=0.8- Diagram of common mode voltage $v_{comI}$ , $v_{comII}$ and $v_{com}$ of 6P_VSI.



Fig 30.   RCMV 4S-PWM technique with Vcommid, m=0.8- Diagram of the line voltage and phase voltage.

Fig 31. RCMV 4S-PWM technique with VcomOpt, m=0.8- Diagram of common mode voltage $V_{comI}$, $V_{comII}$ and $V_{com}$ of 6P_VSI.

## VI. CONCLUSION

This paper presents a novel CBPWM technique to reduce CMV for SPIM VSI. These proposal schemes are implemented simply, require low calculation effort. The average CMV control range is also explained. CBPWM techniques reduce common mode voltage in Vd/6 range. The RCMV POD CBPWM method is very simply implemented. This proposal helps to reduce the common mode voltage and achieves high output voltage quality. However, this scheme has the range of modulation index m only up to m = 0866. RCMV 4S-CBPWM techniques are capable to perform in range entire of the hexagonal voltage vector, mean this proposal schema has modulation index m up to m=1. The calculations show that the responses of the CBPWM technique using VcomMid function produces better output harmonics; while the method using VcomOpt function produce the smallest triple harmonic component or by zero. Because of using the two further vectors, the 4S-CBPWM methods have the responses of output harmonics not good by using the CBPWM POD technique. The theoretical analysis and simulation results have been verified with experiments carry out SPIM drive with the same ratings, parameters, and operating conditions as those of the simulations. In the experiments, the conventional Sin_PD PWM has been implemented. As discussed in Section 3 to compare with three proposal methods Sin_POD PWM, RCMV_4SVmid_PWM, RCMV_4SVmid_PWM technique. The CMV comparison indicates that both RCMV_4S_Vmid_PWM and RCMV_4S_VOpt_PWM technique have low CMV components (Vcom I, VcomII) compared to other methods. With Sin PD-PWM and, the peak component CMV are Vd/2, RCMV_4S_Vmid_PWM,

RCMV_4S_Vmid_PWM methods are Vd/6 reduce approximately 77% to the conventional methods. However, general CMV(Vcom) of Sin POD-PWM method are equal the general CMV(Vcom) of RCMV_4S_Vmid_PWM and RCMV_4S_Vmid_PWM methods by Vd/6 reduce approximately 77% to the conventional methods. As analysed, the use of the six phase induction motor drives are recommended for the high power applications, where the use of multi_level 6P_VSI is considered replacing the two level 6P_VSI in some cases require high voltage. Therefore, it will be interesting if an investigation is carried out to study continue to develop and improve the proposed techniques in this paper to reduce common mode voltage for the six phase drives using the multi_level 6P_VSI in future.

### REFERENCES

[1] M.B.R. Corra, C.B. Jacobina, C.R. da Silva, A.N. Lima, E. R. C. da Silva,"Six-phase AC drive system with reduced common-mode voltage", IEEE International Electric Machines and Drives Conference, pp.1852-1858, 15 July 2003.

[2] Drazen Dujic, Atif Iqbal, Emil Levi "A Space Vector PWM Technique for Symmetrical Six-Phase Voltage Source Inverters", EPE Journal, vol. 17, no. 1, pp-24-32, March 2007.

[3] R. Kianinezhad, B. Nahid-Mobarakeh, L. Baghli, F. Betin, G.A. Capolino, Modeling and Control of Six-Phase Symmetrical Induction Machine Under Fault Condition Due to Open Phases, IEEE Trans. Ind. Elec., vol. 55, no. 5, pp. 1966-1977, May 2008.

[4] R. Kianinezhad1, R. Alcharea2, B. Nahid3, F. Betin2, Analysis and Evaluation of DTC and FOC in Open Phase Fault Operation of Six-Phase Induction Machines G.-A. Capolino2, 978-1-4244-1633-2/08/.00 ©2008 IEEE

[5] R. Bojoi, A. Tenconi, F. Profumo, G. Griva, D. Martinello, "Complete Analysis and comparative study of digital modulation techniques for dual three phase AC motor Drives", pp. 851-857, IEEE PESC 2002.

[6] F. Wang, "Motor shaft voltages and bearing currents and their reduction in multilevel medium-voltage PWM voltage-source-inverter drive applications," IEEE Trans. Ind. Appl., vol. 36, no. 5, pp. 1336 -1341, Sep./Oct. 2000.

[7] R. Naik, T. A. Nondahl, M. Cacciato, A. Consoli, G. Scarcella, and A. Testa, "Reduction of common mode currents in PWM inverter motor drives" IEEE Trans. Ind. Appl., vol. 35, no. 2, pp. 469–476, Mar./Apr. 1999.

[8] G.Oriti, A.L. Julian, T. Lipo, An Inverter/Motor Drive with Common Mode Voltage Elimination, IEEE IAS Proceedings 1997.

[9] J. Huang and H. Shi, "Reducing the common-mode voltage through carrier peak position modulation in an SPWM three-phase inverter,"IEEE Trans. Power Electron., vol. 29, no. 9, pp. 4490-4495, Sep. 2014.

[10] M. H. Hedayati, A. B. Acharya, and V. John, "Common-mode filter design for PWM rectifier-based motor drives," IEEE Trans. Power Electron., vol. 28, no. 11, pp. 5364 – 537, Nov. 2013.

[11] H. Akagi and S. Tamura, "A passive EMI filter for eliminating both bearing current and ground leakage current from an inverter-driven motor," IEEE Trans. Power Electron., vol. 5, no. 5, pp. 1459–1469, Sep. 2006.

[12] Kai Tian, Jiacheng Wang, Bin Wu, Dewei Xu, Zhongyuan Cheng, Navid Reza Zargari, "A Virtual Space Vector Modulation Technique for the Reduction of Common -Mode Voltages in both Magnitude and Third-Order Component", IEEE Transactions on Power Electronics, Vol. 31, No.1, Jan. 2016.

[13] E. Un and A. M. Hava, "A near-state PWM method with reduced switching losses and reduced common-mode voltage for three-phase voltage source inverters," IEEE Trans. Ind. Appl., vol. 45, no. 2, pp. 782-793, Mar./Apr. 2009.

[14] A. M. Hava and E. Un, "A high-performance PWM algorithm for common mode voltage reduction in three-phase voltage source inverters," IEEE Trans. Power Electron., vol. 26, no. 7, pp.1998-2008, Jul. 2011.

[15] M. Pulirenti, G. Scarcella, G. Scelba, and M. Cacciato, "Space Vector Modulation technique for Common Mode Currents reduction in six phase AC drives" DIEES-UNIVERSITY OF CATANIA Viale A. Doria, 6 95125, Catania, Italy

[16] Rutian Wang, Xingjun Mu, Zhiqiang Wu, Lihui Zhu, Qiufeng Chen and Xue Wang; "Carrier-Based PWM Method to Reduce Common-Mode Voltage of Three-to-Five-Phase Indirect Matrix Converter" Hindawi Publishing Corporation Mathematical Problems in Engineering, Vol. 2016, Article ID 6086497, 10 pages.

**Ngoc Thuy Pham** was born in Viet Nam, in 1976. She received the B.Sc degrees in Electrical Engineering from Thai Nguyen University of Technology (TNUT) in 1994, and the M.Sc from Ho Chi Minh City University of Technology ( HCMUT) in 2009. She worked from 2000 in the Faculty of Electrical Engineering, Industrial University of Ho Chi Minh City (IUH). Her current research interests include AC motor drives, active power filters, and PWM techniques for power converters. multiphase induction motor, sensorless control of multiphase induction motor drives.



**Nho Van Nguyen** was born in Vietnam in 1964. He received his M.S. and Ph.D. degrees in Electrical Engineering from the University of West Bohemia, the Czech Republic in 1988 and 1991 respectively. Since 1992, he has been with the Department of Electrical and Electronics Engineering, Ho Chi Minh City University of Technology, Vietnam, where he is currently an Associate Professor. He was with KAIST as a Post-doc Fellow for six months in 2001 and a Visiting Professor for a year in 2003–2004. He was a visiting scholar at the Department of Electrical Engineering, University of Illinois at Urbana-Champaign for a month in 2009. His research interests include modeling and control of switching power supplies, AC motor drives, active power filters, and PWM techniques for power converters. He is a member of the Institute of Electrical and Electronics Engineers (IEEE).

# Performance Analysis of Machine Learning Techniques on Software Defect Prediction using NASA Datasets

Ahmed Iqbal[1], Shabib Aftab[2], Umair Ali[3], Zahid Nawaz[4], Laraib Sana[5], Munir Ahmad[6], Arif Husen[7]

Department of Computer Science, Virtual University of Pakistan, Lahore, Pakistan[1, 2, 3, 4, 6, 7]
Department of Computer Science, Lahore College for Women University, Lahore, Pakistan[5]
Department of Computer Science, COMSATS University Islamabad, Lahore Campus, Pakistan[7]

*Abstract*—Defect prediction at early stages of software development life cycle is a crucial activity of quality assurance process and has been broadly studied in the last two decades. The early prediction of defective modules in developing software can help the development team to utilize the available resources efficiently and effectively to deliver high quality software product in limited time. Until now, many researchers have developed defect prediction models by using machine learning and statistical techniques. Machine learning approach is an effective way to identify the defective modules, which works by extracting the hidden patterns among software attributes. In this study, several machine learning classification techniques are used to predict the software defects in twelve widely used NASA datasets. The classification techniques include: Naïve Bayes (NB), Multi-Layer Perceptron (MLP). Radial Basis Function (RBF), Support Vector Machine (SVM), K Nearest Neighbor (KNN), kStar (K*), One Rule (OneR), PART, Decision Tree (DT), and Random Forest (RF). Performance of used classification techniques is evaluated by using various measures such as: Precision, Recall, F-Measure, Accuracy, MCC, and ROC Area. The detailed results in this research can be used as a baseline for other researches so that any claim regarding the improvement in prediction through any new technique, model or framework can be compared and verified.

*Keywords*—*Software defect prediction; software metrics; data mining; machine learning; classification; class imbalance*

## I. INTRODUCTION

Prediction of defective modules in an early stage of software development is considered as one of most challenging aspect of quality assurance activity [11]. The identification of defects in an early stage is crucial as the cost of correcting these defects increases exponentially in the later phases of software development life cycle (SDLC). In software engineering, testing and bug fixing is very expensive and require huge amount of resources [12]. Predicting defective modules in the developing software has been investigated by many studies since the last two decades. An efficient software defect identification technique depends upon various factors, most importantly the extraction of software metrics from historical data. Various software metrics are used to classify the software instance/class/module as defective or non-defective. [13-16]. Furthermore, many empirical studies have also reflected that the subsets of software metrics can improve the performance of classifiers [17]. The activity of software

defect prediction is necessary in order to enhance the effectiveness of quality assurance process. It can help to develop a qualitative product with limited amount of resources in a limited time period. Machine learning techniques are considered a promising way to predict the software defects in an early stage of SDLC by detecting the hidden pattern in historical software data. The purpose of this paper is to analyze the performance of supervised machine learning techniques on software defect prediction by using NASA datasets. Machine learning techniques used in this research are: Naïve Bayes (NB), Multi-Layer Perceptron (MLP). Radial Basis Function (RBF), Support Vector Machine (SVM), K Nearest Neighbor (KNN), kStar (K*), One Rule (OneR), PART, Decision Tree (DT), and Random Forest (RF). Supervised machine learning techniques need the pre-classified data (training data) for training. During the training process these techniques make rules to classify the unseen data (test data) [18-19], [20-23], [26-27]. In this study, NASA's clean software defect datasets are used for experiments, including: CM1, JM1, KC1, KC3, MC1, MC2, MW1, PC1, PC2, PC3, PC4 and PC5. This research performs a detailed performance analysis of widely used machine learning classification techniques by using the 70:30 proportion of training and test data. The benchmark datasets are used in experiments so that any researcher can compare these results with the results of his/her proposed technique and claim the high accuracy which would be easy to validate for all research community.

Further organization of this paper is as follows. Section II discusses the related work. Section III describes about materials and methods used for the experiments. Section IV reflects the results and findings of the experiments. Section V finally concludes this study.

## II. RELATED WORK

Many researchers have used machine learning techniques to predict the software defects at an early stage of software development, some of the selected studies are discussed here. Researchers in [1] compared six classification techniques by using the data of 27 academic projects. Classification techniques include: Principal Component Analysis (PCA), Discriminant Analysis, Logistic Regression (LR), Holographic Networks, Logical Classification, and Layered Neural Networks model. Back propagation learning technique was used to develop Neural Network. Performance was evaluated

by using Predictive Validity, Verification Cost, Misclassification Rate, and Achieved Quality. According to results, no model performed well in predicting the software defects. Researchers in [2] used SVM for software defect prediction by using four publicly available NASA datasets: PC1, CM1, KC1 and KC3. The performance is compared with eight machine learning and statistical techniques i.e. K-Nearest Neighbours (KNN), Logistic Regression (LR), Multilayer Perceptron (MLP), Decision Trees, Radial Basis Function (RBF), Bayesian Belief Networks (BBN), Naïve Bayes, and Random Forest (RF). Parameters generated from confusion matrix were used for performance evaluation. The results reflected that SVM performed better than some of the other techniques. Researchers in [3] studied and explored the significant software metrics to predict the software defects. Significant metrics were identified through sensitive analysis by ANN model which was trained using the historical data. The identified metrics then used to develop separate Neural Network models to predict the defective modules. The performance was compared with the Gaussian kernel SVM. JM1 dataset was used for experiment from NASA MDP repository. The results reflected that SVM performed better than ANN in binary defect classification. In [4], researchers performed an experiment using three Cost-Sensitive Boosting algorithms and Back-Propagation learning techniques. From these three, two based on weight updating architectures and one based on threshold value. Four NASA datasets were used for experiment and performance was evaluated using Normalized Expected Cost of Misclassification (NECM). According to results, the threshold based Feed Forward Neural Network performed better than other methods particularly for object oriented software modules. Researchers in [5] compared the statistical and machine learning techniques on software defect prediction by using public domain datasets of AR1 and AR6. The techniques included: Artificial Neural Networks, Decision Trees, Cascade Correlation Network, Support Vector Machines, Group Method of Data Handling Method, and Gene Expression programming. Performance was evaluated by using AUC values. Results reflected that Decision Tree achieved 0.8 and 0.9 AUC scores for AR1 and AR6 respectively which were better than other used techniques. Researchers in [6] presented a software defect prediction technique using Conventional Radial Basis Function along with novel Adaptive Dimensional Biogeography-based optimization model. For experiment, five NASA datasets from PROMISE repository were used and the results showed the higher accuracy of proposed technique as compared to early used techniques. In [7], researchers developed a GUI tool in MATLAB for software defect prediction. The proposed tool was based on Bayesian Regularization (BR) technique which reduced the software cost by limiting the squared errors and weights. The performance of used technique was compared with Levenberg Marquardt (LM) Algorithm and according to results BR performed better. Researchers in [8] compared Artificial Neural Network (ANN) and Support Vector Machine (SVM) on software defect prediction. For experiment, seven NASA datasets from PROMISE repository were used. The performance was evaluated in terms of Specificity, Recall, and Accuracy. Results showed that SVM performed better. In [9], the researchers proposed a GUI tool in MATLAB which used CK

(Chidamber and Kemerer) object-oriented metrics for software defect prediction. For experiment, NASA datasets from PROMISE repository were used and performance of Levenberg-Marquardt (LM) algorithm is compared with Polynomial Function-based Neural Network on software defect prediction. According to results the proposed model performed better than other techniques.

## III. MATERIALS AND METHODS

This study analyzes the performance of various machine learning classifiers on software defect prediction by using NASA benchmark datasets. Each dataset includes several features along with known output class. The output/target class is one which is predicted on the basis of other available attributes. The attribute which is predicted is known as dependent attribute whereas other attributes which are used to predict the dependent attribute are known as independent attributes. The selected datasets for this study contains dependent attribute which has values either "Y" or "N". "Y" means the specific software instance or module has tendency to be defective and "N" means it is not defective. In this research, total of 12 cleaned NASA datasets [26] are used in experiments. The datasets includes CM1, JM1, KC1, KC3, MC1, MC2, MW1, PC1, PC2, PC3, PC4 and PC5 (Table I). Each selected dataset represents a NASA's software system, which includes different metrics, closely related to software quality.

Two versions of clean datasets are provided by [26]: DS' (which included duplicated and inconsistent instances) and DS'' (which does not include duplicated and inconsistent instances). These datasets were initially available at [27] but removed later. We have taken these datasets from [28], where backup of NASA datasets are stored. These cleaned datasets are already used and discussed by [29-31]. Table II reflects the cleaning criteria implemented by [26].

The experiments are performed in Weka [10], one of the most popular data mining tools. This tool is developed in Java language at the University of Waikato, New Zealand and widely accepted due to its portability, General Public License and ease of use.

TABLE I.    NASA CLEANED DATASETS [26]

| Dataset | Attributes | Modules | Defective | Non-Defective | Defective (%) |
|---------|-----------|---------|-----------|---------------|---------------|
| CM1 | 38 | 327 | 42 | 285 | 12.8 |
| JM1 | 22 | 7,720 | 1,612 | 6,108 | 20.8 |
| KC1 | 22 | 1,162 | 294 | 868 | 25.3 |
| KC3 | 40 | 194 | 36 | 158 | 18.5 |
| MC1 | 39 | 1952 | 36 | 1916 | 1.8 |
| MC2 | 40 | 124 | 44 | 80 | 35.4 |
| MW1 | 38 | 250 | 25 | 225 | 10 |
| PC1 | 38 | 679 | 55 | 624 | 8.1 |
| PC2 | 37 | 722 | 16 | 706 | 2.2 |
| PC3 | 38 | 1,053 | 130 | 923 | 12.3 |
| PC4 | 38 | 1,270 | 176 | 1094 | 13.8 |
| PC5 | 39 | 1694 | 458 | 1236 | 27.0 |

| Criterion | Data Quality Category | Explanation |
|---|---|---|
| 1. | Identical cases | 'Instances that have identical values for all metrics including class label'. |
| 2. | Inconsistent cases | 'Instances that satisfy all conditions of Case 1, but where class labels differ'. |
| 3. | Cases with missing values | 'Instances that contain one or more missing observations'. |
| 4. | Cases with conflicting feature values | 'Instances that have 2 or more metric values that violate some referential integrity constraint. For example, LOC TOTAL is less than Commented LOC. However, Commented LOC is a subset of LOC TOTAL'. |
| 5. | Cases with implausible values | 'Instances that violate some integrity constraint. For example, value of LOC=1.1' |

## IV. RESULTS AND DISCUSSION

This section aims to analyze the performance of used classification techniques. The performance is analyzed and evaluated through various measures generated from confusion matrix (shown in Fig. 1). A confusion matrix consists of the following parameters:

True Positive (TP): Instances which are actually positive and also classified as positive.

False Positive (FP): Instances which are actually negative but classified as positive.

False Negative (FN): Instances which are actually positive but classified as negative.

True Negative (TN): instances which are actually negative and also classified as negative.

The classification techniques are evaluated through following measures: Precision, Recall, F-measure, Accuracy, MCC and ROC.

Precision is defined as the ratio of True Positive (TP) modules with respect to total number of modules which are classified as positive [2].

$$Precision = \frac{TP}{(TP + FP)} \tag{1}$$



Fig. 1.   Confusion Matrix.

Recall is defined as the ratio of True Positive (TP) modules with respect to the total number of modules that are actually positive [2].

$$Re\,call = \frac{TP}{(TP + FN)} \tag{2}$$

F-measure provides the average of Precision & Recall [2].

$$F\text{-measure} = \frac{Precision * Recall * 2}{(Precision + Recall)} \tag{3}$$

Accuracy indicates that how much the prediction is accurate [2], [32].

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \tag{4}$$

Matthew's Correlation Coefficient (MCC) is defined as a ratio of the observed and predicted binary classifications and ranges from -1 to +1. The results closer to 1 depicts the good prediction whereas closer to or below 0 indicates the bad performance [24], [32].

$$MCC = \frac{TN * TP - FN * FP}{\sqrt{(FP + TP)(FN + TP)(TN + FP)(TN + FN)}} \tag{5}$$

The area under the ROC curve (AUC) is a measure of how well a parameter can distinguish between two classes (defective/non defective) [25], [31].

$$AUC = \frac{1 + TP_r - FP_r}{2} \tag{6}$$

All these performance measures are given by Weka tool. The results of Precision, Recall and F-Measure for each class (Y and N) are reflected in the tables (Table III to Table XIV). These accuracy measures are sensitive to class imbalance problem and reflect the symbol of '?' in case of such issue. Highest scores in each class are highlighted in bold for easy identification.

Results of CM1 datasets are given in Table III. It can be seen that in Precision, NB performed better in both the classes (Y and N). In Recall, NB and DT both performed better in Y class whereas RBF, SVM and PART showed better performance in N class and finally in F-measure, NB showed better performance in Y class whereas RBF, SVM and PART performed better in N class.

Results of JM1 datasets are reflected in Table IV. In precision, PART performed better in Y class whereas kStar performed better in N class. In Recall, kNN performed better in Y class and SVM performed better in N class. In F-measure, kStar outperformed in Y class whereas MLP and RBF outperformed in N class.

Table V reflects the results of KC1 dataset. It can be observed that in precision, SVM performed better in Y class whereas RF performed better in N class. In Recall, kNN and kStar both performed better in Y class and SVM performed better in N class. And finally in F-measure, RF performed better in Y class and RBF outperformed in N class.

TABLE III.    CM1 DATASET RESULTS

| Classifier | Class | Precision | Recall | F-Measure |
|---|---|---|---|---|
| NB | Y | **0.167** | **0.222** | **0.190** |
| | N | **0.919** | 0.888 | 0.903 |
| MLP | Y | 0.000 | 0.000 | 0.000 |
| | N | 0.904 | 0.955 | 0.929 |
| RBF | Y | ? | 0.000 | ? |
| | N | 0.908 | **1.000** | **0.952** |
| SVM | Y | ? | 0.000 | ? |
| | N | 0.908 | **1.000** | **0.952** |
| kNN | Y | 0.067 | 0.111 | 0.083 |
| | N | 0.904 | 0.843 | 0.872 |
| kStar | Y | 0.067 | 0.111 | 0.083 |
| | N | 0.904 | 0.843 | 0.872 |
| OneR | Y | 0.000 | 0.000 | 0.000 |
| | N | 0.903 | 0.944 | 0.923 |
| PART | Y | ? | 0.000 | ? |
| | N | 0.908 | **1.000** | **0.952** |
| DT | Y | 0.118 | **0.222** | 0.154 |
| | N | 0.914 | 0.831 | 0.871 |
| RF | Y | 0.000 | 0.000 | 0.000 |
| | N | 0.907 | 0.989 | 0.946 |

TABLE IV.    JM1 DATASET RESULTS

| Classifier | Class | Precision | Recall | F-Measure |
|---|---|---|---|---|
| NB | Y | 0.537 | 0.226 | 0.318 |
| | N | 0.823 | 0.949 | 0.882 |
| MLP | Y | 0.765 | 0.081 | 0.146 |
| | N | 0.804 | 0.993 | **0.889** |
| RBF | Y | 0.694 | 0.104 | 0.181 |
| | N | 0.807 | 0.988 | **0.889** |
| SVM | Y | ? | 0.000 | ? |
| | N | 0.792 | **1.000** | 0.884 |
| kNN | Y | 0.363 | **0.334** | 0.348 |
| | N | 0.829 | 0.846 | 0.837 |
| kStar | Y | 0.403 | 0.317 | **0.355** |
| | N | **0.830** | 0.876 | 0.853 |
| OneR | Y | 0.378 | 0.151 | 0.216 |
| | N | 0.807 | 0.935 | 0.866 |
| PART | Y | **0.818** | 0.019 | 0.037 |
| | N | 0.795 | 0.999 | 0.885 |
| DT | Y | 0.496 | 0.268 | 0.348 |
| | N | 0.828 | 0.929 | 0.876 |
| RF | Y | 0.572 | 0.189 | 0.284 |
| | N | 0.819 | 0.963 | 0.885 |

Results of KC3 dataset is reflected in Table VI. It is reflected that in Precision, MLP and OneR showed highest performance in Y class whereas NB performed better in N class. In Recall, NB and kNN performed better in Y class and in N class, SVM outperformed the others. In F-measure, NB performed better in Y class whereas SVM performed better in N class.

Results of MC1 dataset are reflected in Table VII. In Precision, kNN and PART showed better performance in Y class whereas NB performed better in N class. In Recall, NB performed better in Y class whereas MLP, RBF, SVM and DT performed better in N class. In F-Measure, kNN and PART performed better in Y class whereas MLP, RBF, SVM and DT performed better in N class.

TABLE V.    KC1 DATASET RESULTS

| Classifier | Class | Precision | Recall | F-Measure |
|---|---|---|---|---|
| NB | Y | 0.492 | 0.337 | 0.400 |
| | N | 0.795 | 0.881 | 0.836 |
| MLP | Y | 0.647 | 0.247 | 0.358 |
| | N | 0.787 | 0.954 | 0.863 |
| RBF | Y | 0.778 | 0.236 | 0.362 |
| | N | 0.789 | 0.977 | **0.873** |
| SVM | Y | **0.800** | 0.045 | 0.085 |
| | N | 0.753 | **0.996** | 0.858 |
| kNN | Y | 0.398 | **0.393** | 0.395 |
| | N | 0.793 | 0.796 | 0.795 |
| kStar | Y | 0.449 | **0.393** | 0.419 |
| | N | 0.801 | 0.835 | 0.817 |
| OneR | Y | 0.444 | 0.180 | 0.256 |
| | N | 0.767 | 0.923 | 0.838 |
| PART | Y | 0.667 | 0.157 | 0.255 |
| | N | 0.771 | 0.973 | 0.861 |
| DT | Y | 0.533 | 0.360 | 0.430 |
| | N | 0.803 | 0.892 | 0.845 |
| RF | Y | 0.615 | 0.360 | **0.454** |
| | N | **0.808** | 0.923 | 0.862 |

TABLE VI.    KC3 DATASET RESULTS

| Classifier | Class | Precision | Recall | F-Measure |
|---|---|---|---|---|
| NB | Y | 0.444 | **0.400** | **0.421** |
| | N | **0.878** | 0.896 | 0.887 |
| MLP | Y | **0.500** | 0.300 | 0.375 |
| | N | 0.865 | 0.938 | 0.900 |
| RBF | Y | 0.000 | 0.000 | 0.000 |
| | N | 0.818 | 0.938 | 0.874 |
| SVM | Y | ? | 0.000 | ? |
| | N | 0.828 | **1.000** | **0.906** |
| kNN | Y | 0.333 | **0.400** | 0.364 |
| | N | 0.870 | 0.833 | 0.851 |
| kStar | Y | 0.300 | 0.300 | 0.300 |
| | N | 0.854 | 0.854 | 0.854 |
| OneR | Y | **0.500** | 0.300 | 0.375 |
| | N | 0.865 | 0.938 | 0.900 |
| PART | Y | 0.250 | 0.100 | 0.143 |
| | N | 0.833 | 0.938 | 0.882 |
| DT | Y | 0.300 | 0.300 | 0.300 |
| | N | 0.854 | 0.854 | 0.854 |
| RF | Y | 0.286 | 0.200 | 0.235 |
| | N | 0.843 | 0.896 | 0.869 |

TABLE VII.    MC1 DATASET RESULTS

| Classifier | Class | Precision | Recall | F-Measure |
|---|---|---|---|---|
| NB | Y | 0.156 | **0.357** | 0.217 |
| | N | **0.984** | 0.953 | 0.968 |
| MLP | Y | ? | 0.000 | ? |
| | N | 0.976 | **1.000** | **0.988** |
| RBF | Y | ? | 0.000 | ? |
| | N | 0.976 | **1.000** | **0.988** |
| SVM | Y | ? | 0.000 | ? |
| | N | 0.976 | **1.000** | **0.988** |
| kNN | Y | **0.400** | 0.286 | **0.333** |
| | N | 0.983 | 0.990 | 0.986 |
| kStar | Y | 0.250 | 0.143 | 0.182 |
| | N | 0.979 | 0.990 | 0.984 |
| OneR | Y | 0.333 | 0.143 | 0.200 |
| | N | 0.979 | 0.993 | 0.986 |
| PART | Y | **0.400** | 0.286 | **0.333** |
| | N | 0.983 | 0.990 | 0.986 |
| DT | Y | ? | 0.000 | ? |
| | N | 0.976 | **1.000** | **0.988** |
| RF | Y | 0.000 | 0.000 | 0.000 |
| | N | 0.976 | 0.998 | 0.987 |

Results of MC2 dataset are reflected in Table VIII. In precision, NB performed better in Y class whereas PART performed better in N class. In Recall, PART performed better in Y class whereas NB and RBF performed better in N class. In F Measure, PART performed better in both the classes.

TABLE VIII.    MC2 DATASET RESULTS

| Classifier | Class | Precision | Recall | F-Measure |
|---|---|---|---|---|
| NB | Y | **0.833** | 0.385 | 0.526 |
| | N | 0.742 | **0.958** | 0.836 |
| MLP | Y | 0.500 | 0.538 | 0.519 |
| | N | 0.739 | 0.708 | 0.723 |
| RBF | Y | 0.800 | 0.308 | 0.444 |
| | N | 0.719 | **0.958** | 0.821 |
| SVM | Y | 0.400 | 0.154 | 0.222 |
| | N | 0.656 | 0.875 | 0.750 |
| kNN | Y | 0.667 | 0.462 | 0.545 |
| | N | 0.750 | 0.875 | 0.808 |
| kStar | Y | 0.400 | 0.308 | 0.348 |
| | N | 0.667 | 0.750 | 0.706 |
| OneR | Y | 0.500 | 0.231 | 0.316 |
| | N | 0.677 | 0.875 | 0.764 |
| PART | Y | 0.727 | **0.615** | **0.667** |
| | N | **0.808** | 0.875 | **0.840** |
| DT | Y | 0.500 | 0.385 | 0.435 |
| | N | 0.704 | 0.792 | 0.745 |
| RF | Y | 0.500 | 0.462 | 0.480 |
| | N | 0.720 | 0.750 | 0.735 |

Table IX reflects the result of MW1 dataset. It can be seen that in Precision, MLP performed better in both the classes. In Recall, MLP performed better in Y class whereas OneR performed better in in N class. In F-measure, MLP performed better in both the classes.

TABLE IX.    MW1 DATASET RESULTS

| Classifier | Class | Precision | Recall | F-Measure |
|---|---|---|---|---|
| NB | Y | 0.333 | 0.625 | 0.435 |
| | N | 0.95 | 0.851 | 0.898 |
| MLP | Y | **0.545** | **0.75** | **0.632** |
| | N | **0.969** | 0.925 | **0.947** |
| RBF | Y | ? | 0.00 | ? |
| | N | 0.893 | 1.000 | 0.944 |
| SVM | Y | ? | 0.000 | ? |
| | N | 0.893 | 1.000 | 0.944 |
| kNN | Y | 0.400 | 0.500 | 0.444 |
| | N | 0.938 | 0.910 | 0.924 |
| kStar | Y | 0.143 | 0.125 | 0.133 |
| | N | 0.897 | 0.910 | 0.904 |
| OneR | Y | 0.500 | 0.125 | 0.200 |
| | N | 0.904 | **0.985** | 0.943 |
| PART | Y | 0.250 | 0.125 | 0.167 |
| | N | 0.901 | 0.955 | 0.928 |
| DT | Y | 0.250 | 0.125 | 0.167 |
| | N | 0.901 | 0.955 | 0.928 |
| RF | Y | 0.333 | 0.125 | 0.182 |
| | N | 0.903 | 0.970 | 0.935 |

TABLE X.    PC1 DATASET RESULTS

| Classifier | Class | Precision | Recall | F-Measure |
|---|---|---|---|---|
| NB | Y | 0.280 | **0.700** | 0.400 |
| | N | 0.983 | 0.907 | 0.944 |
| MLP | Y | **1.000** | 0.300 | 0.462 |
| | N | 0.965 | **1.000** | **0.982** |
| RBF | Y | 0.333 | 0.100 | 0.154 |
| | N | 0.955 | 0.990 | 0.972 |
| SVM | Y | ? | 0.000 | ? |
| | N | 0.951 | **1.000** | 0.975 |
| kNN | Y | 0.273 | 0.300 | 0.286 |
| | N | 0.964 | 0.959 | 0.961 |
| kStar | Y | 0.125 | 0.300 | 0.176 |
| | N | 0.961 | 0.892 | 0.925 |
| OneR | Y | 0.333 | 0.100 | 0.154 |
| | N | 0.955 | 0.990 | 0.972 |
| PART | Y | 0.375 | 0.600 | 0.462 |
| | N | 0.979 | 0.948 | 0.963 |
| DT | Y | 0.389 | **0.700** | **0.500** |
| | N | **0.984** | 0.943 | 0.963 |
| RF | Y | 0.750 | 0.300 | 0.429 |
| | N | 0.965 | 0.995 | 0.980 |

Results of PCI datasets are shown in Table X. It can be seen that in Precision, MLP performed better in Y class whereas DT performed better in N class. In Recall, NB and DT performed better in Y class whereas MLP and SVM both performed better in N class. In F-measure, DT performed better in Y class whereas MLP performed better in N class.

Results of PC2 datasets are shown in Table XI. According to results in Precision, kStar performed well in both the classes. In Recall, kStar performed well in Y class whereas RBF, SVM, DT and RF performed well in N class. In F-measure, kStar performed well in Y class however RBF, SVM, DT and RF performed well in N class.

Results of PC3 dataset is reflected in Table XII. It can be seen that in Precision, OneR and RF performed better in Y class however NB performed better in N class. In Recall, NB performed better in Y class whereas RBF, SVM and PART performed better in N class. In F-measure, DT performed better in Y class whereas OneR and RF performed better in N class.

Results of PC4 dataset are shown in Table XIII. It is reflected that in Precision, SVM performed better in Y class whereas DT performed better in N class. In Recall, DT performed better in Y class whereas SVM performed better in N class. In F-Measure, DT performed better in Y class whereas RF performed better in N class.

Results of PC5 dataset are shown in Table XIV. It can be seen that in Precision, SVM performed better in Y class whereas DT performed better in N class. In Recall, DT performed better in Y class whereas SVM performed better in N Class. In F Measure, DT performed better in Y class whereas RBF performed better in N class.

TABLE XI. PC2 DATASET RESULTS

| Classifier | Class | Precision | Recall | F-Measure |
|---|---|---|---|---|
| NB | Y | 0.000 | 0.000 | 0.000 |
| | N | 0.976 | 0.967 | 0.972 |
| MLP | Y | 0.000 | 0.000 | 0.000 |
| | N | 0.977 | 0.991 | 0.984 |
| RBF | Y | ? | 0.000 | ? |
| | N | 0.977 | **1.000** | **0.988** |
| SVM | Y | ? | 0.000 | ? |
| | N | 0.977 | **1.000** | **0.988** |
| kNN | Y | 0.000 | 0.000 | 0.000 |
| | N | 0.977 | 0.991 | 0.984 |
| kStar | Y | **0.143** | **0.200** | **0.167** |
| | N | **0.981** | 0.972 | 0.976 |
| OneR | Y | 0.000 | 0.000 | 0.000 |
| | N | 0.977 | 0.995 | 0.986 |
| PART | Y | 0.000 | 0.000 | 0.000 |
| | N | 0.977 | 0.991 | 0.984 |
| DT | Y | ? | 0.000 | ? |
| | N | 0.977 | **1.000** | **0.988** |
| RF | Y | ? | 0.000 | ? |
| | N | 0.977 | **1.000** | **0.988** |

TABLE XII. PC3 DATASET RESULTS

| Classifier | Class | Precision | Recall | F-Measure |
|---|---|---|---|---|
| NB | Y | 0.150 | **0.907** | 0.257 |
| | N | **0.929** | 0.190 | 0.316 |
| MLP | Y | 0.346 | 0.209 | 0.261 |
| | N | 0.883 | 0.938 | 0.909 |
| RBF | Y | ? | 0.000 | ? |
| | N | 0.864 | **1.000** | 0.927 |
| SVM | Y | ? | 0.000 | ? |
| | N | 0.864 | **1.000** | 0.927 |
| kNN | Y | 0.480 | 0.279 | 0.353 |
| | N | 0.893 | 0.952 | 0.922 |
| kStar | Y | 0.313 | 0.233 | 0.267 |
| | N | 0.884 | 0.919 | 0.901 |
| OneR | Y | **0.600** | 0.140 | 0.226 |
| | N | 0.879 | 0.985 | **0.929** |
| PART | Y | ? | 0.000 | ? |
| | N | 0.864 | **1.000** | 0.927 |
| DT | Y | 0.500 | 0.279 | **0.358** |
| | N | 0.894 | 0.956 | 0.924 |
| RF | Y | **0.6000** | 0.140 | 0.226 |
| | N | 0.879 | 0.985 | **0.929** |

TABLE XIII. PC4 DATASET RESULTS

| Classifier | Class | Precision | Recall | F-Measure |
|---|---|---|---|---|
| NB | Y | 0.486 | 0.346 | 0.404 |
| | N | 0.901 | 0.942 | 0.921 |
| MLP | Y | 0.676 | 0.481 | 0.562 |
| | N | 0.922 | 0.964 | 0.942 |
| RBF | Y | 0.667 | 0.154 | 0.250 |
| | N | 0.881 | 0.988 | 0.931 |
| SVM | Y | **0.818** | 0.173 | 0.286 |
| | N | 0.884 | **0.994** | 0.936 |
| kNN | Y | 0.477 | 0.404 | 0.438 |
| | N | 0.908 | 0.930 | 0.919 |
| kStar | Y | 0.333 | 0.327 | 0.330 |
| | N | 0.894 | 0.897 | 0.895 |
| OneR | Y | 0.650 | 0.250 | 0.361 |
| | N | 0.892 | 0.979 | 0.933 |
| PART | Y | 0.464 | 0.500 | 0.481 |
| | N | 0.920 | 0.909 | 0.914 |
| DT | Y | 0.515 | **0.673** | **0.583** |
| | N | **0.946** | 0.900 | 0.922 |
| RF | Y | 0.778 | 0.404 | 0.532 |
| | N | 0.912 | 0.982 | **0.946** |

Accuracy results are shown in Table XV. It can be seen that RBF showed better performance with higher accuracy in 5 datasets. MLP, SVM, RF each showed higher accuracy in 4

datasets. On the other hand NB, kStar and kNN did not show higher accuracy in any of the dataset. MCC results are shown in Table XVI, it can be seen that scores in most of the classifiers could not be drawn due to class imbalance. NB and DT each performed better in 3 datasets whereas MLP showed higher performance in 2 datasets. RBF, kNN, kStar, PART and RF each showed higher performance in 1 dataset. SVM and OneR did not perform well in any single dataset. Table XVII shows the results of ROC area, it can be seen that RF reflected high performance in 8 datasets whereas NB, MLP, kStar and PART each showed high performance in 1 dataset. Remaining algorithms did not show high performance in any of the used dataset. It has been noted that Besides the Precision, Recall and F measure, MCC is also sensitive to the class imbalance problem as it could not give the scores in many of the classifiers. However accuracy and ROC both did not show any symbol of class imbalance in the results which make them non sensitive to this issue. Therefore, besides the Accuracy and ROC, other performance measures including Precision, Recall, F-Measure and MCC should be used for effective performance analysis. Class imbalance issue can be resolved in the used datasets with many techniques [31]. However the purpose of this research is to use the exact snapshot of NASA cleaned dataset that is why no preprocessing or class balancing techniques are used.

TABLE XIV. PC5 DATASET RESULTS

| Classifier | Class | Precision | Recall | F-Measure |
|---|---|---|---|---|
| NB | Y | 0.676 | 0.168 | 0.269 |
| | N | 0.759 | 0.970 | 0.852 |
| MLP | Y | 0.560 | 0.204 | 0.299 |
| | N | 0.762 | 0.941 | 0.842 |
| RBF | Y | 0.760 | 0.139 | 0.235 |
| | N | 0.756 | 0.984 | **0.855** |
| SVM | Y | **0.875** | 0.051 | 0.097 |
| | N | 0.740 | **0.997** | 0.850 |
| kNN | Y | 0.500 | 0.496 | 0.498 |
| | N | 0.815 | 0.817 | 0.816 |
| kStar | Y | 0.439 | 0.423 | 0.431 |
| | N | 0.790 | 0.801 | 0.795 |
| OneR | Y | 0.455 | 0.336 | 0.387 |
| | N | 0.776 | 0.852 | 0.812 |
| PART | Y | 0.646 | 0.226 | 0.335 |
| | N | 0.770 | 0.954 | 0.852 |
| DT | Y | 0.537 | **0.526** | **0.531** |
| | N | **0.826** | 0.833 | 0.830 |
| RF | Y | 0.588 | 0.365 | 0.450 |
| | N | 0.794 | 0.906 | 0.846 |

TABLE XV. ACCURACY RESULTS

| Dataset | NB | MLP | RBF | SVM | kNN | kStar | OneR | PART | DT | RF |
|---|---|---|---|---|---|---|---|---|---|---|
| CM1 | 82.6531 | 86.7347 | **90.8163** | **90.8163** | 77.551 | 77.551 | 85.7143 | **90.8163** | 77.551 | 89.7959 |
| JM1 | 79.8359 | 80.3541 | **80.3972** | 79.1883 | 73.9637 | 75.9931 | 77.1589 | 79.4905 | 79.1019 | 80.1813 |
| KC1 | 74.212 | 77.3639 | **78.7966** | 75.3582 | 69.341 | 72.2063 | 73.3524 | 76.5043 | 75.6447 | 77.937 |
| KC3 | 81.0345 | **82.7586** | 77.5862 | **82.7586** | 75.8621 | 75.8621 | **82.7586** | 79.3103 | 75.8621 | 77.5862 |
| MC1 | 93.8567 | **97.6109** | **97.6109** | **97.6109** | 97.2696 | 96.9283 | 97.2696 | 97.2696 | **97.6109** | 97.4403 |
| MC2 | 75.6757 | 64.8649 | 72.973 | 62.1622 | 72.973 | 59.4595 | 64.8649 | **78.3784** | 64.8649 | 64.8649 |
| MW1 | 82.6667 | **90.6667** | 89.3333 | 89.3333 | 86.6667 | 82.6667 | 89.3333 | 86.6667 | 86.6667 | 88.000 |
| PC1 | 89.7059 | **96.5686** | 94.6078 | 95.098 | 92.6471 | 86.2745 | 94.6078 | 93.1373 | 93.1373 | 96.0784 |
| PC2 | 94.47 | 96.7742 | **97.6959** | **97.6959** | 96.7742 | 95.3917 | 97.235 | 96.7742 | **97.6959** | **97.6959** |
| PC3 | 28.7975 | 83.8608 | 86.3924 | 86.39 24 | 86.0759 | 82.5949 | **87.0253** | 86.3924 | 86.3924 | **87.0253** |
| PC4 | 86.0892 | 89.7638 | 87.4016 | 88.189 | 85.8268 | 81.8898 | 87.9265 | 85.3018 | 86.8766 | **90.2887** |
| PC5 | 75.3937 | 74.2126 | 75.5906 | 74.2126 | 73.0315 | 69.8819 | 71.2598 | 75.7874 | 75.000 | **75.9843** |

TABLE XVI. MCC RESULTS

| Dataset | NB | MLP | RBF | SVM | kNN | kStar | OneR | PART | DT | RF |
|---|---|---|---|---|---|---|---|---|---|---|
| **CM1** | **0.097** | -0.066 | ? | ? | -0.037 | -0.037 | -0.074 | ? | 0.041 | -0.032 |
| **JM1** | 0.251 | 0.206 | 0.215 | ? | 0.186 | 0.212 | 0.126 | 0.104 | **0.252** | 0.244 |
| **KC1** | 0.250 | 0.296 | **0.347** | 0.151 | 0.190 | 0.238 | 0.147 | 0.239 | 0.291 | 0.346 |
| **KC3** | **0.309** | 0.295 | -0.107 | ? | 0.218 | 0.154 | 0.295 | 0.056 | 0.154 | 0.111 |
| **MC1** | 0.208 | ? | ? | ? | **0.325** | 0.174 | 0.206 | **0.325** | ? | -0.006 |
| **MC2** | **0.444** | 0.243 | 0.371 | 0.040 | 0.374 | 0.062 | 0.137 | 0.512 | 0.189 | 0.216 |
| **MW1** | 0.367 | **0.589** | ? | ? | 0.373 | 0.038 | 0.211 | 0.110 | 0.110 | 0.150 |
| **PC1** | 0.400 | **0.538** | 0.161 | ? | 0.247 | 0.128 | 0.161 | 0.440 | 0.490 | 0.459 |
| **PC2** | -0.028 | -0.015 | ? | ? | -0.015 | **0.146** | -0.010 | -0.015 | ? | ? |
| **PC3** | 0.088 | 0.183 | ? | ? | 0.294 | 0.173 | 0.245 | ? | **0.304** | 0.245 |
| **PC4** | 0.334 | 0.515 | 0.279 | 0.342 | 0.359 | 0.225 | 0.352 | 0.396 | 0.514 | **0.516** |
| **PC5** | 0.245 | 0.216 | 0.251 | 0.173 | 0.314 | 0.227 | 0.209 | 0.274 | **0.361** | 0.322 |

TABLE XVII. ROC AREA RESULTS

| Dataset | NB | MLP | RBF | SVM | kNN | kStar | OneR | PART | DT | RF |
|---------|-----|-----|-----|-----|-----|-------|------|------|-----|-----|
| **CM1** | 0.703 | 0.634 | 0.702 | 0.500 | 0.477 | 0.538 | 0.472 | 0.610 | 0.378 | **0.761** |
| **JM1** | 0.663 | 0.702 | 0.713 | 0.500 | 0.591 | 0.572 | 0.543 | 0.714 | 0.671 | **0.738** |
| **KC1** | 0.694 | 0.736 | 0.713 | 0.521 | 0.595 | 0.651 | 0.551 | 0.636 | 0.606 | **0.751** |
| **KC3** | 0.769 | 0.733 | 0.735 | 0.500 | 0.617 | 0.528 | 0.619 | 0.788 | 0.570 | **0.807** |
| **MC1** | 0.826 | 0.805 | 0.781 | 0.500 | 0.638 | 0.631 | 0.568 | 0.684 | 0.500 | **0.864** |
| **MC2** | **0.795** | 0.753 | 0.766 | 0.514 | 0.668 | 0.510 | 0.553 | 0.724 | 0.615 | 0.646 |
| **MW1** | 0.791 | **0.843** | 0.808 | 0.500 | 0.705 | 0.543 | 0.555 | 0.314 | 0.314 | 0.766 |
| **PC1** | 0.879 | 0.779 | 0.875 | 0.500 | 0.629 | 0.673 | 0.545 | **0.889** | 0.718 | 0.858 |
| **PC2** | 0.751 | 0.746 | 0.724 | 0.500 | 0.495 | **0.791** | 0.498 | 0.623 | 0.579 | 0.731 |
| **PC3** | 0.773 | 0.796 | 0.795 | 0.500 | 0.616 | 0.749 | 0.562 | 0.79 | 0.664 | **0.855** |
| **PC4** | 0.807 | 0.898 | 0.862 | 0.583 | 0.667 | 0.734 | 0.614 | 0.776 | 0.834 | **0.945** |
| **PC5** | 0.725 | 0.751 | 0.732 | 0.524 | 0.657 | 0.629 | 0.594 | 0.739 | 0.703 | **0.805** |

## V. CONCLUSION

Software defect prediction using machine learning techniques is considered as one of the emerging research areas now days. Identification of defects at the early stage of development can contribute to the delivery of high quality software by using limited amount of resources. This study deals with the detailed performance analysis of various machine learning classification techniques on software defect prediction using 12 widely used and publically available NASA datasets. The classification techniques include: Naïve Bayes (NB), Multi-Layer Perceptron (MLP). Radial Basis Function (RBF), Support Vector Machine (SVM), K Nearest Neighbor (KNN), kStar (K*), One Rule (OneR), PART, Decision Tree (DT), and Random Forest (RF). The performance is evaluated by using various measures extracted from confusion matrix such as: Precision, Recall, F-Measure, Accuracy, MCC, and ROC Area. It is reflected from the results that neither the Accuracy and nor the ROC can be used as an effective performance measure as both of these did not react on class imbalance issue. However, Precision, Recall, F-Measure and MCC reacted to class imbalance problem in the results with the symbol of '?'. The results presented in this research can be used as baseline for other researches so that the results of any proposed technique, model or framework can be compared and easily verified. For future work, class imbalance issue should be resolved in NASA cleaned datasets. Moreover, to improve the performance, feature selection and ensemble learning techniques should also be explored.

### REFERENCES

[1] F. Lanubile, A. Lonigro, and G. Vissagio, "Comparing models for identifying fault-prone software components." Seke, no. July, pp. 312–319, 1995.

[2] K. O. Elish and M. O. Elish, "Predicting defect-prone software modules using support vector machines," J. Syst. Softw., vol. 81, no. 5, pp. 649–660, 2008.

[3] I. Gondra, "Applying machine learning to software fault-proneness prediction," J. Syst. Softw., vol. 81, no. 2, pp. 186–195, 2008.

[4] J. Zheng, "Cost-sensitive boosting neural networks for software defect prediction," Expert Syst. Appl., vol. 37, no. 6, pp. 4537–4543, 2010.

[5] R. Malhotra, "Comparative analysis of statistical and machine learning methods for predicting faulty modules," Appl. Soft Comput. J., vol. 21, pp. 286–297, 2014.

[6] P. Kumudha and R. Venkatesan, "Cost-Sensitive Radial Basis Function Neural Network Classifier for Software Defect Prediction," Sci. World J., vol. 2016, 2016.

[7] R. Mahajan, S. K. Gupta, and R. K. Bedi, "Design of software fault prediction model using BR technique," in Procedia Computer Science, vol. 46, no. Icict 2014, pp. 849–858, 2015.

[8] I. A. and A. Saha, "Software Defect Prediction: A Comparison Between Artificial Neural Network and Support Vector Machine," Adv. Comput. Commun. Technol., pp. 51–61, 2017.

[9] M. Singh and D. Singh Salaria, "Software Defect Prediction Tool based on Neural Network," Int. J. Comput. Appl., vol. 70, no. 22, pp. 22–28, 2013.

[10] I. H. Witten, E. Frank, M. A. Hall, and C. J. Pal, Data Mining: Practical machine learning tools and techniques. Morgan Kaufmann, 2016.

[11] Y. Ma, G. Luo, X. Zeng, and A. Chen, "Transfer learning for cross-company software defect prediction," Inf. Softw. Technol., vol. 54, no. 3, pp. 248–256, 2012.

[12] Y. Singh and R. Malhotra, Object-Oriented Software Engineering. PHI Learning Pvt. Ltd. New Delhi, 2012.

[13] R. Moser, W. Pedrycz, and G. Succi, "A comparative analysis of the efficiency of change metrics and static code attributes for defect prediction," pp. 181, 2008.

[14] E. Giger, M. D'Ambros, M. Pinzger, and H. C. Gall, "Method-level bug prediction," pp. 171, 2012.

[15] S. E. S. Taba, F. Khomh, Y. Zou, A. E. Hassan, and M. Nagappan, "Predicting bugs using antipatterns," IEEE Int. Conf. Softw. Maintenance, ICSM, pp. 270–279, 2013.

[16] K. Herzig, S. Just, A. Rau, and A. Zeller, "Predicting defects using change genealogies," 2013 IEEE 24th Int. Symp. Softw. Reliab. Eng. ISSRE 2013, pp. 118–127, 2013.

[17] S. Moustafa, M. Y. ElNainay, N. El Makky, and M. S. Abougabal, "Software bug prediction using weighted majority voting techniques," Alexandria Eng. J., vol. 57, no. 4, pp. 2763–2774, 2018.

[18] M. Ahmad, S. Aftab, S. S. Muhammad, and S. Ahmad, "Machine Learning Techniques for Sentiment Analysis: A Review," Int. J. Multidiscip. Sci. Eng., vol. 8, no. 3, pp. 27-32, 2017.

[19] M. Ahmad, S. Aftab, I. Ali, and N. Hameed, "Hybrid Tools and Techniques for Sentiment Analysis: A Review," Int. J. Multidiscip. Sci. Eng., vol. 8, no. 3, pp. 29-33, 2017.

[20] M. Ahmad and S. Aftab, "Analyzing the Performance of SVM for Polarity Detection with Different Datasets," Int. J. Mod. Educ. Comput. Sci., vol. 9, no. 10, pp. 29–36, 2017.

[21] M. Ahmad, S. Aftab, and I. Ali, "Sentiment Analysis of Tweets using SVM," Int. J. Comput. Appl., vol. 177, no. 5, pp. 25–29, 2017.

[22] M. Ahmad, S. Aftab, M. S. Bashir, N. Hameed, I. Ali, and Z. Nawaz, "SVM Optimization for Sentiment Analysis," Int. J. Adv. Comput. Sci. Appl., vol. 9, no. 4, pp. 393-398, 2018.

[23] S. Aftab, M. Ahmad, N. Hameed, M. S. Bashir, I. Ali, and Z. Nawaz, "Rainfall Prediction in Lahore City using Data Mining Techniques," Int. J. Adv. Comput. Sci. Appl., vol. 9, no. 4, pp. 254-260, 2018.

[24] N. Farnaaz and M. A. Jabbar, "Random Forest Modeling for Network Intrusion Detection System," Procedia Comput. Sci., vol. 89, pp. 213–217, 2016.

[25] T. Fawcett, "An introduction to ROC analysis," Pattern Recognit. Lett., vol. 27, no. 8, pp. 861–874, 2006.

[26] M. Shepperd, Q. Song, Z. Sun, and C. Mair, "Data quality: Some comments on the NASA software defect datasets," IEEE Trans. Softw. Eng., vol. 39, no. 9, pp. 1208–1215, 2013.

[27] "NASA – Software Defect Datasets [Online]. Available: https://nasa softwaredefectdatasets.wikispaces.com. [Accessed: 01-April-2019].

[28] "NASA Defect Dataset." [Online]. Available: https://github.com/klainfo/NASADefectDataset. [Accessed: 01-April-2019].

[29] B. Ghotra, S. McIntosh, and A. E. Hassan, "Revisiting the impact of classification techniques on the performance of defect prediction models," Proc. - Int. Conf. Softw. Eng., vol. 1, pp. 789–800, 2015.

[30] G. Czibula, Z. Marian, and I. G. Czibula, "Software defect prediction using relational association rule mining," Inf. Sci. (Ny)., vol. 264, pp. 260–278, 2014.

[31] D. Rodriguez, I. Herraiz, R. Harrison, J. Dolado, and J. C. Riquelme, "Preliminary comparison of techniques for dealing with imbalance in software defect prediction," Proc. 18th Int. Conf. Eval. Assess. Softw. Eng. ACM, pp. 1–10, 2014.

[32] A. Iqbal and S. Aftab, "A Feed-Forward and Pattern Recognition ANN Model for Network Intrusion Detection," Int. J. Comput. Netw. Inf. Secur., vol. 11, no. 4, pp. 19–25, 2019.

# Experimental Evaluation of the Virtual Environment Efficiency for Distributed Software Development

Pavel Kolyasnikov[1]
Russian Academy of Education,
Russia

Evgeny Nikulchev[2], Iliy Silakov[3], Dmitry Ilin[4]
MIREA–Russian Technological University
Moscow, Russia

Alexander Gusev[5]
Kuban State University
Krasnodar, Russia

*Abstract*—**At every software design stage nowadays, there is an acute need to solve the problem of effective choice of libraries, development technologies, data exchange formats, virtual environment systems, characteristics of virtual machines. Due to the spread of various kinds of devices and the popularity of Web platforms, lots of systems are developed not for the universal installation on a device (box version), but for a specific architecture with the subsequent provision of web services. Under these conditions, the only way for estimating the efficiency parameters at the design stage is to conduct various kinds of experiments to evaluate the parameters of a particular solution. Using the example of the Web platform of digital psychological tools, the methods for experimental parameter evaluation were developed in the article. The mechanisms and technologies for improving the efficiency of the Vagrant and Docker cloud virtual environment were also proposed in the paper. A set of basic criteria for evaluating the effectiveness of the configuration of the virtual development environment has been determined to be rapid deployment; increase in the speed and decrease in the volume of resources used; increase in the speed of data exchange between the host machine and the virtual machine. The results of experimental estimates of the parameters that define the formulated efficiency criteria are given as: processor utilization involved (percentage); the amount of RAM involved (GB); initialization time of virtual machines (seconds); time to assemble the component completely (Build) and to reassemble the component (Watch) (seconds). To improve the efficiency, a file system access driver based on the NFS protocol was studied in the paper.**

*Keywords*—*Distributed software development; virtual development environment; increase development efficiency; virtual machines; vagrant; Docker; NFS; webpack*

## I. INTRODUCTION

Currently, the virtual cloud development environments are widely used in the team disturbed development of large projects [1]. This technology uses virtualization and virtual machines configuration management tools [2] to apply the necessary parameters and install the required components with the automation of the synchronization process, configuration and launch of the development environment.

The use of virtual development environments allows developers to avoid differences between the local development environments and the final platform, and greatly simplify the installation and configuration of environments on new machines.

The paper examined a system for the preparation of virtual development environments in Vagrant [3], which allows creating reproducible virtual development environments [4] and reduces the number of difficulties that can arise for the reason of the incompatibility of software and hardware used by developers [5].

The use of a development management system facilitates simultaneous distributed work on several components of the developed software [6, 7] and also automates the process of installing and configuring all the necessary components of the development environment [8-10]. The processes of updating and modifying the components used are also simplified [11], since it suffices only to make changes to the configuration files.

The paper contains the results of a study of increasing the efficiency in virtual software development of a digital psychological research platform [12]. At use of the virtual machines structure close to real servers, following performance problems were discovered [13]: low data exchange rate, as well as instability of work.

The aim of the paper is the development of research methodology and the development of mechanisms for increasing the efficiency of the virtual development environment in the distributed development of large software systems.

The paper consists of seven sections: the 1st is the Introduction, the 2nd is devoted to the formulation of the problem of describing the initial version of the development environment, the 3rd section contains the description of research methods and parameters for evaluating the effectiveness; the 4th section presents the results of experimental evaluation of the initial version of the development environment; the 5th section provides the development of alternative options aimed at improving efficiency; the 6th section contains the evaluation of the effectiveness of the alternatives; the 7th section shows the results.

The paper defines a set of criteria, estimation methods and suggests mechanisms for increasing the efficiency of the virtual work environment in distributed software development.

## II.   Source Data

The structure of virtual machines, close to the structure of servers involved in the project, is used as the original development environment (Fig. 1). With this approach, each component of the platform corresponds to a separate virtual machine configuration.

In the original development environment, the following virtual machine parameters were chosen for the components:

- API Server: OC ubuntu/xenial64, version 20180424.0.0, 1 CPU, 1024 MB RAM;

- Researcher Account: OC ubuntu/xenial64, version 20180424.0.0, 2 CPU, 1024 MB RAM;

- Psychotest Player: OC ubuntu/xenial64, version 20180424.0.0, 2 CPU, 512 MB RAM.

When developing using a similar structure of the working environment, there is a need to simultaneously run several virtual machines for simultaneous operation of several interconnected components. During the practical use of the developed structure of the virtual development environment on the computers of developers, a significant reduction in performance was noticed. Further observations revealed the following problems of the development environment:

- high workload on the developers' working machines;

- low data exchange rate of virtual machines with the main system;

- high component build time;

- work instability due to increasing loads;

- the need for manual execution of a large number of operations when starting the environment.

During the development of large projects, the number of components being developed increases. At use of such a structure, it can lead to an even greater decrease in computer performance and a decrease in the efficiency of developers [14]. In addition, the development also requires the use of additional software (development environment, web browser,

network data entry drivers [15]), which also leads to a significant increase in the workload of the developers' machines, a decrease in the stability of the system and an increase in software execution time involved in the development.

The low data exchange speed of the main system and the virtual machine, observed in using the developed structure of the virtual environment, may be due to two factors: a significant increase in the load on the developer's machine, but also the driver used by the virtual machine to access the parent file system.

It is also worth to highlight the problem of the need to perform a large number of repetitive actions manually in the development environment start. It also takes considerable time due to the need to wait for the end of each virtual machine before to start the rest of the components. Due to the increasing load on the developer's work machine, this time can also increase.

The problems described above have a significant impact on the speed of software development due to high performance losses and time consumption [16]. It is worth to note that with the expansion of the overall structure of the project, parallel development of several components becomes almost impossible due to an even greater increase in resource requirements.

Thus, it can be concluded that it is necessary to examine approaches to create development environments for software developers and improve their structure. Based on the problems found during the observations, the main requirements for the desired solution for the software development environment were identified:

- rapid deployment;

- increasing the speed and lowering the resources used;

- increasing the speed of data exchange between the host machine and the virtual machine.

It is necessary to evaluate alternative technologies and develop a solution that meets the requirements described above. In the study, it is also necessary to make an experimental assessment of the used and alternative solutions, and to evaluate the expediency of switching to an alternative structure of the development environment.



Fig. 1.   The Schema of the Original Development Environment.

## III. RESEARCH METHODS

For the experiments described in this paper, measurements of the following indicators were used:

- CPU resources usage (percent);

- the amount of RAM usage (GB);

- virtual machine startup time (seconds);

- component full build time (Build) (seconds);

- component rebuild time (Watch) (seconds).

All experiments were conducted on the same working machine with the following configuration:

- motherboard: Dell 00TMJ3;

- processor: Intel Core i5-5250U;

- RAM: DDR3 12 GB, frequency 800 MHz;

- drive: Samsung SSD 860 EVO 500 GB;

- operating system: Windows 10.

For each component of the virtual development environment, 10 experiments were conducted to measure each of the indicators. To estimate the deviation of the obtained values, Student's coefficient with confidential probability P=0.95 was used.

To measure the processor and RAM resources used, a bash script was developed that compared the figures before the virtual development environment was launched with the figures after the virtual machine started up completely. To get data about the required indicators, the system command "WMIC" was used in the script code. This command provides the possibility of getting the necessary data through a command interface.

All measurements of the indicators were carried out in the waiting state of the virtual machine – after its complete launch and during the work of key components (web server, database server, etc.).

To estimate the virtual machines launch time, the system utility "time" was used. This utility was launched with Vagrant as a prefix in the start up command ("time vagrant up") and provides information on the time spent for execution after completion. If it is necessary to launch several components of the development environment, launching the necessary Vagrant containers was carried out in parallel.

To estimate the execution time of build tasks (complete "build" assembly and reassembly in "watch" mode), the data provided after the task was completed with the tool for building web applications (Webpack) was used.

## IV. ASSESSMENT OF THE ORIGINAL DEVELOPMENT ENVIRONMENT

With the aim of estimating the performance and working speed of the used environment, an experiment was conducted, including measurements of the time taken to start up, as well as an assessment of the CPU resources and RAM used. The results of the experiment are presented in Table I.

According to the data obtained during the experiment, it is concluded about the objectivity of the problems described during the observations, which reinforces the need to find alternative solutions.

TABLE I.        RESULTS OF ASSESSING RESOURCE USAGE BY DEVELOPMENT ENVIRONMENTS

| Configuration | VM count | Startup time, sec | CPU usage, % | RAM usage, GB |
|---|---|---|---|---|
| Initial configuration (API Server) | 2 | 289,4 ± 2,8 | 24,3% ± 6,7% | 1,5 ± 0,1 |
| Initial configuration (Researcher Account) | 1 | 118,5 ± 13 | 21,8% ± 6% | 0,5 ± 0,08 |
| Initial configuration (Psychotest Player) | 1 | 155,3 ± 19,1 | 16,1% ± 3% | 0,7 ± 0,04 |
| Initial configuration (API Server, Researcher Account) | 3 | 347,4 ± 4,4 | 35,4% ± 9,1% | 2 ± 0,2 |
| Initial configuration (API Server, Psychotest Player) | 3 | 290,1 ± 4,5 | 29% ± 5,3% | 2 ± 0,06 |
| Initial configuration (3 components) | 4 | 404,8 ± 4,1 | 49% ± 5% | 2,2 ± 0,18 |

## V. DEVELOPMENT OF MECHANISMS FOR INCREASING THE EFFICIENCY OF THE ENVIRONMENT

Due to the high load arising from the use of virtual machines on the basis of Vagrant, it is advisable to consider alternative technologies for the organizing development environments. As an alternative, the Docker technology [17] was considered, which can also be used to organize synchronized development environments [18].

The Docker system is based on the use of abstraction with the help of the built-in Linux kernel virtualization capabilities for isolating various components used within the stand-alone containers with separate environments [19].

Originally, this approach was aimed at delivering the developed software solutions to the server, but later the platform also began to be used for replacing virtual development environments. At the same time, Docker can act as an independent solution [20], or work as the basis of a Vagrant-based solution, thereby replacing the use of virtual machines.

To estimate Docker as an alternative to Vagrant, a number of criteria were identified for comparing these technologies. The results of comparing are shown in Table II.

It can be concluded that both technologies provide comparable advantages for different approaches for similar tasks execution. Using Docker to synchronize and quickly set up the virtual development environment can potentially be a valid solution due to the advantages of containerization and the solutions used within the Docker implementation. However, this approach also obliges to develop a new structure of components using Docker containers.

TABLE II.        COMPARING VAGRANT AND DOCKER

| Criterion | Vagrant | Docker |
|---|---|---|
| License | MIT | Apache 2.0 |
| Developer | Hashi Corp. | Docker, Inc. |
| Platform supported | Linux, Windows, MacOS | Linux, Windows, MacOS |
| Purpose | Virtual development environment | Containerization of applications and automation of tasks, components delivery to the server |
| Ease of use | Starting the environment with the use of one command | An understanding of the container system and dependencies is necessary; launching the environment with one command in combination with Vagrant |
| Ease of configuration | Configuration file written on Ruby | Native format of configuration file |
| Container structure | The container includes all the dependencies specified in the configuration. The container is only the runtime environment | Each component and its dependencies are separate containers |

In the context of the developed platform, an important task was to maintain maximum proximity to the work systems environment on remote servers. In this case, changing the structure of components for use in the form of Docker containers inevitably causes a discrepancy between the developer's environment and the server system environment, and also makes it difficult to support an existing solution in the context of the developed platform. Therefore, the use of containerization technology becomes impractical.

Another possible solution in this case may save Vagrant as the basis of the chosen solution, but using a single virtual machine to run all the components.

As an alternative solution, it is worth to consider a configuration based on a single virtual environment in Vagrant. In this case, each component runs within one virtual machine, using also a single modular Vagrant configuration (Fig. 2).



Fig. 2.    Alternative Structure of Development Environment.

Use of such a structure should reduce the load on the computer by use of only one virtual machine. At the same time, the use of a modular configuration based on the original solution should simplify support and minimize discrepancies with the server environment.

To compare the original and alternative structures of the development environment, an experimental assessment of the resources consumed was carried out (Table III). As a virtual machine configuration, the alternative solution was based on ubuntu/xenial64 OS, version 20180424.0.0 with 2 CPUs, 2048 MB RAM.

With the aim of increasing the efficiency, it is also worth to examine alternative driver of data exchange with parent system based of NFS protocol [18], since such a solution can significantly increase the speed of working with the file system and the tasks execution [21].

An additional driver is necessary to use NFS on the machines with Windows OS management [22]. In addition, it is necessary to involve bindfs [23] extension that allows to transfer the access rights from parent system.

Due to the architectural feature of the NFS protocol, which does not provide an implementation of system signals in file changes tracking [24], it is also necessary to modify the configuration of the Webpack software [25] used to build web components. The following changes were made:

- to ensure compatibility with NFS, the "watchOptions.poll" option was installed, implementing a workaround of monitored files at a specified time interval;

- to exclude unchangeable library files from the build, the "watchOptions.ignore" option was used.

Additional measurements were taken for the time spent to build a component after configuration changes (Table IV).

TABLE III.        MEASUREMENTS RESULTS OF THE DEVELOPMENT ENVIRONMENTS PERFORMANCE

| Configuration | VM count | Startup time, sec | CPU usage, % | RAM usage, GB |
|---|---|---|---|---|
| Initial configuration (3 components) | 4 | 404,8 ± 4,1 | 49% ± 5% | 2,2 ± 0,18 |
| Alternative configuration (3 components) | 1 | 210,3 ± 4,6 | 23,8% ± 7,7% | 1,3 ± 0,13 |

TABLE IV.        RESULTS OF COMPONENT BUILDING TIME (SEC)

| Configuration | Researcher Account | | Psychotest Player | |
|---|---|---|---|---|
| | Build | Watch | Build | Watch |
| Initial configuration (3 components) | 120,4 ± 2,7 | 5 ± 0,4 | 131 ± 3,1 | 4,4 ± 0,1 |
| Alternative configuration (3 components) | 100,7 ± 2,5 | 3,6 ± 0,3 | 82,4 ± 2,4 | 3,3 ± 0,1 |
| Improved alternative configuration (3 components) | 90,1 ± 2,4 | 1,2 ± 0,3 | 79 ± 1,3 | 1,1 ± 0,2 |

Such configuration changes are also caused an start time increase (from 210.3 to 227.7 seconds in average), associated with the addition of the waiting time for NFS driver starting and the mounting the directory in virtual machine, and also used CPU resources increase (from 24% to 29% in average). However, indicators of resources consumed still remain significantly lower compared with the full launch of all components in original development environment. In this case, a decrease in the time spent on the components build operations was observed in using the proposed improvements.

## VI. RESULTS OF THE EXPERIMENTS

Experiments were conducted to estimate the launch time and the load on the computer for various configurations of the virtual development environment. The results of the experiments are given in Table V.

The results of experiments on evaluating the execution time of the complete component assembly (Build) and component reassembly (Watch) for the Platform components of the Researcher Account and Psychotest Player are presented in Table VI.

From the results of the experiments it can be seen that the proposed alternative configuration exerts a significantly lower load on the CPU (Fig. 3). An improved alternative configuration uses slightly more processor resources, which is associated with the use of additional driver, but even in this case, the load is much lower than with the full launch of the original development environment.

Similar changes are also noticeable in the used RAM comparing (Fig. 4).

TABLE V. THE RESULTS OF THE EVALUATION OF DEVELOPMENT ENVIRONMENT START UP TIME AND PERFORMANCE

| Configuration | VM count | Startup time, sec | CPU usage, % | RAM usage, GB |
|---|---|---|---|---|
| Initial configuration (API Server) | 2 | 289,4 ± 2,8 | 24,3% ± 6,7% | 1,5 ± 0,1 |
| Initial configuration (Researcher Account) | 1 | 118,5 ± 13 | 21,8% ± 6% | 0,5 ± 0,08 |
| Initial configuration (Psychotest Player) | 1 | 155,3 ± 19,1 | 16,1% ± 3% | 0,7 ± 0,04 |
| Initial configuration (API Server, Researcher Account) | 3 | 347,4 ± 4,4 | 35,4% ± 9,1% | 2 ± 0,2 |
| Initial configuration (API Server, Psychotest Player) | 3 | 290,1 ± 4,5 | 29% ± 5,3% | 2 ± 0,06 |
| Initial configuration (3 components) | 4 | 404,8 ± 4,1 | 49% ± 5% | 2,2 ± 0,18 |
| Alternative configuration (3 components) | 1 | 210,3 ± 4,6 | 23,8% ± 7,7% | 1,3 ± 0,13 |
| Improved alternative configuration (3 components) | 1 | 227,7 ± 3,4 | 28,7% ± 5,4% | 1,2 ± 0,09 |

TABLE VI. THE RESULTS OF THE EVALUATION OF THE COMPONENT BUILD TIME

| Configuration | Researcher Account | | Psychotest Player | |
|---|---|---|---|---|
| | Build | Watch | Build | Watch |
| Initial configuration (API Server) | 101,9 ± 2,9 | 3 ± 0,4 | 86,1 ± 2,1 | 3,1 ± 0,1 |
| Initial configuration (Researcher Account) | 100,6 ± 2,5 | 2,7 ± 0,4 | 91,8 ± 2,4 | 3,7 ± 0,1 |
| Initial configuration (Psychotest Player) | 99 ± 2,4 | 3,6 ± 0,4 | 89,3 ± 3 | 3,7 ± 0,1 |
| Initial configuration (API Server, Researcher Account) | 98,7 ± 2,2 | 3,7 ± 0,5 | 119,2 ± 2,4 | 4 ± 0,2 |
| Initial configuration (API Server, Psychotest Player) | 94,6 ± 2,7 | 4,1 ± 0,4 | 137,7 ± 3,6 | 4 ± 0,2 |
| Initial configuration (3 components) | 120,4 ± 2,7 | 5 ± 0,4 | 131 ± 3,1 | 4,4 ± 0,1 |
| Alternative configuration (3 components) | 100,7 ± 2,5 | 3,6 ± 0,3 | 82,4 ± 2,4 | 3,3 ± 0,1 |
| Improved alternative configuration (3 components) | 90,1 ± 2,4 | 1,2 ± 0,3 | 79 ± 1,3 | 1,1 ± 0,2 |



Fig. 3. Average CPU Load, in Percent.



Fig. 4. Average RAM used, GB.

Reducing the CPU load and the number of used virtual machines led to a significant decrease in the average time required to start the virtual development environment (Fig. 5). The improved alternative configuration is lower that alternative configuration without additional modifications, like in case with the load on the processor. This may also be due to the use of an additional driver and the need to wait for its initialization.

However, improved alternative configuration leads in performance in components build (Fig. 6), and also shows significantly better performance indicators in rebuild (Fig. 7). Such improvements are a direct consequence of using the NFS driver, which increases the speed of data exchange with the host machine. It significantly reduces the time required to build components.

Therefore, the acquired improved alternative configuration of the development environment for Vagrant shows a significant reduction in the load on the processor and RAM in comparison with the full launch of the original development environment. It also shows a significant acceleration of the component building process, which reduces the waiting time on the part of the developers and thereby increases their efficiency.



Fig. 5.    Average Startup Time, Sec.



Fig. 6.    Average Execution Time for the Complete Assembly of the Components for the Researcher Account and the Psychotest Player (Build), Sec.



Fig. 7.    Average Execution Time for the Reassembly of the Components for the Researcher Account and the Psychotest Player (Watch), Sec.

## VII. DISCUSSION

In the process of research and development of mechanisms to improve the effectiveness of the original software development environment, it was found that the alternative configuration of Vagrant based on a single virtual machine for all existing platform components provides a significant reduction in both the startup time and computer load.

Both Vagrant and Docker can be used to organize a software development environment. It is reasonable to select between them with respect to the technologies used at the server in order to ensure that the development and the deployment environments are as identical as possible.

Additional research may be conducted to measure the increase in the number of components of the platform and compare the performance of the original and alternative Vagrant configurations. A separate study can provide a more detailed consideration of the Docker container technology as a solution for the organization of the software development environment.

When choosing technologies and preparing the software development environment, it is necessary to consider the whole architecture of the application being developed and the particular features of its operation. It makes sense to choose Vagrant when using virtual machines and the Docker in the case of the container technologies. It should also be borne in mind that the hybrid design can be implemented with both the technologies used. Although, one shouldn't exclude the use of other technologies and configurations. In this regard, to assess the performance of the software development environment, it is necessary to conduct separate experiments and develop one's own testing methodology, which can be completely different from that presented in this article.

Switching from several virtual machines to the single one may not be allowed when configuring certain development environments. For example, if it is imperative to isolate a particular component of the system. In this case, the increase in the productivity of the software development environment might not be achieved and the developer will have to look for alternative approaches to solve this problem.

## VIII. Conclusions

Due to the increase of web projects development complexity and working environments setting up, there is a need to use virtual development environments.

At use of virtual development environment, which structure contains an autonomous virtual machine for each component, some problems with high resource consumption and a reduction in the performance of the developers' working machines were noticed, which made it necessary to consider alternative solutions.

The paper described an analysis of the reasons of performance reduction at use of development environments based on virtualization. The main technologies used for the development of virtual development environments are considered, and an improved structure is proposed. In addition, an experimental assessment of the original and alternative solutions was made.

The assessment of the configuration of the original software development environment showed that using a separate virtual machine for each of the components of the psychological platform to get as close as possible to the server structure was not an effective approach. At the developers' computers, there was a significant drop in performance, a low data exchange rate with the VM and a high component assembly time.

In the search for a solution, the containerization technology based on Docker and the alternative configuration of Vagrant using a single virtual machine for all existing components were considered. However, using Docker requires developing a completely new component structure and will cause the developer's environment to differ from the server environment. Therefore, the alternative configuration of Vagrant is the most preferred option within the framework of the developed platform. The measurements showed that the use of the alternative configuration significantly reduced the startup time and reduced the load on the computer.

As an efficiency enhancement, a driver based on the NFS protocol was applied and the configuration of the Webpack system was modified. The measurements showed that the use of the NFS entailed a slight increase in the time for launching the virtual machine but reduced the time for assembling the components of the platform.

It was shown that in the case of the open digital platform for mass psychological research the application of the alternative Vagrant configuration with the NFS driver and the Webpack optimization provided a significant performance boost compared to the original configuration.

The implemented alternative solution based on a single development environment showed significantly lower resource consumption, as well as a reduction in the tasks building time with the help of the driver for accessing file system based of the NFS protocol.

Therefore, the use of a virtual development environment based on a single virtual machine using the NFS driver can significantly reduce the workload of the developers' computers.

This increases rapidity and reduce time consumption, which improves the developers' efficiency.

Conducted studies, including the stages of parameter estimation, the introduced characteristics and criteria, can be the basis for the formation of a methodology for the experimental evaluation of the software development environment configurations, which would allow choosing effective solutions at the design stage.

## References

[1] Caballer M., Blanquer I., Moltó G., de Alfonso C. (2015) Dynamic management of virtual infrastructures, Journal of Grid Computing, 13(1), 53-70. doi: 10.1007/s10723-014-9296-5

[2] Giannakopoulos I., Konstantinou I., Tsoumakos D., Koziris N. (201) Cloud application deployment with transient failure recovery, Journal of Cloud Computing, 7(1), 11. doi: 10.1186/s13677-018-0112-9

[3] Vagrant, 2019. Available at: https://www.vagrantup.com/ (accessed 27.03.2019).

[4] Spanaki P., Sklavos N. (2018) Cloud Computing: Security Issues and Establishing Virtual Cloud Environment via Vagrant to Secure Cloud Hosts. In Computer and Network Security Essentials. Springer, pp. 539-553. doi: 10.1007/978-3-319-58424-9_31

[5] Hashimoto M. (2013) Vagrant: Up and Running: Create and Manage Virtualized Development Environments. O'Reilly Media Inc, 2013.

[6] Xuan N. P. N., Lim S., Jung S. (2017) Centralized management solution for vagrant in development environment, In Proceedings of the 11th International Conference on Ubiquitous Information Management and Communication. ACM,. art. no. 37. doi: 10.1145/3022227.3022263

[7] Thompson C. (2015) Vagrant virtual development environment cookbook. Packt Publishing Ltd.

[8] Mouat A. (2016) Using Docker: Developing and Deploying Software with Containers. O'Reilly Media Inc.

[9] Sammons G. (22016) Learning Vagrant: Fast programming guide. CreateSpace Independent Publishing Platform.

[10] Peacock, M. (2015) Creating Development Environments with Vagrant. Packt Publishing Ltd.

[11] Iuhasz G., Pop D., Dragan I. (2016) Architecture of a scalable platform for monitoring multiple big data frameworks, Scalable Computing: Practice and Experience, 17(4), 313-321. doi: 10.12694/scpe.v17i4.1203

[12] Nikulchev E., Ilin D., Kolyasnikov P., Belov V., Zakharov I., Malykh S. (2018) Programming Technologies for the Development of Web-Based Platform for Digital Psychological Tools, International Journal of Advanced Computer Science And Applications, 9(8), 34-45. doi: 10.14569/IJACSA.2018.090806

[13] Kashyap S., Min C., Kim T. (2016) Opportunistic spinlocks: Achieving virtual machine scalability in the clouds, ACM SIGOPS Operating Systems Review, 50(1), 9-16. doi: 10.1145/2903267.2903271

[14] Saikrishna P. S., Pasumarthy R., Bhatt N. P. (2017) Identification and multivariable gain-scheduling control for cloud computing systems, IEEE Transactions on Control Systems Technology, 25(3), 792-807. doi: 10.1109/TCST.2016.2580659

[15] Li J., Xue S., Zhang W., Qi Z. (2017) When i/o interrupt becomes system bottleneck: Efficiency and scalability enhancement for sr-iov network virtualization, IEEE Transactions on Cloud Computing, Early Access. Doi: 10.1109/TCC.2017.2712686

[16] Basok B.M., Zakharov V.N., Frenkel S.L. (2017) Iterative approach to increasing quality of programs testing, Russian Technological Journal, 5(4), 43-12.

[17] Docker, 2019. Available at: https://www.docker.com/ (accessed 27.03.2019).

[18] Chen M., Bangera G. B., Hildebrand D., Jalia F., Kuenning G., Nelson H., Zadok E. (2017) vNFS: maximizing NFS performance with compounds and vectorized I/O, ACM Transactions on Storage. 13(3), 21. doi: 10.1145 / 3116213

[19] Kane S.P., Matthias K. (2018) Docker: Up & Running: Shipping Reliable Containers in Production. O'Reilly Media Inc, 2018.

[20] Peinl R., Holzschuher F., Pfitzer F. (2016) Docker cluster management for the cloud-survey results and own solution, Journal of Grid Computing, 14(2), 265-282. doi: 10.1007/s10723-016-9366-y

[21] Krieger, M. T., Torreno, O., Trelles, O., & Kranzlmüller, D. Krieger M. T. et al. (2017) Building an open source cloud environment with auto-scaling resources for executing bioinformatics and biomedical workflows, Future Generation Computer Systems, 67, 329-340. doi:10.1016/j.future.2016.02.008

[22] Vagrant WinNFSd–GitHub, 2019. Available at: https://github.com/ winnfsd/vagrant-winnfsd (accessed 27.03.2019).

[23] Vagrant bindfs – GitHub, 2019. Available at: https://github.com/gael-ian/vagrant-bindfs (accessed 28.03.2019).

[24] Dani S. A. (2017) JavaScript by Example. Packt Publishing.

[25] Webpack, 2019. Available at: https://webpack.js.org/ (accessed 28.03.2019).

# Prediction of Crude Oil Prices using Hybrid Guided Best-So-Far Honey Bees Algorithm-Neural Networks

Nasser Tairan[1], Habib Shah[2], Aliya Aleryani[3]
Department of Computer Science, College of Computer Science
King Khalid University Abha, Saudi Arabia

*Abstract*—The objective of this paper is the use of new hybrid meta-heuristic method called Guided Best-So-Far Honey Bees Inspired Algorithm with Artificial Neural Network (ANN) on the Prediction of Crude Oil Prices of Kingdom of Saudi Arabia (KSA). Very high volatility of crude oil prices is one of the main hurdles for the economic development; therefore, it's the need of the hour to predict crude oil prices, especially for oil-rich countries such as KSA. Hence, in this paper, we are proposing a hybrid algorithm, named: Guided Best-So-Far Artificial Bee Colony (GBABC) algorithm. The proposed algorithm has been trained and tested with ANN for finding the optimal weight values to increase the exploration and exploitation process with balance quantities to obtain the accurate prediction of crude oil prices. The KSA crude oil prices of the five years 2013 to 2017 have been used to train ANN with different topologies and learning parameters of the proposed method for the prediction of the crude oil prices of the next day. The simulation results have been very promising and encouraging of the proposed algorithm when compared and analyzed with ABC, GABC (Gbest Guided ABC) and Best-So-Far ABC methods for prediction purpose. In most cases, the actual prices and predicted crude oil KSA prices are very close, which were obtained by the proposed GBABC method based on the optimal weight values of ANN and minimum prediction error.

*Keywords*—*Bio inspired; best so far; crude oil prices; KSA*

## I. INTRODUCTION

Crude oil, commonly known as petroleum comes from an oil well, is a liquid or solid found within the explored earth comprised of hydrocarbons (compounds composed mainly of hydrogen and carbon), organic compounds and small amounts of metal [1]. It is one of the most important energy resources on earth for humans and machine developments. So far, it remains the world's leading fuel, with nearly one-third of global energy consumption.

Among the other natural resources, crude oil is the "key" treasure for obtaining the stable economic position of the country. It has an important factor affecting the local and global economy of the region as well. For the last two decades, the oil-rich countries have witnessed significant margins and differences in their development rate, economic stability as well as the quality of their organizations. KSA is one of the oil-rich, productive hub resources in the Middle East region. Its cover more than half of the OPEC's total oil exports, and is a major player in setting the oil price in Asia and worldwide [2]. The country, being the second largest producer of petroleum liquids and the largest exporter of crude, has the ability to have a major impact on the global oil industry and economic

stability. KSA has 265.8 billion barrels of crude reserves, the second largest in the world, amounting to 16% of the world's reserves in 2014, which are predicted to rise to 273 billion by the end of 2017 [3]. KSA is also the second largest producer of crude oil in the world. From May to June 2017, the crude Oil Production in Saudi Arabia increased 9880 BBL/D/1K to 10070 BBL/D/1K [4].

It has the highest refining capacity among the OPEC producers and has plans to add 1.2 million bbl/d more by 2020 to the current refining capacity of crude oil of 2.9 million bbl/d. In 2016, the top three crude oil producing countries were Saudi Arabia (10.46m b/d), Russia (10.29m b/d) and the United States (8.88m b/d)[5]. It's come in twelfth largest primary energy consumer in the world like around 1 million bbl/d of oil is used for electricity generation during the heat waves season [2].

The famous upstream and downstream companies operating in KSA are Saudi Aramco, Shell, Total, Chevron, Sinopec, ExxonMobil, Sumitomo and Eni. The country has around 100 oil fields; the five famous in the production are Safaniya, Khurais, Manifa, Shaybah and Ghawar. The Ghawar oil field has the maximum production and estimated proved oil reserves of 70 billion barrels as of 2014. The production of oil in 2014 was 9.7 million bbl/d, which was 32% and 10.5% of the OPEC's and the global crude oil production respectively. International Energy Agency, Energy Information Administration has represented the crude oil production of global [6], OPEC and KSA with the following Fig. 1 [7].

Predicting oil prices brought a considerable attention by scientific researchers and authors to study it from different aspects and different categories who provide a preemptive knowledge in identifying potential candidate forecasting models for crude oil prices. Unfortunately, with inadequate information, uncertain situation, too many variables, the intrinsic complexity of oil market mechanisms, imbalance between production and consumption, weather forecasting, and imprecise elements, the crude oil price system is extremely complex for modeling analyses, and its dynamics are hard to predict its future price [8]. Oil prices are confined between demand and supply framework, oil price volatility analysis and Oil price forecasting [9].

The world crude oil prices from 2004 to 2018 are given in Fig. 2 [10], [11]. Due to the high volatility in the crude oil price, Saudi Arabia has found itself between a rock and a hard place lately. Based on the previous study [9], [10], [11], a relationship between crude oil price, economic activities,

energy supply and distribution and the gross domestic product (GDP) growth rate are asymmetric [12]. The role of oil price in the prediction of economic growth has been investigated [4], [13] which demonstrated a higher level of predictability of 28 developed countries. The petroleum sector accounts for roughly 87% of budget revenues, 42% of GDP, and 90% of export earnings [14].

The annual GDP of KSA (The recent years) is decreasing due to the high volatile of unreliable crude oil price in the global market as given in Fig. 3. Therefore, the economist exploiting new income resources for stable GDP value through different innovative, nationalization, global business strategies, promoting education and energy sector through mission 2030 scheme. Due to its high complexity and price volatility, researchers adopted various statistical, mathematical, artificial intelligence, hardware base, computer science knowledge and software based and engineering methods for the prediction of accurate price of crude oil of KSA. The Autoregressive integrated moving average (ARIMA) and gene expression programming (GEP) techniques used to predict crude oil prices over the period from January 2, 1986 to June 12, 2012 [16].



Fig 1.    Crude Oil Production of 2018 (average per month).



Fig 2.    Fifteen Years (2004-2018) Global Crude Oil Prices[15].



Fig 3.    The GDP of KSA from 2008 to 2017 (Million USD).

The GEP model outperforms the ARIMA has the highest explanatory power as measured by the R-squared statistic. However, the GEP didn't successes in obtaining high prediction accuracy. The hybrid model integrating wavelet and multiple linear regressions (WMLR) was proposed for crude oil price forecasting of WTI obtained higher accuracy than regular LR (Learning Rate), ARIMA, and generalized autoregressive conditional heteroscedasticity (GARCH) model [8], however this method used multiple steps such as Particle Swarm Optimization (PSO) for adopting parameters and principal component analysis for processing subseries data. The ARIMA, GARCH, Belief networks, k-means clustering and an empirical mode decomposition (EMD) methods have used for crude oil forecasting in the last two decades [9], [17]. However, due to the volatility, nonlinearity, and irregularity, the classical and econometric model can lead to the decrease of the accuracy. Therefore, due to the above mentioned limitations of the classical science, mathematical and statistical approaches, soft-computing models can provide powerful solutions to nonlinear crude oil price prediction [18].

Many experiments found that the computational intelligent algorithms which simulate the way humans' reason by incorporating their rate of efficiency and randomness during decision making. They have broadly divided into techniques based on modelling of human mind and nature inspired algorithms [19], [20]. These methods often have some advantages over typical mathematical and statistical-based models [21]. However, these models also have their own drawbacks, such as ANN often suffer from local minima and over-fitting problems, while SVM and GP, including ANN, are sensitive to parameter selection and suitable architectures and dataset behaviors [22], [23]. Also, the standard learning algorithms of the ANN models have the same local minima trapping problem and slow convergence speed [24], [25].

To remedy the above shortcomings of typical models and learning algorithms, meta-heuristic and their hybrid version methods have been used recently to solve time series problems. These are: Artificial Bee Colony (ABC), Cuckoo Search (CS), Bat Algorithm (BA), evolutionary algorithms (EA), genetic algorithm (GA), PSO, Ant Colony Optimization, hybrid algorithms and so on [26]–[30]. These meta-heuristic learning algorithms are more efficient and famous due fast convergence, high efficiency, easy to understand and implementation and robustness for some problems [31], [32]. They can easily find the set of best weight values, through suitable parameter selection, activation function and network structure that will cause the output from the MLP to match the actual target values as closely as possible.

Sections 2 and 3 contain neural networks and crude oil prices forecasting details and honey bees super heuristic optimization methods respectively. The proposed method explained in Section 4. The simulation results and conclusion are added in Sections 5 and 6.

## II.    Neural Networks and Oil Prices Forecasting

Artificial neural network (ANN) often called a "Neural Network" or simply Neural Net (NN) inspired by biological system is an interconnected set of artificial neurons that uses a mathematical model or computational model for information

processing based on external or internal information that flows through the network [33]. The biological neuron structure is given in Fig. 4, which have different functions terminologies such as dendrites, which receive activation from other neurons, some processes which converts incoming activations into output activations, axons act as communication route, synapses, neurotransmitters and nucleus. This field goes by many names, such as parallel distributed processing, Neurocomputing, natural intelligence systems, machine learning algorithms, bio inspired learning methods, deep learning and Multilayer Perceptron (MLP), and others soft computing methods [33], [34], [35].



Fig 4. Structure of biological neuron.

TABLE I. PREVIOUS RESEARCH ON CRUDE OIL PRICES PREDICTION

| ANN Model | Input Patterns | Objectives | Findings |
|---|---|---|---|
| Feedforward with Levenberg-Marquardt algorithm [44] | Daily closing price from Sep 2002 to Aug 2013. | short-term prediction | better than the LSM |
| feed-forward neural network and Radial Basis Function[16] | January 2, 1986 to June 12, 2012. | forecast oil price | Less mean squared error than gene expression programming |
| Artificial Neural Networks-Quantitative with BP [45] | Crude oil dataset of 1984 to February, 2009 with quantitative key factors influencing. | monthly WTI crude oil price prediction | ANN-Q model effective than Hierarchical Conceptual model |
| Empirical Mode Decomposition with Feed-forward Neural Network and Adaptive Linear Neural Network [46]. | The data frequency is daily closing price; from Sep 1996 to Aug 2007 of WTI | short-term prediction up to three days ahead | Intraday data for crude oil prices is not available, overall improvement was insignificant |
| Functional Link Artificial Neural Network And MLP [47] | US dollar index, S&P 500 stock price index, gold spot price, heating oil spot price and US crude oil spot price are employed | predict the next day's spot price of US crude oil | Functional link artificial neural network performed better than standard MLP. |
| Stream learning, Random guess and Forecast combination models [48] | U.S. refiner acquisition cost for crude oil imports and WTI crude oil spot price, | Predict monthly crude oi price | Highest accuracy than no-change and ANN, Random guess and Forecast combination models. |
| factor augmented artificial neural network(FAANN) [40] | South African monthly panel, namely, deposit rate, gold mining share prices and Long-term interest rate, using monthly data over the in-sample period (training set) 1992:1–2006:12. | to forecast 3, 6 and 12 month-ahead forecasts | FAANN model yields substantial improvements over the autoregressive AR benchmark model and standard dynamic factor model (DFM). |
| Time-varying weight combination approach [49] | factors: supply, demand, crack spread, and non-energy commodity prices. West Texas Intermediate crude oil prices Dataset | mid-horizons and long-horizons | Limited horizons |
| deep learning ensemble approach[50] | West Texas Intermediate price series, flow and stock series, and macroeconomic and financial series | Monthly WTI crude oil price | Effective with largest computational cost |
| Deep Learning based Model [51] | historical data of the WTI crude oil market, July 23, 2007 to February 24, 2017 | Daily crude oil price movement prediction | Limited deep learning models |
| convolutional neural network[52] | Brent crude oil generic series of the first month's futures prices, traded on the Intercontinental Exchange (ICE) from 24 June 1988, to 3 November 2018. | short-term crude oil futures prices prediction | Only working for short-term crude oil futures prices. |

Fig 5.    Artificial neural network.

Based on the biological neuron of Fig. 4, McCulloch-Pitts introduced a simplified model called artificial neuron model as given in Fig. 5. The output formula of the artificial neuron model is given in Equation (1).

$$\textbf{Output} = \textbf{sgn}\left( \sum_{i=1}^{n} \textbf{Input i} - \Phi \right) \qquad (1)$$

Initially, it was developed for solving linear problems; later on it has been extended to different models such as Multilayer Feedforward Neural Network (MFFN), Recurrent Neural Network, Probabilistic Neural Network, Pi Sigma Neural Network for solving different optimization problems [33], [34]. MFFN is a famous among the various NN structure most commonly used due to its lower complexity and ability to produce satisfactory results for different problem domains as given in Fig. 6.

The previous study shows that the NN performance can be improved through the selection of suitable structure, activation function, appropriate numbers of input pattern and of course learning algorithms. Besides the common applications, ANN tools are very effective for financial, stationary and non-stationary, meteorological, natural hazards, stock values and crude oil price time series data prediction. For two decades, ANN tools are famous, effective and attractive for prediction time series dataset including crude oil prices, bottom hole pressure in vertical multiphase and stock exchange values [36].

Using MFFN model reached on minimum error and high accuracy through different learning strategies such as supervised, unsupervised and reinforced, such as, backpropagation, genetic algorithm, gradient descent and so on. Typical training algorithm, BP has some drawbacks like slow convergence and trapping in local minima [37], [38]. The Bayesian approach used to predict crude oil price through the various independent variable factors such as world oil demand and supply, the financial situation, upstream costs, and geopolitical events. The results show that the crude oil price is estimated to increase to $169.3/Bbl by 2040 [39].

However, these models used the traditional machine learning algorithms, rely on a fixed set of training data to train a machine learning model and then apply the model to a test set, but may not be effective for non-stationary time series data such as oil price data and other nonlinear complex dataset. From the above-mentioned approaches, that ANN model is more sufficient and effective than standard models; the performance can be easily increased rapidly through robust and

efficient learning algorithm [23], [40], [41]. Bio-Inspired methods are robust and attractive optimization algorithm especially for solving nonlinear complex problems [42], [43]. The previous histories of softcomputing methods for crude oil prices prediction are summarized in Table I.

### III. HONEY BEES META-HEURISTIC OPTIMIZATION METHODS

#### A. Bio-Inspired Artificial Bee Colony

The Artificial Bee Colony (ABC) is a robust bio-inspired learning method, proposed in 2005 by Karaboga [53] to solve the complex nonlinear optimization problem in multivariable functions, is a relatively innovative meta-heuristic optimization algorithm that is based on the social behavior of honey bee colony named: employed bees, onlooker bees and scouts during searching and managing for food sources (FS) [54]. The first half of the colony consists of the employed bees and the second includes the unemployed. From the different numerical and statistical performance measure demonstrate that the ABC algorithm is competitive with other types of meta-heuristic and typical algorithms [55]–[57]. The technical duties of the employed and unemployed artificial bees are given in details.

Each employed bee search around the food source, gathering required information about its quality and position of the onlookers. Then, they carry the information about the position of food source back to the hive and share this information with artificial onlooker bees by dancing in the nearby hive. Onlooker bees: the onlookers tend to choose the best food sources to further exploit, based on information communicated by the employed bees through their dances. Therefore, good food sources attract more onlooker bees compared to the bad ones. The artificial onlooker bees choose the best food source with better quality based on information communicated from those found by employing bees using different ways, such as a probability selection mechanism, greedy selection, fitness function as a proportional of the quality of the food source. The last bees processes are managing by the scout bee group is responsible for the exploration process randomly chose a new good food source to replace the old one. The number of food sources (based on position and quality) which represents a possible solution to the optimization problem and fitness of the associated solution is equal to the number of employed bees and also equal to the number of onlooker bees. The employed, onlooker bees used for exploitation process for a given problem towards best solution space given in Equation (2), while scout bees Equation (3) for exploitation process.

$$v_{ij} = x_{ij} + \Phi_{ij}(x_{ij} - x_{kj}) \qquad (2)$$

Where $v_{ij}$ is a new solution in the neighborhood of $x_{ij}$ for the employed bees, $k$ is a solution in the neighborhood of $i$, $\Phi$ is a random number in the range [-1,1].

$$x_{ij}^{rand} = x_{ij}^{min} + rand(0,1)(x_{ij}^{max} - x_{ij}^{min}) \qquad (3)$$

Although typical ABC is famous due to its robustness and high efficiency for clustering, classification and numerical function optimization problems, however, due to the same and

the random searching approach of exploration cannot guaranty for finding the best food position, also sometimes it trapped in local minima. The researchers improved typical ABC algorithm by different strategies such as, Best-So-Far, discrete, hybrid, gbest guided and quick within employed, onlookers and scout bees. The typical artificial bee colony model which includes three kinds of bees considering the division of labor: employed bees, onlooker bees and scout bees. Each employed bee works on only one food source as given in Fig. 6.

### B. Best-So-Far Artificial Bee Colony Algorithm

The standard ABC algorithm is a unique bio inspired algorithm inspired through the attraction and natural foraging behaviors of honey bees. It has been successfully used for solving different statistical, mathematical, science and engineering problems. The exploration and exploitation are the famous process of ABC by employed and unemployed bees [58]–[62].

Equations (2) and (3) were used for exploitation through employed and onlooker bees and scout bee for exploration process. The first two bee group used the same exploitation process base on random way which cannot guarantee for optimized and fast solution. Therefore, the Best-So-Far ABC proposed [63], is one of the efficient bio inspired algorithm among the ABC improved and hybrid algorithms. It was developed to enhance the exploitation and exploration processes of typical ABC through different strategies: the Best-So-Far method, an adjustable search radius, and an objective-value-based comparison method. Equation (2) of the typical ABC has been reused in the exploration process. The Best-So-Far ABC updated the exploitation process of onlooker bees section by the following Equation (4) as:

$$v_{id} = x_{ij} + \Phi_{ij} f_b (x_{ij} - x_{bj}) \qquad (4)$$

where: $v_{id}$ is the new candidate food source for onlooker bee position $i$ dimension $d$, d = 1,2,…,D; $x_{ij}$ is the selected food source position i in a selected dimension j; $\Phi$ = a random number between −1 and 1; $f_b$ is the fitness value of the best food source so far and $x_{bj}$ is the Best-So-Far food source in selected dimension $j$.



Fig 6.   A Typical Bee Colony Model The artificial bee colony includes three kinds of bees considering the division of labor: employed bees, onlooker bees and scout bees. Each employed bee works on only one food source.

The best food source position is calculated by processing the information received from all employed bees. The onlooker bees using Best-So-Far method will record the best selected food source position so far within the new candidate generation function to enhance the exploitation process. Then the fitness method used to calculate fitness based on that particular best solution, which used to compare with other previous fitness values, the position is updated with Equation (4) based on the fitness level. The second modification made in the adjustable search radius. This is done because of the need to get out of the local optimum solution problem.

$$v_{id} = x_{ij} + \Phi_{ij} \left[ w_{\max} - \left( \frac{iteration}{MCN} \right) (w_{\max} - w_{\min}) \right] x_{ij} \qquad (5)$$

Where $v_{id}$ is a new achievable solution of a scout bee that is modified from the present position of an abandoned food source ($x_{ij}$), the value of $w_{max}$ and $w_{min}$ represent the maximum and minimum percentage of the position adjustment for the scout bee. The last change is finding the minimum objective value, here compare and to select between the old solution and the new solution in each iteration is done by the fitness value of the following Equation (6) as,

$$\text{Fitness(f(x))} = \begin{cases} \dfrac{1}{1 + f(x)} & if \ \ f(x) \geq 0 \\ 1 + |f(x)| & if \ \ f(x) < 0 \end{cases} \qquad (6)$$

The above proposed method is efficient and fast convergence through three modifications for solution update, increase the local search ability of the onlooker bees, maintain the diversity of new food sources by random scout bees, and to resolve round up issues in the computation of the floating point "goodness" value. However, the employed bees section used the random way Equation (2) of typical ABC, which can lead to unbalance of exploration and exploitation process for different complex problems.

### C. Gbest Guide Artificial Bee Colony Algorithm

The typical ABC algorithm used the random searching and selecting methods through employed, onlooker and scout bees, unfortunately no global best solution, best-so-far, local, best or mutation method use to improve and balance the exploration and exploitation process successfully. Based on random searching ways, ABC has global search capability, but poor local search capability. In order to enhance the exploitation capability of the ABC algorithm through candidate solutions like employed bee. Different improved version of typical ABC have been developed for enhancing and balance exploration and exploitation process by introducing different operators, strategies and operators. Gbest Guided Artificial Bee Colony (GABC) algorithm is one of the attractive modified bio inspired algorithm developed to increase the performance of typical ABC [59]. Equation (2) has been modified by the Equation (7) to direct the search path towards global optima.

$$v_{ij} = x_{ij} + \Phi_{ij}(x_{ij} - x_{kj}) + \Psi_{ij}(y_j - x_{ij}) \qquad (7)$$

where $v_{ij}$ (or $x_{ij}$) is a new solution in the neighborhood of $x_{ij}$ for the guided employed bees and  onlooker bee, $\Psi$ij is a random number in the interval between [0,C], and the term $y_j$

represents jth element of the global best solution in current generation. The equation (7) used to enhance the exploitation process through gbest $x_j$ and c values. The proposed gbest guided approach is more effective than typical ABC especially for numerical function optimization, time series prediction and for environmental/ economic dispatch considering wind power and classification task [64], [65].

## IV. PROPOSED GBEST GUIDED BEST-SO-FAR ABC ALGORITHM

Bio inspired agents have motivated many researchers to develop various mathematical and computational approaches for solving complex optimization problems [28], [43], [66], [67]. Theoretically, it is because of their unique movements of searching, gathering, sharing, dancing, foraging, selection, flying, managing, building, communication, social, emotional, their foraging behavior, their mating and reproduction behavior, their pheromone laying behavior and navigation behavior, self-adapting and self-organizing characteristics [68], [69].

These bio inspired artificial agents or methods have successfully applied to various applications and linear, nonlinear complex problems such as classification, clustering, time series prediction, numerical function optimization and other combinatorial problems [70]. The performance of these algorithms depends on exploration, exploitation, balancing, convergence and global optimum position. The typical ABC easily trapped in local minima which lead to slow convergence. Beside ABC, the GABC and Best-So-Far ABC is the modern improved examples in success history as mentioned in the above sections [58], [60].

Taking the advantages of high exploration and exploitation process from GABC and Best-So-Far algorithm, a new hybrid method is proposed called a Guided Best ABC algorithm. The proposed GBABC algorithm will increase the effectiveness of typical ABC, GABC and Best-So-Far ABC algorithm [46], [62], [71], [72]. The GBABC will first use the gbest strategy through employing bee to increase the exploitation process according to the given problem, while the onlookers and scout bees will adapt the Best-So-Far method for exploration process with balance quantity.

The searching strategy of employed bees given by Equation (7) will increase the exploitation around the, current best ever food source through the guided strategy to enhance the real foods source position. These food source positions will be shared with the best so far onlooker and scout bees which will further select the most appropriate position for the given problem. In order to accelerate the convergence speed of the GBABC algorithm, the idea of focusing the search around the current best ever food source gbest, best so far, the best fitness formula was proposed. These bees will repeat their intelligence procedures until a predetermined maximum number of cycles (MCN) or the best food source achieved so far. Of course, the enough and balance exploitation and exploration process will successfully escape local optima trapping and slow convergence difficulties. The proposed GBABC method has been used to train the feed-forward neural networks for the crude oil prediction purpose; the details are added in the following section. Beside the proposed GBABC approach, the

main difference in the implementation phase which is executing based on max cycle numbers instead of max fitness number of evaluations, which are used by the typical ABC, GGABC, Best-so-far ABC and so many others algorithms.

**The pseudo-code of GBABC algorithm**
**Start**
**// Initialization**
Initialize the control variables and food source positions;
Evaluate the nectar amount of food sources; /* by using the following equation*/

$$fitness_i = \begin{cases} \dfrac{1}{1+f_i} & , \ if \ \ f_i \geq 0 \\ 1+|f_i| & , \ if \ f_i < 0 \end{cases} \qquad (8)$$

Cycle=1
**Repeat** */(the termination conditions are not met) /* (Max Cycle Numbers=as mentioned in Table)
 **/*Gbest Guided Employed Bees' Phase*/**
**FOR** (each gbest guided employed bee)
Produce a new food source following the Eq. (7);
Evaluate the fitness of the new food source;
Calculate the fitness of the explored food source best so far employed bees.

$$fit_i = \frac{fit_i}{\sum\limits_{i=1}^{SN} fit_i} \qquad (9)$$

Calculate the fitness values and normalize Pi values into [0,1] using Eq. (8)
Memorize the best solution of Gbest Guided Employed Bee So far;
**END FOR**
Calculate the probability P for each food source following the Eq. (10);

$$p_i = \frac{fitness_i}{\sum\limits_{j=1}^{SN} fitness_j} \qquad (10)$$

**/* Best-So-Far Onlooker Bees' Phase*/**
**FOR** (each Best-So-Far onlooker bee)
Send Best-So-Far onlooker to food sources depending on P;
Produce a new food source following the Eq. (4);
Evaluate the fitness of the new food source;
Apply greedy selection on the new FS and the old one;
Memorize the best solution of Best-So-Far Onlooker So far;
**END FOR**
**/* Best-So-Far Scout Bees' Phase*/**
**IF** (an gbest guided employed bee becomes into a Best-So-Far scout bee)
Send the Best-So-Far scout bee to a new randomly produced food source;
**END IF**
Determine the abandoned solution (source), if exists, replace it with a new randomly produced solution $x_i$ for the Best-So-Far scout bee using the equation (5),

Memorize the best solution achieved so far

Cycle=Cycle+1

Until Cycle=MCN

**End**

## V. EXPERIMENTAL EVALUATION AND ANALYSIS

In this research paper, crude oil prices time series data set of KSA used for one step ahead prediction of the year 2013, 2014, 2015, 2016 and 2017 respectively, which are used in the original form for prediction tasks [73]. The oil prices on Monday, Tuesday, Wednesday and Thursday were used to predict the next week values. Here, the crude oil settlement prices parameter used for predicting the next value versus days. The details of the above dataset have been given in Table II [73]. The various ANN topologies and algorithms parameters were set according to Table III. The main objectives of the ABC, GABC, BABC and GBABC algorithms are to train an FFNN for the prediction of crude oil price through optimal weights, bias, and hidden layer neurons. FFFNN trained with varied number of hidden layer neurons so that the best number of neurons could be obtained through the exploitation and exploration process of the Best-So-Far ABC, GBAC and GBABC algorithms. Different methods used for simulation analysis, which are Mean Square Error (MSE), Normalized Mean Square Error (NMSE), Mean Absolute Percentage Error (MAPE), Root Mean Square (RMSE), accuracy and success rate based on runtimes number Equations (12, 13 and 14). The FFNN trained with 75 % and 25 % of the whole data set with sigmoid function.

TABLE II. CRUDE OIL PRICES DATASET BEFORE PREPROCESSING

| Dataset durations | Total days | Dataset |
|---|---|---|
| 19-08-2013   to 31-12-2013 | 97 | 94 |
| 01-01-2014   to 31-12-2014 | 261 | 252 |
| 01-01-2015   to 31-12-2015 | 261 | 252 |
| 01-01-2016   to 31-12-2016 | 262 | 253 |
| 03-01-2017   to 01-08-2017 | 155 | 149 |
| Total= 19-08-2013   to   01-08-2017 | 1036 | 1000 |

$$MSE = \frac{1}{N} \sum_{i=1}^{n} (d_i - y_i)^2 \qquad (11)$$

$$NMSE = \frac{\sum_{i=1}^{n} (d_i - y_i)^2}{\sum_{i=1}^{n} (d_i - \bar{d}_i)^2} \qquad (12)$$

$$Accuraccy = \left( \frac{1}{n} \sum \frac{|d_i - y_i|}{|d_i|} \right) * 100 \qquad (13)$$

$$MAPE = \frac{1}{n} \sum_{t=1}^{n} \left| \frac{y_t - \hat{y}t}{y_t} \right| \qquad (14)$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{t=1}^{n} (y_t - \hat{y}t)^2} \qquad (15)$$

TABLE III. SETTING OF ANN TOPOLOGIES AND ALGORITHM PARAMETERS

| Algorithm | No of Inputs | Hidden Nodes | Food Sources | Limit |
|---|---|---|---|---|
| ABC | 4 | 2-10 | 10-20 | 5,10,15,20 |
| GABC | 4 | 2-10 | 10-20 | 5,10,15,20 |
| Best-So-Far | 4 | 2-10 | 10-20 | 5,10,15,20 |
| GBABC | 4 | 2-10 | 10-20 | 5,10,15,20 |

The performance of the ABC and their various updated algorithms depends on the control parameters: such as Colony Size (which can affect the processing speed and optimality values), limit. Through the proper selection of these parameters' values, the exploration and exploitation can be achieved in a balance way with global optima as well.

The simulation results of GBABC compared with ABC, GBAC and Best so far ABC algorithms. The average MSE training and testing, NMSE, accuracy and success rate are given from Table IV to Table VIII, respectively. In terms of MSE training, the ABC, GABC and Best-So-Far reached to 0.000371, 7.12E-05, 3.93E-06 respectively, while the proposed GBABC algorithm success to obtain the most minimum value than typical algorithms with error 3.13E-11 as given in Table IV. In case of 4-8-1, the Best-So-Far method reached to 1.10E-08 which is the minimum error than obtained by ABC, GABC methods.

From the NMSE results as given in Table V, Best-So-Far outperformed than ABC and GABC and proposed GBABC algorithm outperformed than all typical algorithms except Best-So-Far with 4-7-1 NN topology. The proposed GBABC obtained the best MSE out of sample data set for crude oil prices prediction than others bio inspired methods as shown in Table VI. Here, again the Best-So-Far method successes in obtaining least testing error when the hidden nodes reached to seven. Based on overall result, the GBABC model is reliable and promising for the accurate prediction of crude oil prices.

The average accuracy and success rate of the proposed and typical methods are given in Table VII and Table VIII. Where the success rate implies the percentage of getting success each time when the program is run with different control parameters. The GBABC prediction accuracy reached at 97.98, 98.70, 99.15, 99.10, 99.87 and with 100% success rate are higher than ABC, GABC and Best-So-Far ABC algorithms. Therefore, the performance accuracy of GBABC can be considered consistent. The ANN structure is based on Dimension (D) value which represents the number of input layer (with nodes), hidden Layer (with nodes), biases values and output node as well. The best MSE training results when CS=8, D=13 and C=1.0, CS=10, D=21 and C=1.2 and CS=10, D=21 and C=1.2 are given in Tables IX, X and XI. The Best-So-Far and proposed methods obtained the best MSE training error from ABC and GABC algorithms when CS=10, D=21 and C=1.5 and 7, 8, 9 hidden nodes.

TABLE IV.    AVERAGE MSE TRAINING FOR CRUDE OIL PRICES PREDICTION

| NN Structure | ABC | GABC | Best-So-Far | GBABC |
|---|---|---|---|---|
| 4-2-1 | 0.001234 | 1.16E-03 | 9.73E-05 | **3.17E-07** |
| 4-3-1 | 0.001202 | 1.61E-03 | 9.26E-05 | **3.99E-07** |
| 4-4-1 | 0.001034 | 1.09E-03 | 9.80E-05 | **3.12E-09** |
| 4-5-1 | 0.001042 | 1.01E-03 | 7.65E-05 | **1.39E-09** |
| 4-6-1 | 0.000836 | 1.23E-04 | 4.53E-05 | **1.02E-09** |
| 4-7-1 | 0.000491 | 1.01E-04 | 3.01E-05 | **1.01E-09** |
| 4-8-1 | 0.000402 | 9.81E-05 | **1.10E-08** | 8.11E-10 |
| 4-9-1 | 0.000371 | 7.12E-05 | **3.93E-06** | 3.13E-11 |

Table XI and XII contains the best NMSE and MSE out of sample for oil prices prediction when CS=10, D=21 and C=1.2 and CS=10, D=21 and C=1.2. From these tables, the error is rapidly decreasing when the numbers of hidden nodes are decreasing.

TABLE V.    AVERAGE NMSE FOR CRUDE OIL PRICES PREDICTION

| NN Structure | ABC | GABC | Best-So-Far | GBABC |
|---|---|---|---|---|
| 4-2-1 | 0.001131 | 1.16E-04 | 1.07E-06 | **1.01E-07** |
| 4-3-1 | 0.001044 | 1.61E-04 | 1.93E-06 | **1.20E-08** |
| 4-4-1 | 0.001012 | 1.09E-04 | 3.18E-06 | **3.12E-09** |
| 4-5-1 | 0.001042 | 1.01E-04 | 4.77E-06 | **1.39E-09** |
| 4-6-1 | 0.000498 | 1.02E-04 | 6.45E-06 | **1.07E-09** |
| 4-7-1 | 0.000209 | 1.01E-05 | **1.02E-07** | 1.00E-06 |
| 4-8-1 | 0.000329 | 1.08E-05 | 1.82E-06 | **1.81E-10** |
| 4-9-1 | 0.000109 | 8.71E-06 | 3.73E-07 | **3.93E-11** |

TABLE VI.    AVERAGE MSE OUT OF SAMPLES FOR CRUDE OIL PRICES PREDICTION

| NN Structure | ABC | GABC | Best-So-Far | GBABC |
|---|---|---|---|---|
| 4-2-1 | 3.20E-05 | 1.83E-06 | 9.99E-06 | **9.99E-08** |
| 4-3-1 | 3.10E-06 | 9.12E-05 | 9.81E-06 | **8.12E-08** |
| 4-4-1 | 2.01E-03 | 8.13E-05 | 1.92E-06 | **8.00E-08** |
| 4-5-1 | 1.10E-05 | 9.22E-06 | **1.23E-07** | 9.05E-07 |
| 4-6-1 | 1.09E-05 | 9.11E-06 | 9.78E-07 | **2.00E-08** |
| 4-7-1 | 2.76E-05 | 1.90E-06 | **1.93E-08** | 1.91E-07 |
| 4-8-1 | 3.12E-05 | 3.00E-06 | 7.92E-06 | **9.01E-09** |
| 4-9-1 | 2.01E-05 | 1.59E-06 | 9.46E-06 | **1.82E-10** |

TABLE VII.    AVERAGE ACCURACY FOR CRUDE OIL PRICES PREDICTION

| NNs Structure | ABC | GABC | Best-So-Far | GBABC |
|---|---|---|---|---|
| 4-5-1 | 90.13 | 93.72 | 95.74 | **97.98** |
| 4-6-1 | 91.01 | 92.24 | 96.92 | **98.70** |
| 4-7-1 | 90.38 | 94.81 | 96.27 | **99.15** |
| 4-8-1 | 91.63 | 94.78 | 96.31 | **99.10** |
| 4-9-1 | 92.12 | 94.31 | **98.35** | 98.17 |

TABLE VIII.    SUCCESS RATE OF ALL LEARNING ALGORITHMS FOR CRUDE OIL PRICES PREDICTION

| NN Structure | ABC | GABC | Best-So-Far | GBABC |
|---|---|---|---|---|
| 4-2-1 | 80 % | 87 % | 98  % | **100 %** |
| 4-6-1 | 95 % | 98 % | **100 %** | **100 %** |
| 4-7-1 | 89 % | **100 %** | **100 %** | **100 %** |
| 4-8-1 | **100 %** | **100 %** | **100 %** | **100 %** |
| 4-9-1 | **100 %** | **100 %** | **100 %** | **100 %** |

TABLE IX.    BEST MSE TRAINING FOR CRUDE OIL PRICES PREDICTION WHEN CS=8, D=13 AND C=1.0

| NN Structure | ABC | GABC | Best-So-Far | GBABC |
|---|---|---|---|---|
| 4-2-1 | 0.0012001 | 5.16E-03 | 8.13E-05 | **2.12E-06** |
| 4-3-1 | 0.0011915 | 1.10E-03 | 4.24E-05 | **3.15E-07** |
| 4-4-1 | 0.0010021 | 3.10E-03 | 2.50E-04 | **2.18E-06** |
| 4-5-1 | 0.0009323 | 1.01E-04 | 1.65E-05 | **1.59E-07** |
| 4-6-1 | 0.0008102 | 3.27E-04 | 3.34E-05 | **4.08E-08** |
| 4-7-1 | 0.0000961 | 1.02E-05 | **3.01E-07** | 1.16E-08 |
| 4-8-1 | 0.0001223 | 9.21E-05 | **1.10E-08** | 6.15E-08 |
| 4-9-1 | **0.0001002** | **7.12E-06** | 2.93E-06 | 8.12E-09 |

TABLE X.    BEST MSE TRAINING FOR CRUDE OIL PRICES PREDICTION WHEN CS=10, D=21 AND C=1.2

| NN Structure | ABC | GABC | Best-So-Far | GBABC |
|---|---|---|---|---|
| 4-2-1 | 0.0001281 | 2.10E-04 | 1.10E-05 | 8.72E-06 |
| 4-3-1 | 0.0013123 | 3.04E-04 | 2.84E-06 | 6.95E-07 |
| 4-4-1 | 0.0001902 | 3.05E-04 | 3.56E-06 | 5.17E-07 |
| 4-5-1 | 0.0009091 | 4.05E-04 | 3.84E-06 | 1.09E-07 |
| 4-6-1 | 0.0001230 | 2.01E-04 | 4.54E-06 | **4.08E-07** |
| 4-7-1 | 0.0001294 | **1.02E-06** | **6.01E-07** | 7.06E-07 |
| 4-8-1 | **0.0001928** | **3.23E-06** | **8.10E-07** | **3.10E-07** |
| 4-9-1 | **7.05E-05** | **7.12E-06** | **9.93E-07** | **3.10E-08** |

TABLE XI.    BEST NMSE TESTING FOR CRUDE OIL PRICES PREDICTION WHEN CS=10, D=21 AND C=1.2

| NN Structure | ABC | GABC | Best-So-Far | GBABC |
|---|---|---|---|---|
| 4-2-1 | 0.001081 | 1.16E-05 | 1.07E-06 | 1.21E-07 |
| 4-3-1 | 0.000131 | 1.61E-05 | 1.73E-06 | 1.30E-08 |
| 4-4-1 | 0.000105 | 1.09E-05 | 3.18E-06 | 3.10E-09 |
| 4-5-1 | 0.000118 | 1.01E-05 | 4.17E-06 | 1.30E-09 |
| 4-6-1 | 0.000290 | 1.02E-06 | 6.40E-07 | 1.09E-09 |
| 4-7-1 | **0.000829** | **1.01E-06** | **1.03E-07** | **1.00E-12** |
| 4-8-1 | **0.000128** | **1.08E-06** | **1.32E-08** | **1.81E-12** |
| 4-9-1 | **0.000017** | **8.71E-07** | **3.23E-08** | **3.93E-13** |

TABLE XII.    BEST OUT OF SAMPLE FOR CRUDE OIL PRICES PREDICTION WHEN CS=10, D=21 AND C=1.2

| NN Structure | ABC | GABC | Best-So-Far | GBABC |
|---|---|---|---|---|
| 4-2-1 | 1.09E-03 | 1.09E-05 | 1.09E-04 | 1.09E-06 |
| 4-3-1 | 9.12E-04 | 9.12E-05 | 9.12E-04 | 9.12E-06 |
| 4-4-1 | **9.78E-04** | **9.78E-05** | **9.78E-05** | **9.78E-06** |
| 4-5-1 | **2.00E-05** | **2.00E-06** | **2.00E-06** | **2.00E-08** |
| 4-6-1 | 2.00E-05 | 2.00E-06 | 2.00E-06 | 2.00E-08 |
| 4-7-1 | 2.36E-04 | 2.36E-07 | 2.36E-06 | **2.36E-08** |
| 4-8-1 | **9.99E-06** | **9.99E-07** | **9.99E-06** | **9.99E-10** |
| 4-9-1 | **7.12E-05** | **7.12E-07** | **7.12E-06** | **7.12E-10** |
| 4-9-1 | 2.76E-05 | **2.76E-07** | **2.76E-07** | **2.76E-11** |

TABLE XIII.    AVERAGE MAPE FOR CRUDE OIL PRICES PREDICTION

| NN Structure | ABC | GABC | Best-So-Far | GBABC |
|---|---|---|---|---|
| 4-2-1 | 0.0015 | 0.00139 | **0.000102** | **0.000107** |
| 4-3-1 | 0.00431 | 1.30E-03 | **4.31E-06** | **1.21E-07** |
| 4-4-1 | 0.00273 | 1.01E-03 | **9.10E-05** | **3.00E-08** |
| 4-5-1 | 0.00383 | 1.00E-03 | **7.01E-05** | **1.02E-09** |
| 4-6-1 | 0.00091 | 1.01E-04 | **4.35E-05** | **1.01E-09** |
| 4-7-1 | 0.00022 | 1.00E-04 | **3.00E-05** | **1.00E-09** |
| 4-8-1 | 0.00041 | 9.21E-05 | **1.00E-07** | **2.11E-10** |

TABLE XIV. AVERAGE RMSE FOR CRUDE OIL PRICES PREDICTION

| NN Structure | ABC | GABC | Best-So-Far | GBABC |
|---|---|---|---|---|
| 4-2-1 | 0.01131812 | 0.01449137 | 0.003316625 | 0.002952965 |
| 4-3-1 | 0.03622568 | 0.01743559 | 0.00168523 | **0.000833667** |
| 4-4-1 | 0.01379130 | 0.01746424 | 0.001886796 | **0.000719027** |
| 4-5-1 | 0.03015128 | 0.02012461 | 0.001959592 | **0.000330151** |
| 4-6-1 | 0.01109053 | 0.01417744 | 0.002130728 | **0.000638749** |
| 4-7-1 | 0.01137541 | **0.00100995** | **0.000775242** | 0.000840238 |
| 4-8-1 | 0.01388524 | **0.00179722** | **0.00091902** | 0.000556776 |
| 4-9-1 | **0.00839642** | 0.00266833 | 0.000996494 | **0.000176068** |

The average MAPE values obtained through the abovementioned methods for crude oil prices are shown in Table XIII. Among the various method, the proposed GBABC achieved the lowest (the best) values with more hidden nodes, while ABC and GABC achieved the highest (the worst) values in average cases. The MAPE values obtained by Best-So-Far method are also lowest when compared with ABC and GABC algorithms.

The values of RMSE obtained by ABC, GABC, Best-So-Far and proposed GBABC algorithms are mentioned in Table XIV. In term of RMSE values, the typical two methods (GABC and Best-So-Far) and particularly GBABC algorithm obtained the highest numbers of minimum error values for the prediction of crude oil prices.

The best average simulation results using all above-mentioned algorithms for crude oil price prediction are given Fig. 7 to 15, respectively. The analysis of original and predict prices has presented in convergence, training and out of sample testing figures. Different values of the dimension, hidden layer nodes C and colony size have obtained different prediction results. From the following figures, ABC and GABC are not stable in convergence speed while Best so far and GBABC have fast convergence speed.

From Fig. 7, 8 and 9, the proposed GBABC converged quickly, which signifies that the proposed method is robust, promising and effective in the prediction of crude oil prices. The ABC and GABC algorithms are failed to convergence quickly as given in Fig. 7. The Best-So-Far method also fails to convergence quickly with dimension value 41 as shown in Fig. 9, means that GABC, ABC, and Best-So-Far are not stable in convergence speed for crude oil rice prediction. The GBABC has very fast conversion speed when D=21, 35 and 41 and colony size is 7 and 9.

Using the above bio inspired learning techniques explore the best values of weights of each connection in order to reduce the training error for the crude oil prices prediction task. After repeating this process for a sufficiently large number of learning cycles the network will usually converge to any state, where the error of the calculations is small with high predicted performance. The prediction curves on 75 % dataset during training phase are given in Fig. 10, 12, 13, 14 and 15 with different D, CS and C values.

In case of ABC, GABC and Best-So-Far methods, the predicted signal are not close and stable to original; these methods did not success to predict the future crude oil prices as shown in Fig. 10, 11 and 12, respectively. However, Fig. 13

and 15 are clearly showing that crude oil prices predicted by the GBABC algorithm are very close to actual oil prices.



Fig 7. Average convergence curves of typical and proposed algorithm (CS=8 and Limit 10).



Fig 8. Average convergence curves of typical and proposed algorithm (CS=10 and Limit 20).

Fig 9.   Average convergence curves of typical and proposed algorithms (CS=10, limit 10 & 20).



Fig 11.   Best average Curves by GABC of crude oil prices prediction (limit 10 & 15).



Fig 10.   Best average prediction Curves by ABC of crude oil prices prediction (limit 10 & 15).



Fig 12.   Best average Curves by Best-So-Far ABC of crude oil prices prediction (limit 10 & 20).

Fig 13.  Best average prediction Curves by GBABC of crude oil prices (limit 20 for both).



Fig 15.  Best average prediction curves out of sample by GBABC algorithm (limit 15 & 20).

The explored weight values have used to test GBABC performance on FFNN with the rest 25% dataset of crude oil price. The best average crude oil prediction curves out of samples have presented in Fig. 14 and 15. Again, the proposed GBABC algorithm predicted prices are very close to the original crude oil prices with different FFNN structures. From the above-mentioned all tables values and figures, the performance of the proposed GBABC algorithm success to reach to minim training and testing, prediction error, fast convergence, high success rate, and high prediction accuracy for crude oil prices. Based on the above simulation results and analysis, the proposed GBABC algorithm has the capability to predict the accurate future crude oil prices. Based upon on the above-mentioned proposed bio inspired GBABC algorithm outstanding simulation, the Saudi Arabia crude oil price can be easily predicted based on the past values.

From the abovementioned tables and figures, the prediction results have been affected through the control parameters specially colony size and limit. In ABC, the best values found for the crude oil prices prediction were limit 10 and 15, CS 10 and 12 and D=21 and 39 as shown in the abovementioned figures and tables. For the proposed GBABC algorithm, the best control parameters values were limit 20, D 29 and 41 and CS 10 as mentioned in the simulation results. The proper selection of the control parameters can increase the effectiveness of the ABC, GABC, Best So far and GBABC algorithms.





Fig 14.  Best average prediction curves out of sample by ABC and GABC algorithms (limit 20 & 10).

TABLE XV.    NULL HYPOTHESIS SIGNIFICANT TEST AMONG THE ORIGINAL AND PREDICTED OIL PRICES

| Groups | Count | Sum | Average | Variance | | |
|---|---|---|---|---|---|---|
| **Actual Dataset** | 200 | 20044.5 | 100.223 | 18.1112 | | |
| **Predicted by GBABC** | 200 | 20022.8 | 100.1400 | 18.1002 | | |
| **Predicted by ABC** | 200 | 19756.9 | 98.7845 | 17.2231 | | |
| **Predicted by Best So Far** | 200 | 20001.5 | 100.0075 | 17.9151 | | |
| **Predicted by GABC** | 200 | 19826.5 | 99.1325 | 16.1102 | | |
| **ANOVA** | | | | | | |
| *Source of Variation* | *SS* | *df* | *MS* | *F* | *P-value* | *F crit* |
| **Between Groups** | 245.786 | 4 | 61.4464 | 3.38183 | 0.00928 | 2.38088 |
| **Within Groups** | 18078.7 | 995 | 18.1696 | | | |
| **Total** | 18324.5 | 999 | | | | |

TABLE XVI.   F-TEST TWO-SAMPLE FOR VARIANCES

| | **Original Prices Dataset** | **Predicted by GBABC** |
|---|---|---|
| Mean | 100.2225 | 99.98792 |
| Variance | 18.11023 | 17.98520 |
| Observations | 200 | 200 |
| df | 199 | 198 |
| F | 1.027519027 | |
| P(F<=f) one-tail | 0.494170446 | |
| F Critical one-tail | 1.263340341 | |

The results obtained by these four algorithms were tested using the Null Hypothesis Significance Testing (NHST) to know the significant similarity of the proposed and typical algorithms for crude oil prices prediction tasks. In this regards, the original dataset of the first 200 days were selected to know the critical difference. The selected dataset were compared with the results obtained by ABC, GABC, Best so far ABC and proposed GBABC algorithms. The results obtained with the 200 MCN, C1 and C2 (1.5 and 2.5), Limit 20, CS=30, D=23 and alpha=0.05 as given in the following Tables XV and XVI.

A part from MSE, NMSE, success rate and accuracy, the ANOVA test analysis of Tables XV and XVI shows the ratio of similarity/difference ratio between the original crude oil prices dataset and the predicted prices using ABC, GABC, GBABC and best so far ABC algorithms. Table XVI, where the overall original crude oil prices of the first 200 days, and the predicted crude oil prices through the typical and proposed method show that the proposed method GBABC algorithm along with ANN has a very least difference (in term of Sum, Average and Variances) in when compare to GABC, ABC and Best so far ABC algorithms. The mean and variance of the original prices and predicted prices by GBABC algorithm are much close as shown in the Table XV. Overall, the p value less than 0.5 so we reject the null hypothesis. We can then conclude that the average of the dependent variable is not the same for all methods. According to the ANOVA, if the null hypothesis is rejected, then all we know is that at least two methods:

GABC and ABC are different from each other and original as well.

Furthermore, the predicted prices obtained by Best so far methods are also close to the original crude oil prices. From Table XV, when the results obtained by the proposed GBABC algorithm analysis with original crude oil prices of the first 200 days, where again the original and predicted values are very close. Based on the T test, as shown that the P value is greater than 0.05 so the null hypothesis can't be rejected. This shows that the crude oil prices obtained by the proposed method are very close to the actual prices' dataset

## VI. CONCLUSION

The FFNN was trained and tested through the bio inspired learning methods such as ABC, GABC, Best-So-Far and GBABC and it is a computer based mathematical simulation model that predicts high accurate crude oil prices. Different FFNN structures of hidden layer nodes with different learning methods, parameters such as CS, D and C values are used as given in Table III. The proposed GBABC algorithm performed well in term of MSE, NMSE, Accuracy and success rate. The convergence speed of GBABC method is faster than typical methods. Courtesy of the strong exploration and exploitation process, the proposed method has enough decreased global trapping problem based on guided employed bees, Best-So-Far onlooker and scout bees' strategies. Therefore, the proposed method has a high accuracy, less prediction error and high success rate. Thus, the proposed GBABC is a robust and

efficient method for crude oil prices prediction. However, it is difficult to set the parameters of the proposed algorithm and ANN model appropriately. Based on high accurate predictions, the crude oil investors, analyst or producer can easily predict the volatility of the crude oil prices in different times and the environment as well. The GBABC algorithm can be extended to different application such as clustering, numerical function optimization and other complex problem solver.

REFERENCES

[1] V. Guerriero, S. Mazzoli, A. Iannace, S. Vitale, A. Carravetta, and C. Strauss, "A permeability model for naturally fractured carbonate reservoirs," Mar. Pet. Geol., vol. 40, pp. 115–134, 2013.

[2] J. Chamber, "Saudi Arabian Sector Report – Oil and Gas July 2015," Jeddah Chamber, Jeddah, Kingdom of Saudi Arabia, 2015.

[3] E. McAleavey., "Saudi Arabia oil and gas market to 2023," 2014. .

[4] T. Economics, "Saudi Arabia Crude Oil Production," vol. 2017. 2017.

[5] Argaam, "Top 10 oil producing countries in the world in 2016," 2017. .

[6] I. E. Agency, "Oil Market Report," Oil Market Report, vol. 2017. https://www.iea.org/, 2015.

[7] R. de Best, "Total OPEC crude oil production from 1998 to 2017 (in 1,000 barrels per day)," 2018. [Online]. Available: https://www.statista.com/chart/14298/black-and-gold/. [Accessed: 12-Feb-2019].

[8] A. Shabri and R. Samsudin, "Crude Oil Price Forecasting Based on Hybridizing Wavelet Multiple Linear Regression Model, Particle Swarm Optimization Techniques, and Principal Component Analysis," Sci. World J., vol. 2014, p. 8, 2014.

[9] W. Xie, L. Yu, S. Xu, and S. Wang, "A New Method for Crude Oil Price Forecasting Based on Support Vector Machines," in Computational Science – ICCS 2006: 6th International Conference, Reading, UK, May 28-31, 2006, Proceedings, Part IV, V. N. Alexandrov, G. D. van Albada, P. M. A. Sloot, and J. Dongarra, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2006, pp. 444–451.

[10] O. of the P. E. Countries, "OPEC Crude Oil Price," vol. 2017. Organization of the Petroleum Exporting Countries, 2017.

[11] index mundi, "http://www.indexmundi.com/commodities/?commodity=crude-oil&months=60." 2016.

[12] J. Hamilton, "What is an oil shock?," J. Econom., vol. 113, no. 2, pp. 363–398, 2003.

[13] P. K. Narayan, S. Sharma, W. C. Poon, and J. Westerlund, "Do oil prices predict economic growth? New global evidence," Energy Econ., vol. 41, pp. 137–146, 2014.

[14] I. Mundi, "Crude Oil Price." 2017.

[15] F. Fuels, "OPEC oil price annually 1960-2018," Statistical Report, 2019. [Online]. Available: https://www.statista.com/statistics/262858/change-in-opec-crude-oil-prices-since-1960/. [Accessed: 02-Nov-2019].

[16] M. M. Mostafa and A. A. El-Masry, "Oil price forecasting using gene expression programming and artificial neural networks," Econ. Model, vol. 54, pp. 40–53, 2016.

[17] B. Abramson and A. Finizza, "Using belief networks to forecast oil prices," Int. J. Forecast., vol. 7, no. 3, pp. 299–315, 1991.

[18] B. R. Das, S. Sahoo, C. S. Panda, and S. Patnaik, "Part of Speech Tagging in Odia Using Support Vector Machine," Procedia Comput. Sci., vol. 48, no. Supplement C, pp. 507–512, 2015.

[19] Z. Pawlak, "Rough sets," Int. J. Comput. Inf. Sci., 1982.

[20] J. J. Buckley and L. J. Jowers, "Fuzzy sets," Stud. Fuzziness Soft Comput., 2008.

[21] C. N. Babu and B. E. Reddy, "A moving-average filter based hybrid ARIMA-ANN model for forecasting time series data," Appl. Soft Comput. J., 2014.

[22] A. Khashman and N. I. Nwulu, "Intelligent prediction of crude oil price using Support Vector Machines," in 2011 IEEE 9th International Symposium on Applied Machine Intelligence and Informatics (SAMI), 2011, pp. 165–169.

[23] M. Khashei and M. Bijari, "An artificial neural network (p,d,q) model for timeseries forecasting," Expert Syst. Appl., vol. 37, no. 1, pp. 479–489, 2010.

[24] M. a Otair and W. a Salameh, "Speeding Up Back-Propagation Neural Networks," Pros. 2005 Informing Sci. IT Educ. Jt. Conf. Speeding, 2005.

[25] P. D. McNelis et al., "Ant Colony Optimization," Neural Networks Financ., vol. 2, no. 1, p. 12, Dec. 2016.

[26] A. Panakkat and H. ADELI, "NEURAL NETWORK MODELS FOR EARTHQUAKE MAGNITUDE PREDICTION USING MULTIPLE SEISMICITY INDICATORS," Int. J. Neural Syst., vol. 17, no. 01, pp. 13–33, 2007.

[27] S. Dehuri, S. Patnaik, A. Ghosh, and R. Mall, "Application of elitist multi-objective genetic algorithm for classification rule generation," Appl. Soft Comput., vol. 8, no. 1, pp. 477–487, 2008.

[28] X.-S. Yang, "Bat Algorithm and Cuckoo Search: A Tutorial," in Artificial Intelligence, Evolutionary Computing and Metaheuristics: In the Footsteps of Alan Turing, X.-S. Yang, Ed. Berlin, Heidelberg: Springer Berlin Heidelberg, 2013, pp. 421–434.

[29] K.-L. Du and M. N. S. Swamy, "Ant Colony Optimization," in Search and Optimization by Metaheuristics: Techniques and Algorithms Inspired by Nature, Cham: Springer International Publishing, 2016, pp. 191–199.

[30] X. S. Y. and S. Deb, "Cuckoo Search via Lévy flights," in World Congress on Nature & Biologically Inspired Computing (NaBIC), 2009, pp. 210–214.

[31] S. Binitha and S. Siva Sathya, "A Survey of Bio inspired Optimization Algorithms," Int. J. Soft Comput. Eng., 2012.

[32] R. Irani and R. Nasimi, "Application of artificial bee colony-based neural network in bottom hole pressure prediction in underbalanced drilling," J. Pet. Sci. Eng., 2011.

[33] S. Haykin, Neural Networks: A Comprehensive Foundation. Prentice Hall PTR, 1998.

[34] F. A. Rosenblatt, Probabilistic Model for Information Storage and Organization in the Brain. 65 Cornell Aeronautical Laboratory .

[35] Y. Shin and J. Ghosh, "The pi-sigma network: an efficient higher-order neural network for pattern classification and function approximation," in IJCNN-91-Seattle International Joint Conference on Neural Networks, 1991, vol. i, pp. 13–18 vol.1.

[36] J. W. S. Hu, Y. C. Hu, and R. R. W. Lin, "Applying neural networks to prices prediction of crude oil futures," Math. Probl. Eng., 2012.

[37] D. Sarkar, "Methods to speed up error back-propagation learning algorithm," ACM Comput. Surv., vol. 27, no. 4, pp. 519–544, 1995.

[38] H. Chiroma et al., "A Review on Artificial Intelligence Methodologies for the Forecasting of Crude Oil Price," Intell. Autom. Soft Comput., vol. 22, no. 3, p. 14, 2016.

[39] H. Naser and F. Alaali, "Can oil prices help predict US stock market returns? Evidence using a dynamic model averaging (DMA) approach," Empir. Econ., vol. 55, no. 4, pp. 1757–1777, Dec. 2018.

[40] A. Babikir and H. Mwambi, "Factor Augmented Artificial Neural Network Model," Neural Process. Lett., vol. 45, no. 2 LB-Babikir2017, pp. 507–521, 2017.

[41] H. Shah, R. Ghazali, T. Herawan, N. Khan, and M. S. M. S. Khan, Hybrid guided artificial bee colony algorithm for earthquake time series data prediction, vol. 414. .

[42] H. Garg, M. Rani, and S. P. Sharma, "An approach for analyzing the reliability of industrial systems using soft-computing based technique," Expert Syst. Appl., 2014.

[43] J. Kennedy, R. C. Eberhart, and Y. Shi, "chapter seven - The Particle Swarm," in Swarm Intelligence, San Francisco: Morgan Kaufmann, 2001, pp. 287–325.

[44] M. S and W. W., "Prediction Model for Crude Oil Price Using Artificial Neural Networks," Appl. Math. Sci., vol. 8, no. 80, p. 12, 2014.

[45] S. N. Abdullah and X. Zeng, "Machine learning approach for crude oil price prediction with Artificial Neural Networks-Quantitative (ANN-Q) model," in The 2010 International Joint Conference on Neural Networks (IJCNN), 2010, pp. 1–8.

[46] S. Kulkarni and I. Haidar;, "Forecasting Model for Crude Oil Price Using Artificial Neural Networks and Commodity Futures Prices," Int. J. Comput. Sci. Inf. Secur., vol. 1, no. 2, p. 8, 2009.

[47] M. Hamdi, C. Aloui, and S. kumar N. N. V.- 4, "Comparing Functional Link Artificial Neural Network And Multilayer Feedforward Neural Network Model To Forecast Crude Oil Prices," 2016.

[48] S. Gao and Y. Lei, "A new approach for crude oil price prediction based on stream learning," Geosci. Front., vol. 8, no. 1, pp. 183–187, 2017.

[49] X. Yin, J. Peng, and T. Tang, "Improving the forecasting accuracy of crude oil prices," Sustain., 2018.

[50] Y. Zhao, J. Li, and L. Yu, "A deep learning ensemble approach for crude oil price forecasting," Energy Econ., 2017.

[51] Y. Chen, K. He, and G. K. F. Tso, "Forecasting Crude Oil Prices: A Deep Learning based Model," in Procedia Computer Science, 2017.

[52] Z. Luo, X. Cai, K. Tanaka, T. Takiguchi, T. Kinkyo, and S. Hamori, "Can We Forecast Daily Oil Futures Prices? Experimental Evidence from Convolutional Neural Networks," J. Risk Financ. Manag., vol. 12, no. 1, p. 9, 2019.

[53] D. Karaboga, "2005. An Idea Based on Honey Bee Swarm for Numerical Optimization. Kayseri," Erciyes University, 2005.

[54] D. Karaboga, B. Akay, and C. Ozturk, "Artificial Bee Colony (ABC) Optimization Algorithm for Training Feed-Forward Neural Networks," in Modeling Decisions for Artificial Intelligence: 4th International Conference, MDAI 2007, Kitakyushu, Japan, August 16-18, 2007. Proceedings, V. Torra, Y. Narukawa, and Y. Yoshida, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2007, pp. 318–329.

[55] J. P. T. Y. Harvey Jake G. Opeña, "Automated Tomato Maturity Grading Using ABC-Trained Artificial Neural Networks," Malaysian J. Comput. Sci., vol. 30, no. 1, p. 11.

[56] D. Karaboga and B. Gorkemli, "A quick artificial bee colony (qABC) algorithm and its performance on optimization problems," Appl. Soft Comput., vol. 23, pp. 227–238, 2014.

[57] H. Shah, N. Tairan, H. Garg, R. Ghazali, and H. G. and R. G. Habib Shah, Nasser Tairan, "A Quick Gbest Guided Artificial Bee Colony algorithm for stock market prices prediction," Symmetry (Basel)., vol. 10, no. 7, p. 15, 2018.

[58] W. Gao and S. Liu, "A modified artificial bee colony algorithm," Comput. Oper. Res., 2012.

[59] G. Zhu and S. Kwong, "Gbest-guided artificial bee colony algorithm for numerical function optimization," Appl. Math. Comput., 2010.

[60] M. Tuba, N. Bacanin, and N. Stanarevic, "Guided artificial bee colony algorithm," Proceedings of the 5th European conference on European computing conference. World Scientific and Engineering Academy and Society (WSEAS), Paris, France, pp. 398–403, 2011.

[61] N. Veček, M. Mernik, and M. Črepinšek, "A chess rating system for evolutionary algorithms: A new method for the comparison and ranking of evolutionary algorithms," Inf. Sci. (Ny)., 2014.

[62] H. Shah, N. Tairan, H. Garg, and R. Ghazali, "A Quick Gbest Guided Artificial Bee Colony Algorithm for Stock Market Prices Prediction," Symmetry (Basel)., vol. 10, no. 7, p. 292, 2018.

[63] A. Banharnsakun, T. Achalakul, and B. Sirinaovakul, "The best-so-far selection in Artificial Bee Colony algorithm," Appl. Soft Comput., vol. 11, no. 2, pp. 2888–2901, 2011.

[64] H. T. Jadhav and R. Roy, "Gbest guided artificial bee colony algorithm for environmental/economic dispatch considering wind power," Expert Syst. Appl., vol. 40, no. 16, pp. 6385–6399, 2013.

[65] N. A. Husaini, R. Ghazali, N. M. Nawi, L. H. Ismail, M. M. Deris, and T. Herawan, "PI-SIGMA NEURAL NETWORK FOR A ONE-STEP-AHEAD TEMPERATURE FORECASTING," Int. J. Comput. Intell. Appl., vol. 13, no. 04, p. 1450023, Sep. 2014.

[66] T. Wu, M. Yao, and J. Yang, "Dolphin swarm algorithm," Front. Inf. Technol. Electron. Eng., vol. 17, no. 8 LB-Wu2016, pp. 717–729, 2016.

[67] H. Shah, R. Ghazali, N. M. Nawi, and M. M. Deris, Global hybrid ant bee colony algorithm for training artificial neural networks, vol. 7333 LNCS, no. PART 1. 2012.

[68] E. Hancer and D. Karaboga, "A comprehensive survey of traditional, merge-split and evolutionary approaches proposed for determination of cluster number," Swarm Evol. Comput., vol. 32, pp. 49–67, 2017.

[69] H.-B. Duan, C.-F. Xu, and Z.-H. Xing, "A hybrid artificial bee colony optimization and quantum evolutionary algorithm for continuous optimization problems.," Int. J. Neural Syst., vol. 20, no. 1, pp. 39–50, 2010.

[70] X. S. Yang and Suash Deb, "Cuckoo Search via Lévy flights," in World Congress on Nature & Biologically Inspired Computing (NaBIC), 2009, pp. 210–214.

[71] D. Karaboga and C. Ozturk, "A novel clustering approach: Artificial Bee Colony (ABC) algorithm," Appl. Soft Comput., vol. 11, no. 1, pp. 652–657, 2011.

[72] D. Karaboga and B. Akay, "A survey: Algorithms simulating bee swarm intelligence," Artif. Intell. Rev., 2009.

[73] Kabbani and B. F., "Gulf Base," Price Performance Charts, 2017. [Online]. Available: http://www.gulfbase.com/.

# Designing Model of Serious Game for Flood Safety Training

Nursyahida Mokhtar[1], Amirah Ismail[2], Zurina Muda[3]

Faculty of Information Science and Technology,
Universiti Kebangsaan Malaysia
43600 Bangi Selangor, Malaysia

*Abstract*—Serious games have the potential to increase motivation among users in the aspect of safety training. Additionally, serious games can also positively impact training outcomes when the knowledge and skills acquired during a serious game training are transferred to a real-world application. The development of a serious game is based on the game elements and theories determined according to the goals or objectives of the game being developed. There are existing serious games that have been used for training but less usage of scenarios and feedback element render the game less effective for training purposes. Besides that, existing serious games for training purposes fail in delivering domain content to achieve the game objectives since they are more focused on entertainment. This is because the games do not involve experts in providing game domain content. The objective of this paper is to design a serious game model for flood safety training. Preliminary study and literature review are used in this study as the research method. The result of this study is a model of a serious game for flood safety training. In conclusion, this study focuses on the design of a serious game model for flood safety training that includes the elements of serious game identified and adapted to psychology readiness based on the flood training module by Malaysian Defense Force (APM). This makes the serious game more attractive and can give intrinsic motivation to players. For future studies, every single element serious game and theory of psychology readiness in the model developed in this study will be validated with the expert game and expert psychology.

*Keywords*—*Serious game model; flood safety training; flood awareness; intrinsic motivation; psychology readiness*

## I. INTRODUCTION

The rapid growth of multimedia technologies over the last 20 years means that today's children and young adults were born in a computerized world and are used to handling all kinds of software products and games [1][2]. A serious game is one the multimedia technologies [3][4] defined as an interactive computer application with or without significant hardware components, that has a challenging goal, fun to play and engaging, incorporates some scoring mechanism, and supplies the user with skills, knowledge, or attitudes that are useful in reality.

One of the uses of serious games is for training purposes. The main purpose in a training context is to help players to reach their learning objectives and provide an alternative method of training to gain and maintain skills. Besides that, serious games for training aim to keep the interest in training activities high by keeping players entertained, yet interest and

motivation often wane over games. There are many domains in which people need to invest their time and effort in a training activity to see future benefits. One of the instances where serious games have been used is in natural disaster training. Natural disaster is a phenomenon which occurs all around the world in various ways such as floods.

Floods are a natural disaster of concern to society as floods negatively impact humans and infrastructures [5]. When a flood occurs, besides physical readiness, psychological readiness is also important for monitoring and adjusting individual responses.

According to [6], the awareness level of civilians on the importance of preparing for floods is still low. This statement is supported by the APM. Although various methods have been conducted by APM to prepare the civilians physically or emotionally, they still do not show a positive response to do flood preparation. APM requires a more interactive medium of disseminating necessary information. This is because existing mediums such as flyer distribution, television advertisement, APM official social media and video on the official APM website do not give a positive impression to the civilians.

Therefore, serious games are suitable to disseminate information [7] and can be used as a medium for training because serious games allow users to experience situations that are difficult to achieve in reality due to factors such as cost, time and safety.

Serious games are already being used for training purposes in real-life [8]; however, they do not gain much popularity as less usage of scenario and feedback element are integrated into the games, making the games less effective for training. Therefore, serious games must be developed in such a way that it can show more than one scenario to generate feedbacks more often. Because the role of feedback in serious games for training is important as it can give intrinsic motivation to players to keep playing [9]. In addition, existing serious games fail to deliver domain content to achieve game objectives. But, are more focused on entertainment [10][11]. This happens because it does not involve expert in providing game domain content.

The objective of this study is to design a serious game model for flood safety training which has three flood situation scenarios equipped with elements and psychological readiness theory that can provide intrinsic motivation to players. Besides, this game also involves experts to provide game domain

content. This serious game helps the APM to convey information on flood safety based on flood training modules to the civilians. Furthermore, this serious game can give the awareness to the civilians of the importance of preparing for flood situations.

A serious game for flood safety training focuses on training objectives in the context of games that resemble real-world scenarios. Players are briefed on the real flood situations to explain the actions that should be taken when floods indeed occur [12]. Therefore, the tasks in the game are designed according to the appropriateness of the flood training modules to help players adapt to the game environment. This study is divided into several sections including introduction, research method, research findings, design model, discussion, conclusions and future works.

## II. METHOD

This qualitative research uses two research methods: i) Preliminary study involving interview session conducted at the Malaysia Civil Defence Force Training Academy (ALPHA) in Bangi, Selangor, with six members of APM to obtain data, and ii) literature review as shown in Fig. 1.



Fig. 1. Method for design a serious game model for flood safety training.

Preliminary study was conducted using the interview method to obtain data analysis of flood awareness training preparation among civilians, the existing methods of dissemination of flood awareness information and the need for technology requirements to carry out the flood awareness training preparation from the experts' perspective.

Literature review was conducted to identify the elements of a serious game for training in various domains used in previous studies and then to identify the appropriate elements that can provide intrinsic motivation to the players which can be applied into the serious game for flood safety training. Additionally, literature review was done to identify the suitable theories and approaches to be applied in flood preparedness training.

### A. Preliminary Study

The preliminary study for this research has three phases. Phase 1 is setup to identify the objective of the preliminary study, preparation of the instrument, and to identify the informants. Phase 2 is data collection which was carried out by

conducting interview sessions with the informants while Phase 3 is data analysis to analyse the information gathered from the interview. The content of the data analysis for the interview method has been discussed in the preliminary study of flood awareness training preparation using serious games [13].

### B. Literature Review

The literature review method which is adapted from [14] has four phases. Phase 1 identifies the objectives of the literature review to retrieve articles relevant to this research. Phase 2 analyses the retrieved articles. Phase 3 determines the quality of the articles analysed by ensuring that articles are downloaded from reliable databases. Finally, Phase 4 discusses the findings before making conclusion.

*1) Identifying the objectives and downloading related articles*: In this phase, the primary purpose of literature review was identified to ensure the title and contents of the journal articles selection conform to the topic of research. The articles were selected based on two methods. In the first method, the articles were identified using the search of title, domain keywords, and abstract, where selected articles were from years 2009-2018.

The keywords used were "Serious Games", "Serious Game Training", "Disaster", "Flood", "Psychology Readiness" and "Intrinsic Motivation". The databases used for downloading articles were Google Scholar, IEEE Xplore, ACM Digital Library, and Science Direct.

The selected articles must be written in English. According to [15], the first search method involved articles focusing specifically on the title, known as the 'primary' articles. The second method involved articles that were not related to the title; thus, the search needed to be wider and known as the secondary source. For the secondary source, the article titles did not discuss the domain but there were related contents. The search was done with a quick reading of the content of the articles and then focused on the relevant contents.

Meanwhile, another approach used is called snowballing [16] (i.e. retrieving relevant papers based on target papers' references list or paper citing). Articles downloaded in the form of PDF files were structured in the Mendeley program to facilitate the arrangement of references. From the 100 related articles downloaded, 24 studies were related to the elements of serious game for training, one study was related to readiness psychology, and eight studies were related to intrinsic motivation.

*2) Analysing related articles*: The criteria and specifications of the downloaded articles were analysed. First, articles that stated the elements of game were selected. Next, the elements that can give intrinsic motivation to players were identified. Finally, articles related to the psychology readiness of the society related to flood are identified. According to [15], the critical aspects observed during the analysis of literature studies include definitions, objectives, characteristics, historical analysis, success factors, failure/problem factors, research methods, theory, further studies, and contents.

*3) Evaluating the quality of study*: In this evaluation phase, articles that were not documented in detail are not selected. The selected articles must also be retrieved from a reliable database.

*4) Discussions and making conclusions*: The result of the review method is the elements of serious game for training, elements that were identified to provide intrinsic motivation to users, and psychology readiness. The detail result will be discussed further in the next section.

## III. RESULT

In this section, the results are discussed details based on the findings from the preliminary study, literature review. Next, as a final result, result from both of these two methods is combined to produce a design model of serious game for flood safety training.

### A. Result Systematic Review

The result of preliminary study has been discussed in the previous paper preliminary study of flood awareness training preparation using serious games [13].

Furthermore, APM advised that game need to be developed based on the APM flood safety training module to deliver safety information to the public. The flood safety training module is provided by the APM headquarters. This module is used as a reference by the ALPHA to train volunteers and APM staff who undergo training at ALPHA. The module describes the safety measures in dealing with a flood situation in a well-planned manner. In addition, this module briefly explains the actions the civilians should take and the role of APM in a flood disaster. The flood safety training module is divided into three flood situations namely before, during and after the flood. The module aims to provide a complete guideline on the measures that should be taken to avoid loss of life and property.

### B. Result Literature Review

The result of the literature review method is the selected element of serious game for training and theory of psychology readiness.

*1) Element of serious game*: Total of 24 previous studies on serious games for training have been analyzed. From the 24 studies, there are four domains involved, namely natural disasters, medical, safety, and education. Seven studies in the domain of natural disasters, eight studies in the medical domain, six studies in safety domains and three studies in the education domain were done, giving a total of 24 studies. From the analysis of the 24 studies, a total of 26 serious game elements for training were identified. Elements of serious game for training as listed in Table I.

The elements obtained were analyzed to identify the appropriate elements for the game to be developed in this serious game. For objective purpose, there must be elements of the game that can enhance the intrinsic motivation of the players in order to attract players to continue playing.

TABLE I. ELEMENT OF SERIOUS GAME FOR TRAINING

| Bil | Element / Domain | Disaster | Medical | Safety | Education | Total |
|-----|------------------|----------|---------|--------|-----------|-------|
| 1 | Scenarios | 7 | 6 | 4 | 0 | 17 |
| 2 | Interactive | 5 | 4 | 5 | 2 | 16 |
| 3 | Immersive | 5 | 5 | 5 | 0 | 15 |
| 4 | Goal | 3 | 6 | 4 | 1 | 14 |
| 5 | Feedback | 2 | 6 | 3 | 2 | 13 |
| 6 | Challenges | 4 | 3 | 1 | 3 | 11 |
| 7 | Rules | 4 | 3 | 2 | 0 | 9 |
| 8 | Involvement | 4 | 2 | 3 | 0 | 9 |
| 9 | Storyline | 2 | 1 | 2 | 1 | 6 |
| 10 | Reward | 1 | 2 | 0 | 1 | 4 |
| 11 | Players objective | 2 | 0 | 2 | 0 | 4 |
| 12 | Game Outcome | 1 | 0 | 2 | 0 | 3 |
| 13 | Enjoyment | 1 | 0 | 1 | 0 | 2 |
| 14 | Game mechanics | 0 | 0 | 2 | 0 | 2 |
| 15 | Probability | 1 | 0 | 1 | 0 | 2 |
| 16 | Game Procedures | 1 | 0 | 0 | 0 | 1 |
| 17 | Resources | 1 | 0 | 0 | 0 | 1 |
| 18 | Conflict | 1 | 0 | 0 | 0 | 1 |
| 19 | Boundary | 1 | 0 | 0 | 0 | 1 |
| 20 | Game Play | 1 | 1 | 0 | 0 | 1 |
| 21 | Interaction Modes | 0 | 1 | 0 | 0 | 1 |
| 22 | Frame story | 0 | 1 | 0 | 0 | 1 |
| 23 | Focus attention | 0 | 1 | 0 | 0 | 1 |
| 24 | Playability | 0 | 1 | 0 | 0 | 1 |
| 25 | Game fullness, | 0 | 1 | 0 | 0 | 1 |
| 26 | A sense of control | 0 | 1 | 0 | 0 | 1 |

According to [17], intrinsic motivation can encourage user engagement with games. When users engage with the game, they start to feel attracted and will continue to repeat the game until it can influence their behavior [18]. Hence, this serious game can give users an awareness of the importance of flood preparation training.

Based on the analysis, elements that were identified to provide intrinsic motivation to users are the elements of scenario, feedback, challenge, reward and enjoyment. Table II shows the elements that identified.

TABLE II.     ELEMENT SRIOUS GAME FOR FLOOD SAFETY TRAINING

| Bil | Elemen / Domain | Disaster | Medical | Safety | Education | Total |
|---|---|---|---|---|---|---|
| 1 | Scenarios | 7 | 6 | 4 | 0 | 17 |
| 2 | Feedback | 2 | 6 | 3 | 2 | 13 |
| 3 | Challenges | 4 | 3 | 1 | 3 | 11 |
| 4 | Reward | 1 | 2 | 0 | 1 | 4 |
| 5 | Enjoyment | 1 | 0 | 1 | 0 | 2 |

*a) Scenarios*

From analysis of the 24 studies, the usage of scenario element is the highest in serious games for training as it was used in 17 studies. From the 17 studies, seven studies were in the domain of natural disasters, six studies in the medical domain, and four studies in the safety domain. However, the scenario element in the education domain was not present. This shows the scenario element is important to describe the situation in a game [19]. All the 17 studies just use one senario. According to [8], by using scenario elements, game progressions can be seen from time to time through events that occur during a game play. In a serious game for flood safety training, there are three flood scenarios, namely before, during and after flood. With different training scenarios, it can create resilient trainers because they can respond to different situations [20]. The events that occur in these scenarios will explain the use of the elements of feedback, challenges, rewards and enjoyment.

*b) Feedback*

A total of 13 studies used feedback element, implying the element is widely used in serious games for training. From the 13 studies, two studies used feedback element in the domain of natural disasters while six studies used the element in medical domain. Meanwhile, three studies represented the safety domain and two studies represented the educational domain. From this analysis, it was found that the feedback element is important in a serious game for training since the frequency of its usage appeared in 13 studies. Nevertheless, this element is less used in the domain of natural disasters. According to [9] and [21], the role of the feedback element in serious games for training is important as it can give intrinsic motivation to players to continue playing. In a serious game for flood safety training, the use of feedback element is to show game feedbacks on actions taken by player.

*c) Challenges*

Element of challenge is also an element that is used often in studies of serious game for training. 11 out of 24 studies used the challenge element with four studies in the domain of natural disasters, three in the medical domain, one in the safety domain and three in the education domain. According to [22], playing a challenging game can meet the needs of intrinsic motivation of players in terms of efficiency and autonomy. Challenge, in this context, is defined as the challenge of how players make decision to take necessary actions. If the player erred in decision making, the action taken will result in them being unable to continue the game. Therefore, the element of challenge is important in serious games of flood safety training. This is to inform the player of the actions that should and should not be taken during floods.

*d) Reward*

Element of reward was found to be less used in serious games for training. Only 4 studies used this element with one study representing the domain of natural disasters, two representing the medical domain and one representing the education domain. However, there was no application of reward element in the safety domain. Although this element is less applied in previous serious games for training, but the reward element should be used in serious games to encourage user to continue playing. According to [23], to attract users to keep playing, users need to be rewarded to increase the intrinsic motivation of users.

*e) Enjoyment*

For element of enjoyment is the least used in serious games for training. Only two out of all 24 studies applied this element, with each one study in the domain of natural disasters and the domain of education. In the medical and safety domains, the enjoyment element was not used. Even though the enjoyment element is less applied in previous studies, this element needs to be applied to give users a fun time while playing. According to [24], an exciting game provides a platform for players to experience fun and can lead to increased intrinsic motivation. The enjoyment element can be applied through the use of music or sound effects that can enable users to have fun and enjoy the game until the end.

*2) Psychology readiness*: The study was adapted and modified from a study by [25] in relation to the community's preparedness from a psychological perspective related to flood disaster. This study emphasized on the community preparedness to cope with the disaster. From the perspective of psychology, psychology readiness regarding preparation is important in discussing the level of readiness [26]. Psychology readiness explains the need for internal and external interaction of human preparation. Psychology readiness are based on three theory, namely cognitive theory, affective theory and psychomotor theory. All these theories influence each other and produce human behavior towards something as it happens in the issue of human readiness for flood disaster.

*C. Combination Result of the Two Study Methods*

The results of the two study methods are combined to produce a model of serious game for flood safety training that can provide intrinsic motivation to the players. Results are shown in Table III.

TABLE III.     RESULT FROM TWO RESEARCH METHOD OF THE STUDY

| No | Preliminary Study | Literature review | |
|---|---|---|---|
| | | Element Game | Psychology Readiness |
| | Awareness of flood training preparation among civilians | Scenario | Cognitive Affective Psychomotor |
| | Campaign Used for Flood Training Preparation | Feedback Challenge | Cognitive |
| 3 | Existing Method of Dissemination of Flood Awareness Information | Reward Enjoyment | Affective |
| 4 | Technology Requirements for Flood Awareness Training Preparation | Game Features | Psychomotor |

*1) Awareness of flood training preparation among civilians (Scenario)*: Since the level of awareness among the civilians about the importance of preparing for flood is still low, flood preparedness exercises in terms of physical and emotional aspects needs to be given to them. To make this training more effective, a technology is required to help the APM convey safety information to the civilians so that they can raise their awareness. Therefore, it is proposed that a serious game is the technology that can help the APM convey information about flood safety. The proposed serious game can provide awareness to the civilians of the importance of preparing for flood situations. For that purpose, the serious game must have a scenario element that can describe the actual flood situation. This game has to be developed based on the APM flood training module. Using the module, a simple storyline is built. There are three flood scenarios in this serious game: before flood, during flood and after flood. The sequence of scenarios is in the order of actual flood events. Each scenario has different training activities and focuses on important tasks in training. Besides that, cognitive theory, affective theory, and psychomotor theory are applied in scenario elements to illustrate the decision, feeling and action of the players when floods occur.

*2) Campaign used for flood training preparation (feedback and challenge)*: The APM has organized campaigns in areas affected by flood disasters. However, all the campaigns are still less effective to civilians because they do not seem to have awareness in preparing for the flood disaster. Using serious game, information can be delivered more attractively and effectively. The use of feedback elements can deliver information directly to players in a style that can attract players to know the information they want to convey. When a player takes an action in the game, they will get a feedback from the game about the actions taken. Through this feedback, players will be given information related to floods. Additionally, the challenge element is also important to help the civilians to make decisions during emergencies. Through the use of challenge element, players can be more aware of their actions in whether the actions are right or wrong in a flood situation. Cognitive theory are applied in feedback and challenge element when players must make decisions.

*3) Existing method of dissemination of flood awareness information (reward and enjoyment)*: The existing method needs to be upgraded using an effective technology so that it can stimulate the minds of the civilians to better understand the information delivered by APM. The existing medium to disseminate information is less effective because it does not attract interest of civilians. The use of serious game can attract the interest of the civilians. The use of reward and enjoyment elements can motivate players to continue playing the game. The affective theory is applied in reward and enjoyment element when the player feels happy to get marks.

*4) Technology requirements for flood awareness training preparation (game features)*: Based on the suggestions on game features that APM has provided, a serious game with a role play game feature will be developed. Role play game is a computer game in which the players assume the roles of characters or take control of an avatar in a fictional setting. According to [27] role-playing games (RPGs) are digital games that strongly emphasize narrative, alternating action with episodes of exploration and dialogue, and with intricate reward mechanisms.

Among the game components available in this game is a scoring system to reward players [28] and using a pop-up box to give directions and information to players. Besides that, players can print out important information about flood safety. The game's user interface is also interesting and easy to understand as it displays a background illustrating the situation of the flood. The serious game also contains interactive multimedia elements as requested by the APM such as text, graphics, animation, and audio. Psychomotor theory can be seen throughout the game process.

Overall, the acquired elements were matched with the preliminary study result and were used for designing the game model. The description of the elements chosen for the game model is explained in the next sub-topic, which is the design model.

## IV. MODEL DESIGN

Five game elements are applied in this game of scenario, feedback, challenge, reward, and enjoyment. Other than element of serious game, psychology readiness theory is also involved in developing this serious game for preparation in terms of player emotion which is cognitive theory, affective theory and psychomotor theory. Each of the element of serious game and psychology readiness which is used elaborated further through the design of the model in Fig. 2.



Fig. 2. Serious Game Model for Flood Safety Training.

### A. Game Element

*1) Scenario*: The Scenario in a game is defined as the "context" in which the game takes place [29], the desired game development over time, and some events that occur when playing the game to improve the training. The general characteristic of a scenario is to describe the chronology of events. Each scenario has the goal of providing a lesson to the player. According to, scenarios in serious games must involve experts and follow the predetermined content of the study

domain so that the objective of the game can be conveyed to players. Therefore, in this serious game, there are three scenarios describe the situation of the flood: scenario before flood, scenario during flood and scenario after flood. The three scenarios are designed based on the module derived from the ALPHA authority as a study reference.

*2) Feedback*: Feedback is the information received by the player base on the action taken [30]. Feedback occurs when a player performs a task [31]. Through the feedback received, the player can find out whether the action taken is correct or wrong. Furthermore, through the feedback received, the user can assess whether the task can be solved successfully or not [21]. According to [32], feedback is necessary to create interactions between the players and the game. Therefore, players need to understand the characteristics of a game, so that players can understand how the game gives feedback. Feedbacks can be provided in the form of visual and audio.

This game uses a type of visual feedback known as explanatory feedback. This feedback is also known as corrective feedback. According to [27], corrective feedback allows players to get information from the game about their game performance in terms of whether the action they took is right or wrong. In other words, this feedback only tells the user whether the action was taken correctly or wrongly. But explanatory feedback occurs when players take action in the game and receive feedback that clearly explains their performance quality. Explanatory feedback will show the answers by explaining why the answer is correct or wrong [33].

In the serious game for flood safety training, explanatory feedback is selected to give an explanation to the player in a visual form regarding the actions taken by the player either correct or incorrect when in a flood situation.

*3) Challenges*: Challenge is a task in a game that gives players a problem to continue the game [34]. A challenge occurs when a problem arises, the player can plan and choose different strategies to solve the problem. Each player has different abilities to solve problems in each game. Challenges in a game will test the player's skill level to make a choice. Challenges in every game are important to make the game more interesting. Additionally, this element makes the player not bored and not to give up easily. In this game, the challenge element in decision making is applied to allow the player to decide on the correct action that should be taken during an emergency. According to [35] [36]decision making is an important challenge in any game because it's difficult to predict how a person will react in an emergency crisis due to many factors involved in decision making.

*4) Reward*: The reward is a game element that gives users a sense of satisfaction and encourages them to continue playing to achieve more rewards [37]. According to [38], to attract the players to continue playing, players need to be rewarded. This incentive aims to increase the intrinsic motivation of players to keep playing. In this game model, rewards are given in the form of scores when players can complete each given task properly.

*5) Enjoyment*: The primary emotion of playing a game is enjoyment [39]. Enjoyment can be classified as an attitude towards experiences in entertainment, complete with cognitive and affective psychology [40]. The element of enjoyment is an important component that keeps players entertained while playing games. Without the enjoyment element in games, players can easily feel bored [24] and will be less interested in continuing the game. Interesting games provide platforms for players to experience fun and can lead to an increase in intrinsic motivation. Therefore, it is the responsibility of game designers to develop games that can give players a sense of enjoyment [39]. In order to make a game more enjoyable, there should be new additions to the game [38].

In this serious game, the element of enjoyment is highlighted in every flood scenario. The transition from one scenario to another with the addition of new things allows players to enjoy more and keep playing the game. Not only that, the enjoyment element is applied in this game through audio. Appropriate background music and sound effects will be applied in this game. For example, if a player scores well, a happy sound effect that can raise the player's spirits to keep playing will be played.

*B. Psychology Readiness*

*1) Cognitive theory:* Cognitive theory involves an individual's mental readiness to understand, think and reason [41] (make judgments and evaluations using intellect or logic) in any situation. When in an emergency, an individual's cognitive performance will be disturbed. The individual will make an unreasonable decision, and consequently the opportunity to live diminishes.

In the serious game for flood safety training, cognitive theory is highlighted when players face challenges to make decisions. For example in the scenario after the floods, the player enters the toilet to check the condition but there is a venomous animal (snake). In this instance, the player uses his cognitive ability to make the right decision and does not endanger himself. The player needs to decide whether to catch the snake himself or to contact 999. If the player decides to catch the snake himself, the player will die of a snake bite (game ends). Conversely, if a player calls 999, APM members will arrive to capture the snake and the player will score. In this situation, the player needs to think carefully to make the right decision because the decision made will influence the decision of the game.

*2) Affective theory*: Affective theory involves feelings [39][42]. This condition can be formed in a flood victim when faced with floods. In a serious game of flood safety training, there are two feelings of a player that show affective theory which are the feelings of fear and joy. Fear comes in the scenario of during floods wherein the players are on the way to the relief center and face anxiety when flood water starts to rise dramatically. In this situation, players are anxious to ask for help. The second scene is a venomous animal scene in a scenario of after floods. The affective theory in both scenes is highlighted through the use of sound effects. The sound

effects used can make the player feel the anxiety. With this method, the game becomes more interesting. On the other hand, a happy feeling happens when players get scores. The affective theory of happy feeling is shown using a sound effect that indicates the increase of scores. This can give players a sense of excitement and make them more eager to continue the game.

*3) Psychomotor theory*: Psychomotor theory is the potential of physical maturity or preparation and cooperation to carry out a work. Psychomotor theory is described through physical action taken by players based on cognitive and affective theories. High psychomotor preparation enables an individual to act efficiently and effectively. This psychomotor theory can be seen throughout the game process. As long as the user plays the game, the player performs training and preparations for the flood.

*C. Intrinsic Motivation*

According to [43], human motivation refers to one's inspiration to act. Therefore, users need to be motivated to play a game to allow them to act according to the requirements of the game. One of the motivations that show a user really enjoys the game is intrinsic motivation [38]. Intrinsic motivation exists in individuals when doing activities [40] that can provide satisfaction to themselves and meet individual psychological needs naturally. In addition, according to [44], intrinsic motivation is present when the user strives to pursue the game because the game is interesting and fun. Digital games can give intrinsic motivation to players to continue playing [17] because these games can give them an excitement feeling. When players have intrinsic motivation, it makes them feel more positive to continue playing in the present time and in the future [17][45].

## V. DISCUSSION

The preliminary study method is used to identify the problem and the domain needs. While literature review is used to identify problem regarding preview serious game in training and identify a solution. Based on the result of both methods, it used to design serious game for flood safety training.

Based on preliminary study result, flood preparation is important to determine the civilians' ability to experience in this situation. But the civilians are still unaware of the need to prepare for flood crisis. In the era of rapidly developing technology, people are more likely to receive information and conduct training in the form of interactive technology to enable them to interact with the technology. Because of that APM need a new technology which can attract civilians to do a preparation.

A serious game can be useful as a training tool for people who have to act in emergencies [46][47]. Therefore, preparations in the form of safety training using serious game were developed to train the civilians so that they will know the necessary measures that must be taken if floods occur. Serious games in training are could contribute to discipline human behaviors. The innovation of these serious games brings motivations, encourages good moral values and natures positive responses about flood preparation. It also can provide

players with experience and making them enjoyable in a safe and reasonable situation these enable them to gain the knowledge, skills, and competencies that can be applied to life [48].

Based on literature review result serious game for flood safety training needs to develop more than one scenario in such a way can generate feedbacks more often. Furthermore serious games with other element such as challenges, reward, enjoyment making the serious game more interesting, and challenging tasks or quest, can influence the user experience of players and encourage them to further explore the game.

Other than preparation in physical training, emotional preparations should also be applied. Therefore, a serious game for flood safety training is carried out physically and emotionally. Through physical exercises, the emotions of civilians can be trained at the same time [49]. Emotions are related to an individual's psychology and mood that can be shown in physical behavior.

## VI. CONCLUSIONS AND FUTURE WORKS

In conclusion, to attract civilians to do preparation training, it requires an interactive approach. Therefore serious games have great potential to be used for training. Serious games aim to improve training processes by providing attractive, motivating and effective tools that may also create positive situations among trainer and civilians. Training activities utilizing serious games in 3D animation environments are used increasingly to create training scenarios.

The use of more than one scenario to allow trainees to learn techniques for various stressors which can help build awareness. This study focuses on the design of a serious game model for flood safety training that can be useful as a training tool for people who have to act in emergencies. The design model includes the elements of a serious game that have been identified and adapted to readiness psychology and will be used to develop a serious game based on the flood training modules set by the Civil Defence Force of Malaysia (APM). This makes the serious game to be more attractive and can give intrinsic motivation to players to keep playing. Limitation of this study focuses on element game for training only, not the overall element in the serious game. Besides that the content of flood management is from Malaysia only. For future studies, every single element serious game and theory of psychology readiness in the model developed in this study will be validated with the expert game and expert psychology using Inter-Rater Reliability method.

### REFERENCES

[1] J. Ecalle, A. Magnan, and U. L. Lyon, "Serious games as new educational tools : How effective are they ? A meta-analysis of recent studies," J. Comput. Assist. Learn., no. June, 2013.

[2] N. A. M. Zin and W. S. Yue, "History educational games design," Proc. 2009 Int. Conf. Electr. Eng. Informatics, ICEEI 2009, vol. 1, no. July, pp. 269–275, 2009.

[3] S. Raihan, Z. Abidin, S. Fadzilah, M. Noor, and N. S. Ashaari, "Guidelines of Brain-based Learning through Serious Game for Slow Reader Students," 978-1-5386-0475-5/17/$31.00 ©2017 IEEE, 2017.

[4] N. A. M. Zin and W. S. Yue, "History educational games design," Proc. 2009 Int. Conf. Electr. Eng. Informatics, ICEEI 2009, vol. 1, no. July, pp. 269–275, 2009.

[5] V. Anindhita and D. P. Lestari, "Designing Interaction for Deaf Youths by Using User-centered Design Approach," 2016.

[6] E. A. A. R. Asiah Sarji, Fahmi Mahamood, Hirwan Jasbir, Norrafidah, "Musibah banjir. menghadapinya gaya perlis," Procedia Semin. R. Inst., no. SIRaj II, pp. 1–11, 2014.

[7] D. Sutherland and R. Dennick, "Exploring culture, language and the perception of the nature of science," Int. J. Sci. Educ., vol. 24, no. 1, pp. 1–25, 2002.

[8] A. J. Q. Tan, C. C. S. Lee, P. Y. Lin, S. Cooper, L. S. T. Lau, W. L. Chua, and S. Y. Liaw, "Designing and evaluating the effectiveness of a serious game for safe administration of blood transfusion: A randomized controlled trial," Nurse Educ. Today, vol. 55, no. April, pp. 38–44, 2017.

[9] C. Burgers, A. Eden, M. D. Van Engelenburg, and S. Buningh, "Computers in Human Behavior How feedback boosts motivation and play in a brain-training game," Comput. Human Behav., vol. 48, pp. 94–103, 2015.

[10] S. Luz, M. Masoodian, R. R. Cesario, and M. Cesario, "Using a serious game to promote community-based awareness and prevention of neglected tropical diseases," Entertain. Comput., vol. 15, pp. 43–55, 2016.

[11] G. Rebolledo-Mendez, K. Avramides, S. de Freitas, and K. Memarzia, "Societal impact of a serious game on raising public awareness," Proc. 2009 ACM SIGGRAPH Symp. Video Games - Sandbox '09, p. 15, 2009.

[12] J. Tixier, A. Dandrieux, and P. Slangen, "Training decision-makers : Existing strategies for natural and technological crisis management and specifications of an improved simulation-based tool," Saf. Sci., 2016.

[13] N. Mokhtar, A. Ismail, and Z. Muda, "Preliminary Study : Flood Awareness Training Preparation Using Serious Games," Asia-Pacific J. Inf. Technol. Multimed. J. Teknol. Mklm. dan Multimed. Asia-Pasifik Vol. 7 No. 2-2, December 2018 13 - 26 e-ISSN 2289-2192, vol. 7, no. 2, pp. 13–26, 2018.

[14] Ng, K. H., Bakri, A. and Abdul Rahman, A., Effects of persuasive designed courseware on children with learning difficulties in learning Malay language subject, Journal Of Information Systems Research And Innovation 2015, 56–65. doi:10.1007/s10639-015-9391-7.

[15] W. Bandara, S. Miskon, and E. Fielt, "A Systematic, Tool-Supported Method For Conducting Literature Reviews In Information Systems," in European Conference on Information Systems (ECIS), 2011, p. 221.

[16] C. Wohlin, "Guidelines for Snowballing in Systematic Literature Studies and a Replication in Software Engineering," Proc. 18th Int. Conf. Eval. Assess. Softw. Eng. 1-10., 2014.

[17] R. M. Ryan, C. S. Rigby, and A. Przybylski, "The motivational pull of video games: A self-determination theory approach," Motiv. Emot., vol. 30, no. 4, pp. 347–363, 2006.

[18] E. Uhlmann and J. Swanson, "Exposure to violent video games increases automatic aggressiveness," J. Adolesc., vol. 27, no. 1, pp. 41–52, 2004.

[19] E. Prasolova-førland, M. Fominykh, and A. I. Mørch, "Training Cultural Awareness in Military Operations in a Virtual Afghan Village : A Methodology for Scenario Development," 2013 46th Hawaii Int. Conf. Syst. Sci. Train., pp. 903–912, 2013.

[20] B. Kolen, B. Thonus, K. M. Zuilekom, and E. De Romph, "Evacuation a serious game for preparation," 2011 Int. Conf. Networking, Sens. Control. ICNSC 2011, no. April, pp. 317–322, 2011.

[21] P.-H. Tan, S. Ling, and C. Ting, "Adaptive digital game-based learning framework," DIMEA '07 Proc. 2nd Int. Conf. Digit. Interact. media Entertain. arts, no. 1, pp. 142–146, 2007.

[22] A. K. Przybylski, R. M. Ryan, and C. S. Rigby, "The motivating role of violence in video games," Personal. Soc. Psychol. Bull., vol. 35, no. 2, pp. 243–259, 2009.

[23] M. V Birk, R. L. Mandryk, and C. Atkins, "The Motivational Push of Games : The Interplay of Intrinsic Motivation and External Rewards in Games for Training," CHI Play '16, Oct. 16-19, 2016, Austin, TX, USA © 2016 ACM. ISBN 978-1-4503-4456-2/16/10…$15.00 DOI http//dx.doi.org/10.1145/2967934.2968091, 2016.

[24] B. Kang and S. H. Tan, "Interactive Games: Intrinsic and Extrinsic Motivation, Achievement, and Satisfaction," J. Manag. Strateg., vol. 5, no. 4, pp. 110–116, 2014.

[25] T. Pah and R. Syed, "Kesediaan diri anggota masyarakat daripada perspektif psikologi berkaitan bencana banjir di daerah segamat," J. techno Soc. issue Soc. Gov., vol. 7, no. No 2, pp. 1–24, 2015.

[26] J. C. Turner and P. J. Oakes, "The significance of the social identity concept for social psychology with reference to individualism, interactionism and social influence," Br. J. Soc. Psychol., vol. 25, no. 3, pp. 237–252, 1986.

[27] F. Cornillie, G. Clarebout, and P. Desmet, "The role of feedback in foreign language learning through digital role playing games," Procedia - Soc. Behav. Sci., vol. 34, pp. 49–53, 2012.

[28] H. Wang, "Game Reward Systems : Gaming Experiences and Social Meanings," Proc. DiGRA 2011 Conf. Think Des. Play., no. March, 2016.

[29] C. Hartog, "Scenario design for serious gaming," no. January, 2009.

[30] G. S. Alder, "Examining the relationship between feedback and performance in a monitored environment : A clarification and extension of feedback intervention theory," J. High Technol. Manag. Res. 17 157–174 Examining, vol. 17, pp. 157–174, 2007.

[31] S. R. Serge, H. A. Priest, P. J. Durlach, and C. I. Johnson, "The effects of static and adaptive performance feedback in game-based training," Comput. Human Behav., vol. 29, no. 3, pp. 1150–1158, 2013.

[32] M. Schmierbach, Q. Xu, A. Oeldorf-Hirsch, and F. E. Dardis, "Electronic Friend or Virtual Foe: Exploring the Role of Competitive and Cooperative Multiplayer Video Game Modes in Fostering Enjoyment," Media Psychol., vol. 15, no. 3, pp. 356–371, 2012.

[33] C. I. Johnson, S. K. T. Bailey, and W. L. Van Buskirk, "Designing Effective Feedback Messages in Serious Games and Simulations : A in Serious Games and Simulations : A Research Review," Springer Int. Publ. Switz. 2017, no. September, 2017.

[34] T. G. Toolkit, "Game Elements," pp. 1–16, 2010.

[35] R. C. Allen, "Discovery of an embodying self: Cancer, identities, narratives.," Diss. Abstr. Int. Sect. B Sci. Eng., vol. 56, no. 9–B, p. 5192, 1996.

[36] N. El Mawas, J. Cahier, and I. C. D. Tech-cico, "Co-designing a serious game to train Emergency Medical Services," 978-1-4673-6404-1/13/$31.00 ©2013 IEEE, pp. 588–593, 2013.

[37] S. H. Hsu, J.-W. Chang, and C.-C. Lee, "Designing Attractive Gamification Features for Collaborative Storytelling Websites," Cyberpsychology, Behav. Soc. Netw., vol. 16, no. 6, pp. 428–435, 2013.

[38] M. V Birk, R. L. Mandryk, and C. Atkins, "The Motivational Push of Games : The Interplay of Intrinsic Motivation and External Rewards in Games for Training," CHI Play '16, Oct. 16-19, 2016, Austin, TX, USA © 2016 ACM. ISBN 978-1-4503-4456-2/16/10…$15.00 DOI http//dx.doi.org/10.1145/2967934.2968091, 2016.

[39] E. Boyle, T. M. Connolly, and T. Hainey, "The role of psychology in understanding the impact of computer games," Entertain. Comput., vol. 2, no. 2, pp. 69–74, 2011.

[40] R. L. Nabi and M. Krcmar, "Conceptualizing media enjoyment as attitude: Implications for mass media effects research," Commun. Theory, vol. 14, no. 4, pp. 288–310, 2004.

[41] S. J. Robinson and S. J. Robinson, "ScienceDirect ScienceDirect ScienceDirect How can Psychology inform disaster research ? How can Psychology inform disaster research ?," Procedia Eng., vol. 212, pp. 1083–1090, 2018.

[42] D. Öztürk, N. Çal, Z. Gocmen, A. Karada, and H. Karabulut, "Nurse Education Today Determining the effect of periodic training on the basic psychomotor skills of nursing students," 2014 Elsevier Ltd. All rights Reserv., vol. 35, pp. 402–407, 2015.

[43] R. M. Ryan and E. L. Deci, "Intrinsic and Extrinsic Motivations : Classic Definitions and New Directions," Contemp. Educ. Psychol. 25, 54–67, vol. 67, pp. 54–67, 2000.

[44] E. D. Mekler, F. Brühlmann, A. N. Tuch, and K. Opwis, "Computers in Human Behavior Towards understanding the effects of individual gami fi cation elements on intrinsic motivation and performance," Comput. Human Behav., 2015.

[45] R. Tamborini, N. D. Bowman, A. Eden, M. Grizzard, and A. Organ, "Defining Media Enjoyment as the Satisfaction of Intrinsic Needs," J. Commun., vol. 60, no. 4, pp. 758–777, 2010.

[46] C. McGregor, B. Bonnis, B. Stanfield, and M. Stanfield, "A Method for Real-Time Stimulation and Response Monitoring Using Big Data and Its

[47] J. Tixier, A. Dandrieux, and P. Slangen, "Training decision-makers : Existing strategies for natural and technological crisis management and specifications of an improved simulation-based tool," Saf. Sci., 2016.

[48] T. S. M. T. W. & N. A. H. M. Z. Hairul Fahmi Md. Muslim, "Kajian awal permainan serius untuk warga emas mengidap diabetes," Proceeding Glob. Summit Educ. GSE 2014 (E-ISBN 978-967-11768-5-6). 4-5 March 2014, Kuala Lumpur, MALAYSIA. Organ. by WorldConferences.net, vol. 2014, no. March, pp. 553–559, 2014.

[49] M. T. P. Adam, "A Serious Game Using Physiological Interfaces For Emotion Regulation Training In The Context Of Financial Decision-," Assoc. Inf. Syst. AIS Electron. Libr. ECIS 2012 Proc. Eur., 2012.

Application to Tactical Training," 2015 IEEE 28th Int. Symp. Comput. Med. Syst., pp. 169–170, 2015.

# Comparison of Reducing the Speckle Noise in Ultrasound Medical Images using Discrete Wavelet Transform

Asim ur Rehman Khan[1]

Electrical Engineering Department
National University of Computer and
Emerging Sciences
Karachi, Pakistan

Farrokh Janabi-Sharifi[2],
Mohammad Ghahramani[3]

Mechanical and Industrial
Engineering
Ryerson University
Toronto, Ontario, Canada

Muhammad Ahsan Rehman
Khan[4]

Department of Medicine
Dr. Ruth K. M. Pfau, Civil Hospital
Karachi, Pakistan

*Abstract*—**Speckle noise in ultrasound (US) medical images is the prime factor that undermines its full utilization. This noise is added by the constructive / destructive interference of sound waves travelling through hard- and soft-tissues of a patient. It is therefore generally accepted that the noise is unavoidable. As an alternate researchers have proposed several algorithms to somewhat undermine the effect of speckle noise. The discrete wavelet transform (DWT) has been used by several researchers. However, the performance of only a few transforms has been demonstrated. This paper provides a comparison of several DWT. The algorithm comprises of a pre-processing stage using Wiener filter, and a post-processing stage using Median filter. The processed image is compared with the original image on four metrics: two are based on full-reference (FR) image quality assessment (IQA), and the remaining two are based on no-reference (NR) IQA metrics. The FR-IQA are peak signal-to-noise ratio (PSNR) and mean structurally similarity index measure (MSSIM). The two NR-IQA techniques are blind pseudo-reference image (BPRI), and blind multiple pseudo-reference images (BMPRI). It has been demonstrated that some of these wavelet transforms outperform others by a significant margin.**

*Keywords*—*Discrete wavelet transform; image quality assessment; ultrasound medical image*

## I. INTRODUCTION

An ultrasound (US) medical image helps in an early diagnosis of kidney stones. These stones cause severe pain in situations where they become large or block the flow of urine. In rare situations, a small stone is stuck in the ureter. The ureter is a small tube connecting kidney and bladder. As per statistics, 1 in 11 persons in USA suffer from kidney stones [1]. An early treatment can save someone from severe pain, cost and medical complexities. The US imaging is quick, non-invasive, cost effective, and has no known side effects to the best of our knowledge. The medical complexities involving kidney stones are considered high-risk illnesses. These are life threatening if left untreated for a long time.

The size and location of kidney stone is also important. A patient may experience no symptom at all to severe, incapacitating pain in the loin requiring urgent treatment. At times the patient is complaining burning sensation during passage of urine, blood, or small stone debris in their urine. Stones can block the main outflow of urine from the kidney leading to irreversible kidney damage, disturbances in the biochemical balance of the body and eventually death. Thus a safe, economical and rapid imaging technique requiring no prior preparation is invaluable in saving many lives.

The practical issue of speckle noise is not new. This has been addressed by several researchers, some as early as 1980's. Jain had initially approximated the speckle noise as multiplicative [2]. He suggested to apply homomorphic filter. The pre- and post-processing were performed using Wiener, and the Median filter, respectively. Chien-Min used Bayesian approach for removing the speckle noise [3]. Perona and Malik introduced an edge preserving approach of anisotropic diffusion (AD) in [4]. The AD filter was subsequently used by several researchers for speckle noise reduction [5]-[8]. The estimation of signal and noise using Kuan's filter has been used for speckle reduction anisotropic diffusion (SRAD) in [5]. SRAD using filtering across image contours and principle curvature directions are given in [6]. The probabilistic model-based SRAD is discussed in [7]. A comparison of SRAD and several other schemes are presented in [8].

The introduction of wavelets during the early 90's resulted in several papers on speckle noise. Mallat introduced multichannel decomposition of images using wavelet transform in [9]. The wavelet theory for image coding was developed in [10]. The application of wavelet packets was presented in [11]. The Bayesian maximum a posteriori (MAP) estimator based design were illustrated in [12]-[13]. The rational-dilation wavelet transform (RADWT) and non-linear bilateral filter based approaches were proposed in [14]. The genetic algorithm based solution was introduced in [15]. A quantum-inspired de-speckling method was discussed in [16]. The wavelet and fuzzy theory based approach is given in [17].

During the last two decades, several new wavelets have been derived. These wavelets offer numerous benefits. Some of them are preferred over the others in terms of number of stages, linearity, ease of use, etc. The performance of Symlet wavelets has been discussed in [18]. A comparison of five wavelets Haar, Daubechies, Symlet, Coiflet and biorthogonal

wavelets for removing the speckle noise has been given in [19]. The applications of various wavelets for identification of bone fracture has been discussed in [20]. A comparative study of Birge-Massart strategy for setting a threshold for image compression is given in [21]. The identification and classification of colonic polyps using wavelets transforms has been demonstrated in [22].

This paper compares the performance of seven wavelets for reducing the speckle noise in US medical images. The selected discrete wavelets are Haar, Daubechies, Symlet, Coiflet, biorthogonal, reverse biorthogonal, and discrete Mayer wavelets. The pre- and post-processing is performed by using Weiner and Median filters, respectively. The introduction is followed by evaluation criteria in Section 2. The description of selected discrete wavelet transforms is given in Section 3. The performances of these wavelets are given in simulations in Section 4. Section 5 concludes this paper.

## II. EVALUATION CRITERIA

There are generally two ways of image quality assessment (IQA). The first approach compares the results with the original image. This is termed as full-reference (FR) IQA [23]. The second approach is more recent in which the assumption is that no reference image is available. This is referred to as blind or no-reference (NR) IQA [24]-[25]. In case of an US image, a truly clean image without speckle noise is not really available. In this regard, the NR-IQA seems to be a better metrics for US images.

### A. Full-Reference IQA (FR-IQA)

The performance of various wavelets is tested using two FR-IQA. The first is based on peak signal to noise ratio (PSNR) and the second is mean structurally similarity index measure (MSSIM). The mean square error (MSE) is used as a criterion for comparing the original image with the processed image as given by

$$MSE = \frac{1}{MN}\sum_{i=1}^{M}\sum_{j=1}^{N}(y_{ij} - \hat{y}_{ij})^2 \qquad (1)$$

where, $y_{ij}$ is the pixel value of the original image and $\hat{y}_{ij}$ is the estimated value of $y_{ij}$. The image row and column numbers are given by M, and N, respectively. The peak signal-to-noise ratio in decibels $PSNR_{dB}$ is found by using,

$$PSNR_{dB} = 20\ log_{10}\frac{2^B - 1}{\sqrt{MSE}} = 10\ log_{10}\frac{(255)^2}{MSE} \qquad (2)$$

The superscript B is the number of bits in a pixel. The value of B is taken as 8, resulting in 256 grey shades. Each pixel value therefore varies in the range of 0-255.

The structurally similarity index measure (SSIM) was proposed in 2004 [23]. It compares the mean and variance of two images. The two images are considered by x and y. The SSIM is given as,

$$SSIM(x,y) = \frac{\big(2\mu_x\mu_y + c_1\big)\big(2\sigma_{xy} + c_2\big)}{\big(\mu_x^2 + \mu_y^2 + c_1\big)\big(\sigma_x^2 + \sigma_y^2 + c_2\big)} \qquad (3)$$

where $\mu_x$, $\mu_y$ are the average of images x and y, respectively. $\sigma_x^2$, and $\sigma_y^2$ are the variances of images x, and y, respectively. $c_1 = (k_1L)^2$, $c_2 = (k_2L)^2$ are two variables to stabilize the division with weak denominator. L is the dynamic range of the pixel value 255. The $k_1$ and $k_2$ are equal to 0.01 and 0.03 (taken by default), respectively. Another single parameter used for images is MSSIM given as [23].

$$MSSIM = \frac{1}{M}\sum_{j=1}^{N} SSIM(x_i, y_j) \qquad (4)$$

The SSIM is the similarity index calculated over a small region of an image. The MSSIM is the mean value of SSIM across all windows. The MATLAB code for MSSIM is available in [26].

### B. No-Reference IQA (NR-IQA)

The no-reference (NR) IQA assumes that a clean image is not available, and the image quality is assessed based on the available noisy image. In this paper two NR-IQA techniques are used to test the quality of US images. Both the techniques are essentially based on the use of a pseudo-reference image (PRI). A PRI is generated from the distorted image. In the absence of a clean image, the PRI is used as a reference image for the targeted US image quality assessment.

The first NR assessment technique is blind PRI (BPRI) [24]. This technique measures the image quality in terms of block effects, sharpness, and noise. The block effect is found by using pseudo structure similarity (PSS) index. The local binary pattern (LBP) is used for sharpness, and noise measurement. The PRI-based sharpness and noise is used to derive the local structure similarity (LSS) index. The second NR assessment technique is multiple PRI (MPRI) [25]. In this technique, the image is distorted with an aggregation of four types of distortions, based on JPEG compression, JPEG2000 compression, Gaussian blur (GB), and white Gaussian noise (WGN).

## III. WAVELET SELECTION

The Fourier and Laplace transforms have been used extensively for extracting the significant frequency components of a noisy image. These transforms perform very well by translating the time domain signal into frequency domain. The main limitation of these transforms is that the local information is lost. In an analogue or digital signal transmission system, the location of noise is generally not very critical. This is different in image processing, where the perceived quality of an image depends on the location of noise. As an example, the loss of signal at the edges of an image can be acceptable, but the loss of fine details at the centre of an image, or around the critical regions is quite unacceptable. A significant advantage of wavelet transform over the previously available transforms is that the wavelet not only translates time-domain into frequency-domain, it also preserves the physical location of noise present in an image. This advantage has no parallel in the other transforms.

During the last two decades, several wavelets have been derived with different set of properties. Before proceeding on wavelets, it is important to review few basics of wavelet

transforms. An image decomposed by wavelet has two functions. These are scaling function and wavelet function, sometimes refer to as the mother wavelet, and the father wavelet. The scaling component gives lower frequency components corresponding to variations in the grey shades. The wavelet function gives high frequency components like edges. The scaling function, $\varphi(t)$ and the wavelet function, $\psi(t)$ are given by,

$$\varphi(t) = \sum_n h(n)\sqrt{2}\ \varphi(2t - n)$$

$$\psi(t) = \sum_n g(n)\sqrt{2}\ \varphi(2t - n) \tag{5}$$

The $h(n)$ and $g(n)$ are low-pass and high-pass filters, respectively. The n is the periodic shift that implements the filter coefficient index. Both filters are related by,

$$g(n) = (-1)^n h(N - n - 1) \tag{6}$$

where N is the number of vanishing moments. A wavelet with N vanishing moment has at least a polynomial of order N-1. The vanishing moments represent level of differentiability of a function. A wavelet with vanishing moment N is defined as multi-scale differential of order N. This, in essence, defines the local irregularity of a signal. A smaller value of N is therefore preferred over the larger values. An N vanishing moment corresponds to 2N taps in the filter bank. The filter is implemented as a finite-impulse response (FIR) filter. A smaller value of N, therefore, corresponds to a shorter filter with less number of taps.

TABLE I.      WAVELETS USED FOR US IMAGE ANALYSIS

| S. No. | Wavelet Families | MATLAB Functions | |
|---|---|---|---|
| 1 | Haar | Haar | Haar |
| 2 | Daubechies | dbN | db2, db3, db4, db5, db6, db7, db8, db9, db10, db11, db12, db13, db14, db15, db16, db17, db18, db19, db20, db21, db22, db23, db24, db25, db26 |
| 3 | Symlet | symN | sym1, sym2, sym3, sym4, sym5, sym6, sym7, sym8, sym9, sym10 |
| 4 | Coiflet | coifN | coif1, coif2, coif3, coif4, coif5 |
| 5 | Biorthogonal | bior $N_r.N_d$ | bior1.1, bior2.2, bior3.1, bior3.3, bior4.4, bior5.5, bior6.8 |
| 6 | Reverse Biorthogonal | rbio $N_r.N_d$ | rbio1.1, rbio2.2, rbio3.1, rbio3.3, rbio4.4 |
| 7 | Discrete Meyer | dmey | dmey |

Haar wavelet is the oldest and the simplest of all wavelets. This is the only orthogonal wavelet having linear phase. Haar wavelet decomposes the discrete signal into two sub-signals of half its length. One sub-signal provide the trend, while the second sub-signal gives difference or fluctuations. The main advantage of Haar wavelet is that it is fast, memory efficient, and conceptually simple to implement.

The biorthogonal wavelets have linear phase. These filters have a pair of scaling functions and an associated scaling filters used for analysis and synthesis. The analysis and synthesis filters can be designed to have different order of vanishing moments. It is possible to have greater number of vanishing moments for sparse representation analysis and a smoother wavelet for reconstruction. In MATLAB notation, the biorthogonal wavelets are designated as 'biorN$_r$.N$_d$'. Similarly, the reverse biorthogonal wavelets are represented as 'rbioN$_r$.N$_d$'. The 'N$_r$' represents the effective number of reconstruction filter, and the 'N$_d$' represents effective number of decomposition filters. Table I gives several choices of 'N$_r$' and 'N$_d$'.

The Daubechies wavelet has several versions represented by vanishing moments, N. In Symlet wavelet, the value of 'N' varies from 1 to 10. The Coiflet wavelet has 5 variations represented by the value of 'N' that equals 1 to 5. In all the above wavelets, the number of taps in synthesis and analysis filters are same. The discrete Meyer wavelet has a single transform. A comprehensive mathematical analysis of wavelet theory is given in [27].

The Haar and Daubechies wavelets are orthogonal wavelets, while biorthogonal and reverse biorthogonal wavelets are biorthogonal in nature. The discrete Meyer wavelet is simply a discrete version of the continuous Meyer wavelet. A compactly supported wavelet function restricts itself to within certain limits and as a consequence the signal is also restricted to within some limits. The Haar, Daubechies, and reverse orthogonal wavelets have compactly supported functions. Both the Daubechies and reverse biorthogonal wavelets show an arbitrary number of vanishing moments. All the selected wavelets have finite impulse response (FIR). The FIR has an advantage of having only a few non-zero coefficients.

## IV. SIMULATIONS

The algorithm is tested on six US medical images of abnormal kidneys with stones. These images are downloaded from the US imaging database [28]. The database comprises of large collection of US images that are categorized as fetal, kidney, renal calculi, appendix, urinary bladder, liver, spleen, chest and vascular system. The images in each category consists of high resolution samples. The low resolution images are good for fast processing; however, they are not appropriate as the fine details are lost. Also, the outcome may have little practical value. The high-resolution images provide enough details but they require more processing power as well as time. If these images are being used during an operation, then the speed of processing needs to be sufficiently higher to give real-life response to the ongoing activities, like surgery and other diagnostic treatments.

Fig 1.   Flow diagram of US image processing.

The selected kidney images are US images of patients complaining about minor to severe pain in the left region of their abdomen. All samples are for male patients. The initial diagnosis recommended US imaging for further treatment. The US images clearly showed stones, but the number of stones, and their sizes are not clearly identified as the images contain significant amount of speckle noise. The stones are more identifiable if they are either in larger size or present close enough to form a larger area in concerned region. Unfortunately, in most cases the kidney stones have relatively smaller sizes, and they are scattered across the whole active region of a kidney. In US images, kidney stones usually appear in 'white' or lighter grey shades. A distinct feature of the stones is that they always have a long 'shadow' that originates from the stone, and spreads out towards the outer edge of the kidney. These shadows are quite visible in most of the US images; however, in few cases the shadow is much lighter and may be overlooked. Image enhancement helps in improving the visibility of these shadows in such situations. In general, there are two primary objectives. The first is to reduce the speckle noise. The second is to improve the contrast level of regions containing significant amount of information like stones.

Fig. 1 gives various steps involving US medical image processing. These are pre-processing, filtering, and post-processing. In pre-processing, Wiener filter helps in smoothing the image. This has been demonstrated that although this causes sharp edges to be slightly smoothed out, but as a result the overall visibility improves. The pre-processing is followed by filtering, using several discrete wavelet transforms. The selected transforms are Haar, Daubechies, Symlet, Coiflet, biorthogonal, reverse biorthogonal, and discrete Meyer wavelets. The Haar and discrete Meyer wavelets are distinct in nature. The remaining have several combinations. A list of selected wavelets is given in Table I. The post-processing is based on Median filter that helps in restoring the sharp edges. The window size used in the Wiener filter is (3 x 3) pixels. It has been observed that this is a preferred window size than the larger window size of (5 x 5), or (7 x 7) pixels. The two-dimensional Median filter is applied using (5 x 5) pixel window. The subjective tests have supported the use of (5 x 5) pixel windows against the possible (3 x 3) pixel or (7 x 7) pixel windows.

The processing of wavelet is performed only at the first level of decomposition. The coefficients in the horizontal, vertical, and the diagonal directions are filtered using a hard-limiter as,

$$h(\omega, \lambda) = \begin{cases} \omega & |\omega| > \lambda \\ 0 & |\omega| \leq \lambda \end{cases} \qquad (7)$$

where $\lambda$ is the threshold. The value of $\lambda$ in the horizontal, vertical, and the diagonal directions are selected as 200. This high value of threshold essentially removes most of the effects present in the above three directions. The low resolution coefficient cut-off is dynamically calculated by taking the mean pixel value of the selected image. It has been observed that a higher threshold value of low resolution coefficients result in severe quality degradation.

The qualitative analysis of the above algorithms is tested on six US images of various granularity levels. The unprocessed, original images are given in the left column of Fig. 4. The location of stones and their shadows are marked with arrows on original images. The processed images are given on the right side. The histogram of original and the processed images are given on right side of the successive image. A quick comparison of the results reveal that a significant amount of speckle noise has been removed. At the same time, the contrast of the images has also improved. The histogram of all images show that the lower pixel values (darker grey) have been filtered out. The grey shades at higher values (lighter grey) have smoothed out, but the general shape has remained same. This has resulted in making the shadow region more prominent. The shadows in original images as shown in Fig. 4(a)-(e) are quite visible. These have become clearer in the processed images given on the right side. The shadows in Fig. 4(b), and f are not very clear in the original images, but they have become quite visible in the processed images.

The quantitative analysis of the images are given in Table II through Table V. As mentioned earlier, the performance of various wavelet transforms are analyzed on four criterions broadly categorized as FR-IQA, and NR-IQA. The FR-IQA techniques are PSNR, and MSSIM, while the NR-IQA techniques are BPRI, and BMPRI. The PSNR and MSSIM metrics are given in Table II and Table III, respectively. The PSNR is given in decibels (dB), while the similarity index using MSSIM is given in a value between 0 and 1. A value closer to 1 corresponds to a higher similarity across images. It is clear that the PSNR value depends on image contents. This is the reason that the PSNR of six selected images varies. However, the PSNR of any one image processed through a wavelet has smaller variation in terms of decibels (dB). This is also observed that a higher value of MSSIM corresponds to a higher value of PSNR. The metrics of BPRI, and BMPRI techniques are given in Table IV and Table V, respectively. The BPRI index is in the range of 0 to 1, while the BMPRI is in the range of 20 to 60.

TABLE II.        PEAK SIGNAL-TO-NOISE RATIO (PSNR)

| MATLAB function | (a) | (b) | (c) | (d) | (e) | (f) | MATLAB function | (a) | (b) | (c) | (d) | (e) | (f) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| haar | 31.55 | 27.12 | 23.72 | 26.41 | 28.33 | 29.71 | sym1 | 31.55 | 27.12 | 23.72 | 26.41 | 28.33 | 29.71 |
| db2 | 32.16 | 28.07 | 24.05 | 27.35 | 28.80 | 30.08 | sym2 | 32.16 | 28.07 | 24.05 | 27.35 | 28.80 | 30.08 |
| db3 | 32.26 | 28.33 | 24.16 | 27.65 | 29.15 | 30.37 | sym3 | 32.26 | 28.33 | 24.16 | 27.65 | 29.15 | 30.37 |
| db4 | 32.23 | 28.24 | 24.14 | 27.47 | 29.11 | 30.39 | sym4 | 32.23 | 28.36 | 24.19 | 27.55 | 29.05 | 30.28 |
| db5 | 32.31 | 28.43 | 24.21 | 27.52 | 29.20 | 30.48 | sym5 | 32.25 | 28.26 | 24.19 | 27.20 | 28.91 | 30.24 |
| db6 | 32.40 | 28.60 | 24.26 | 27.71 | 29.21 | 30.49 | sym6 | 32.27 | 28.47 | 24.22 | 27.61 | 29.12 | 30.34 |
| db7 | 32.40 | 28.42 | 24.23 | 27.44 | 29.10 | 30.38 | sym7 | 32.37 | 28.54 | 24.26 | 27.73 | 29.23 | 30.47 |
| db8 | 32.32 | 28.39 | 24.23 | 27.35 | 29.06 | 30.37 | sym8 | 32.29 | 28.55 | 24.24 | 27.63 | 29.14 | 30.37 |
| db9 | 32.40 | 28.56 | 24.28 | 27.65 | 29.17 | 30.46 | sym9 | 32.27 | 28.45 | 24.25 | 27.35 | 28.98 | 30.34 |
| db10 | 32.44 | 28.62 | 24.29 | 27.67 | 29.24 | 30.53 | sym10 | 32.31 | 28.58 | 24.25 | 27.64 | 29.14 | 30.39 |
| db11 | 32.46 | 28.58 | 24.28 | 27.59 | 29.24 | 30.52 | coif1 | 32.19 | 28.09 | 24.06 | 27.49 | 28.97 | 30.18 |
| db12 | 32.40 | 28.57 | 24.28 | 27.53 | 29.22 | 30.52 | coif2 | 32.24 | 28.46 | 24.20 | 27.69 | 29.12 | 30.34 |
| db13 | 32.42 | 28.65 | 24.28 | 27.45 | 29.26 | 30.53 | coif3 | 32.26 | 28.56 | 24.23 | 27.73 | 29.15 | 30.38 |
| db14 | 32.50 | 28.66 | 24.29 | 27.49 | 29.21 | 30.51 | coif4 | 32.27 | 28.60 | 24.25 | 27.75 | 29.16 | 30.40 |
| db15 | 32.51 | 28.61 | 24.28 | 27.45 | 29.17 | 30.47 | coif5 | 32.28 | 28.62 | 24.25 | 27.75 | 29.17 | 30.41 |
| db16 | 32.47 | 28.60 | 24.29 | 27.51 | 29.18 | 30.47 | bior1.1 | 31.55 | 27.12 | 23.72 | 26.41 | 28.33 | 29.71 |
| db17 | 32.42 | 28.58 | 24.30 | 27.62 | 29.24 | 30.50 | bior2.2 | 32.20 | 28.12 | 24.11 | 27.55 | 29.05 | 30.28 |
| db18 | 32.47 | 28.66 | 24.29 | 27.57 | 29.29 | 30.52 | bior3.1 | 32.41 | 28.32 | 24.18 | 27.38 | 29.25 | 30.54 |
| db19 | 32.51 | 28.66 | 24.29 | 27.55 | 29.29 | 30.54 | bior3.3 | 32.43 | 28.39 | 24.23 | 27.51 | 29.28 | **30.59** |
| db20 | 32.51 | 28.68 | 24.30 | 27.61 | **29.29** | 30.57 | bior4.4 | 32.17 | 28.50 | 24.20 | 27.73 | 29.10 | 30.32 |
| db21 | 32.46 | 28.70 | 24.30 | 27.60 | 29.24 | 30.56 | bior5.5 | 32.12 | 28.60 | 24.23 | 27.61 | 29.07 | 30.30 |
| db22 | 32.46 | 28.66 | 24.29 | 27.47 | 29.20 | 30.55 | bior6.8 | 32.27 | 28.59 | 24.25 | **27.77** | 29.17 | 30.41 |
| db23 | 32.51 | 28.67 | 24.30 | 27.45 | 29.19 | 30.51 | rbio1.1 | 31.55 | 27.12 | 23.72 | 26.41 | 28.33 | 29.71 |
| db24 | **32.54** | 28.67 | **24.32** | 27.51 | 29.22 | 30.52 | rbio2.2 | 32.20 | 27.88 | 23.94 | 26.99 | 28.86 | 30.06 |
| db25 | 32.50 | 28.65 | 24.31 | 27.59 | 29.25 | 30.54 | rbio3.1 | 31.91 | **28.96** | 24.02 | 26.65 | 28.91 | 30.22 |
| db26 | 32.48 | 28.65 | 24.30 | 27.52 | 29.27 | 30.54 | rbio3.3 | 31.97 | 28.72 | 23.85 | 26.37 | 28.75 | 30.13 |
| | | | | | | | rbio4.4 | 32.30 | 28.37 | 24.17 | 27.65 | 29.14 | 30.34 |
| | | | | | | | dmey | 32.29 | 28.67 | 24.27 | 27.74 | 29.17 | 30.43 |

TABLE III.     Mean Structurally Similarity Index Measure (MSSIM)

| MATLAB function | (a) | (b) | (c) | (d) | (e) | (f) | MATLAB function | (a) | (b) | (c) | (d) | (e) | (f) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **haar** | 0.932 | 0.842 | 0.727 | 0.911 | 0.861 | 0.792 | **sym1** | 0.932 | 0.842 | 0.727 | 0.911 | 0.861 | 0.792 |
| **db2** | 0.935 | 0.853 | 0.752 | 0.919 | 0.874 | 0.805 | **sym2** | 0.935 | 0.853 | 0.752 | 0.919 | 0.874 | 0.805 |
| **db3** | 0.936 | 0.855 | 0.764 | 0.923 | 0.880 | 0.810 | **sym3** | 0.936 | 0.855 | 0.764 | 0.923 | 0.880 | 0.810 |
| **db4** | 0.936 | 0.855 | 0.766 | 0.922 | 0.882 | 0.811 | **sym4** | 0.936 | 0.855 | 0.763 | 0.924 | 0.881 | 0.811 |
| **db5** | 0.936 | 0.856 | 0.768 | 0.923 | 0.883 | 0.814 | **sym5** | 0.936 | 0.855 | 0.760 | 0.925 | 0.880 | 0.811 |
| **db6** | 0.937 | 0.857 | 0.769 | 0.925 | 0.925 | 0.815 | **sym6** | 0.936 | 0.856 | 0.766 | 0.926 | 0.882 | 0.812 |
| **db7** | 0.937 | 0.857 | 0.764 | 0.925 | 0.884 | 0.815 | **sym7** | 0.937 | 0.857 | 0.770 | 0.926 | 0.883 | 0.814 |
| **db8** | 0.936 | 0.857 | 0.764 | 0.925 | 0.883 | 0.813 | **sym8** | 0.936 | 0.857 | 0.767 | 0.926 | 0.882 | 0.813 |
| **db9** | 0.937 | 0.858 | 0.768 | 0.927 | 0.885 | 0.816 | **sym9** | 0.936 | 0.856 | 0.764 | 0.926 | 0.881 | 0.813 |
| **db10** | 0.937 | 0.859 | 0.772 | 0.926 | 0.886 | 0.818 | **sym10** | 0.936 | 0.857 | 0.767 | 0.926 | 0.883 | 0.813 |
| **db11** | 0.937 | 0.859 | 0.773 | 0.926 | 0.886 | 0.818 | **coif1** | 0.935 | 0.853 | 0.756 | 0.920 | 0.876 | 0.807 |
| **db12** | 0.937 | 0.858 | 0.772 | 0.926 | 0.885 | 0.817 | **coif2** | 0.936 | 0.856 | 0.766 | 0.925 | 0.881 | 0.812 |
| **db13** | 0.937 | 0.858 | 0.771 | 0.926 | 0.886 | 0.817 | **coif3** | 0.936 | 0.856 | 0.768 | 0.926 | 0.882 | 0.813 |
| **db14** | 0.938 | 0.859 | 0.769 | 0.926 | 0.887 | 0.818 | **coif4** | 0.936 | 0.857 | 0.768 | 0.927 | 0.882 | 0.813 |
| **db15** | 0.938 | 0.859 | 0.767 | 0.926 | 0.887 | 0.818 | **coif5** | 0.936 | 0.857 | 0.769 | 0.927 | 0.883 | 0.814 |
| **db16** | 0.937 | 0.859 | 0.769 | 0.927 | 0.886 | 0.817 | **bior1.1** | 0.932 | 0.842 | 0.727 | 0.911 | 0.861 | 0.792 |
| **db17** | 0.938 | 0.859 | 0.772 | 0.926 | 0.886 | 0.818 | **bior2.2** | 0.935 | 0.854 | 0.761 | 0.922 | 0.879 | 0.810 |
| **db18** | 0.938 | 0.860 | 0.772 | 0.926 | 0.888 | 0.819 | **bior3.1** | 0.937 | 0.857 | 0.767 | 0.921 | 0.884 | 0.818 |
| **db19** | 0.938 | 0.860 | 0.773 | 0.926 | 0.888 | 0.819 | **bior3.3** | 0.937 | 0.858 | 0.772 | 0.924 | 0.887 | **0.821** |
| **db20** | 0.938 | 0.860 | 0.773 | 0.927 | **0.887** | 0.819 | **bior4.4** | 0.935 | 0.855 | 0.766 | 0.925 | 0.880 | 0.811 |
| **db21** | 0.938 | 0.860 | 0.772 | 0.926 | 0.887 | 0.819 | **bior5.5** | 0.936 | 0.855 | 0.765 | 0.924 | 0.879 | 0.809 |
| **db22** | 0.938 | 0.860 | 0.770 | 0.926 | 0.887 | 0.820 | **bior6.8** | 0.936 | 0.857 | 0.769 | **0.927** | 0.882 | 0.813 |
| **db23** | 0.938 | 0.860 | 0.769 | 0.927 | 0.888 | 0.820 | **rbio1.1** | 0.932 | 0.842 | 0.727 | 0.911 | 0.861 | 0.792 |
| **db24** | **0.938** | 0.860 | **0.771** | 0.926 | 0.888 | 0.820 | **rbio2.2** | 0.935 | 0.851 | 0.746 | 0.914 | 0.871 | 0.801 |
| **db25** | 0.938 | 0.860 | 0.773 | 0.926 | 0.888 | 0.820 | **rbio3.1** | 0.935 | **0.855** | 0.776 | 0.923 | 0.884 | 0.813 |
| **db26** | 0.938 | 0.860 | 0.773 | 0.926 | 0.888 | 0.820 | **rbio3.3** | 0.935 | 0.853 | 0.755 | 0.912 | 0.874 | 0.804 |
| | | | | | | | **rbio4.4** | 0.936 | 0.856 | 0.765 | 0.924 | 0.882 | 0.812 |
| | | | | | | | **dmey** | 0.936 | 0.857 | 0.769 | 0.927 | 0.883 | 0.814 |

TABLE IV.    BLIND PSEUDO-REFERENCE IMAGE (BPRI) METRICS

| MATLAB function | (a) | (b) | (c) | (d) | (e) | (f) | MATLAB function | (a) | (b) | (c) | (d) | (e) | (f) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| haar | 0.058 | 0.034 | 0.033 | 0.046 | 0.040 | 0.033 | sym1 | 0.058 | 0.034 | 0.033 | 0.046 | 0.040 | 0.033 |
| db2 | 0.119 | 0.065 | 0.072 | 0.088 | 0.079 | 0.069 | sym2 | 0.119 | 0.065 | 0.072 | 0.088 | 0.079 | 0.069 |
| db3 | 0.125 | 0.069 | 0.077 | 0.091 | 0.084 | 0.072 | sym3 | 0.125 | 0.069 | 0.077 | 0.091 | 0.084 | 0.072 |
| db4 | 0.125 | 0.069 | 0.076 | 0.090 | 0.083 | 0.072 | sym4 | 0.122 | 0.068 | 0.076 | 0.091 | 0.083 | 0.071 |
| db5 | 0.126 | 0.070 | 0.077 | 0.091 | 0.084 | 0.073 | sym5 | 0.125 | 0.069 | 0.076 | 0.091 | 0.083 | 0.073 |
| db6 | 0.129 | 0.071 | 0.079 | 0.091 | 0.087 | 0.074 | sym6 | 0.124 | 0.069 | 0.077 | 0.091 | 0.084 | 0.072 |
| db7 | 0.124 | 0.070 | 0.075 | 0.089 | 0.083 | 0.073 | sym7 | 0.128 | 0.071 | 0.079 | 0.093 | 0.086 | 0.075 |
| db8 | 0.125 | 0.069 | 0.074 | 0.088 | 0.084 | 0.074 | sym8 | 0.126 | 0.070 | 0.077 | 0.092 | 0.084 | 0.072 |
| db9 | 0.127 | 0.071 | 0.078 | 0.088 | 0.085 | 0.073 | sym9 | 0.127 | 0.070 | 0.077 | 0.091 | 0.085 | 0.074 |
| db10 | 0.126 | 0.070 | 0.078 | 0.088 | 0.084 | 0.074 | sym10 | 0.127 | 0.070 | 0.078 | 0.092 | 0.085 | 0.073 |
| db11 | 0.126 | 0.070 | 0.076 | 0.088 | 0.085 | 0.073 | coif1 | 0.115 | 0.064 | 0.073 | 0.087 | 0.077 | 0.066 |
| db12 | 0.126 | 0.070 | 0.077 | 0.088 | 0.084 | 0.073 | coif2 | 0.126 | 0.069 | 0.079 | 0.096 | 0.085 | 0.072 |
| db13 | 0.125 | 0.069 | 0.076 | 0.088 | 0.084 | 0.071 | coif3 | 0.128 | 0.070 | 0.080 | 0.095 | 0.086 | 0.074 |
| db14 | 0.124 | 0.070 | 0.075 | 0.087 | 0.083 | 0.073 | coif4 | 0.127 | 0.070 | 0.078 | 0.092 | 0.085 | 0.073 |
| db15 | 0.125 | 0.070 | 0.076 | 0.088 | 0.085 | 0.074 | coif5 | 0.126 | 0.070 | 0.077 | 0.090 | 0.084 | 0.073 |
| db16 | 0.126 | 0.070 | 0.077 | 0.086 | 0.085 | 0.073 | bior1.1 | 0.058 | 0.034 | 0.033 | 0.046 | 0.040 | 0.033 |
| db17 | 0.123 | 0.070 | 0.076 | 0.085 | 0.084 | 0.072 | bior2.2 | 0.118 | 0.064 | 0.074 | 0.088 | 0.078 | 0.066 |
| db18 | 0.124 | 0.069 | 0.076 | 0.086 | 0.083 | 0.073 | bior3.1 | 0.125 | 0.069 | 0.079 | 0.092 | 0.084 | 0.073 |
| db19 | 0.125 | 0.070 | 0.077 | 0.088 | 0.085 | 0.073 | bior3.3 | 0.124 | 0.069 | 0.077 | 0.090 | 0.083 | **0.072** |
| db20 | 0.126 | 0.070 | 0.076 | 0.088 | **0.085** | 0.073 | bior4.4 | 0.117 | 0.065 | 0.074 | 0.089 | 0.080 | 0.068 |
| db21 | 0.123 | 0.069 | 0.076 | 0.085 | 0.082 | 0.072 | bior5.5 | 0.121 | 0.067 | 0.076 | 0.090 | 0.082 | 0.071 |
| db22 | 0.122 | 0.069 | 0.075 | 0.086 | 0.083 | 0.073 | bior6.8 | 0.124 | 0.069 | 0.078 | **0.091** | 0.084 | 0.072 |
| db23 | 0.125 | 0.070 | 0.077 | 0.088 | 0.085 | 0.073 | rbio1.1 | 0.058 | 0.034 | 0.033 | 0.046 | 0.040 | 0.033 |
| db24 | **0.125** | 0.070 | **0.077** | 0.085 | 0.085 | 0.073 | rbio2.2 | 0.115 | 0.064 | 0.071 | 0.086 | 0.077 | 0.066 |
| db25 | 0.123 | 0.069 | 0.076 | 0.085 | 0.084 | 0.072 | rbio3.1 | 0.074 | **0.045** | 0.037 | 0.051 | 0.047 | 0.042 |
| db26 | 0.122 | 0.069 | 0.075 | 0.085 | 0.083 | 0.073 | rbio3.3 | 0.119 | 0.066 | 0.065 | 0.087 | 0.076 | 0.066 |
| | | | | | | | rbio4.4 | 0.124 | 0.069 | 0.077 | 0.092 | 0.083 | 0.072 |
| | | | | | | | dmey | 0.125 | 0.069 | 0.075 | 0.087 | 0.082 | 0.072 |

TABLE V.    BLIND MULTIPLE PSEUDO-REFERENCE IMAGE (BMPRI) METRICS

| MATLAB function | (a) | (b) | (c) | (d) | (e) | (f) | MATLAB function | (a) | (b) | (c) | (d) | (e) | (f) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| haar | 27.94 | 11.58 | 9.83 | 18.31 | 15.13 | 11.35 | sym1 | 27.94 | 11.58 | 9.83 | 18.31 | 15.13 | 11.35 |
| db2 | 54.17 | 31.68 | 36.29 | 40.57 | 38.55 | 32.92 | sym2 | 54.17 | 31.68 | 36.29 | 40.57 | 38.55 | 32.92 |
| db3 | 56.59 | 32.64 | 38.95 | 42.97 | 40.34 | 36.78 | sym3 | 56.59 | 32.64 | 38.95 | 42.97 | 40.34 | 36.78 |
| db4 | 56.82 | 32.75 | 39.37 | 42.84 | 39.26 | 35.94 | sym4 | 55.97 | 32.25 | 38.52 | 42.66 | 40.96 | 34.50 |
| db5 | 57.33 | 33.57 | 38.76 | 43.05 | 41.77 | 36.79 | sym5 | 56.77 | 33.24 | 39.95 | 44.06 | 41.70 | 35.59 |
| db6 | 58.52 | 34.21 | 39.68 | 42.57 | 43.30 | 37.96 | sym6 | 56.63 | 32.49 | 39.74 | 42.51 | 41.19 | 34.89 |
| db7 | 57.61 | 33.83 | 38.68 | 42.69 | 41.21 | 36.58 | sym7 | 58.40 | 34.02 | 41.18 | 44.39 | 41.78 | 38.13 |
| db8 | 56.92 | 33.27 | 38.67 | 42.55 | 42.15 | 36.72 | sym8 | 57.22 | 33.23 | 39.27 | 42.71 | 41.63 | 35.81 |
| db9 | 58.25 | 34.19 | 40.97 | 42.95 | 43.25 | 37.61 | sym9 | 57.53 | 33.67 | 39.94 | 43.92 | 42.54 | 35.98 |
| db10 | 58.45 | 34.27 | 40.11 | 43.04 | 42.60 | 38.17 | sym10 | 57.64 | 33.97 | 39.59 | 43.72 | 42.11 | 36.31 |
| db11 | 58.46 | 33.72 | 40.00 | 43.93 | 42.71 | 37.73 | coif1 | 53.33 | 31.74 | 36.13 | 40.42 | 37.10 | 31.68 |
| db12 | 57.65 | 33.30 | 39.30 | 41.94 | 42.34 | 37.11 | coif2 | 56.89 | 32.54 | 39.78 | 43.06 | 41.39 | 35.39 |
| db13 | 57.52 | 32.91 | 40.45 | 42.80 | 42.81 | 37.13 | coif3 | 57.23 | 33.61 | 40.11 | 44.28 | 42.03 | 36.23 |
| db14 | 57.86 | 34.69 | 39.73 | 42.49 | 42.33 | 37.25 | coif4 | 57.10 | 33.23 | 39.56 | 44.45 | 42.10 | 36.23 |
| db15 | 58.09 | 34.44 | 39.33 | 43.75 | 43.85 | 37.70 | coif5 | 57.14 | 33.67 | 40.02 | 42.97 | 41.99 | 35.86 |
| db16 | 58.00 | 34.16 | 39.93 | 42.42 | 42.99 | 37.32 | bior1.1 | 27.94 | 11.58 | 9.83 | 18.31 | 15.13 | 11.35 |
| db17 | 58.00 | 34.16 | 39.93 | 42.42 | 42.99 | 37.32 | bior2.2 | 54.13 | 30.58 | 36.31 | 40.22 | 38.91 | 31.55 |
| db18 | 58.09 | 34.33 | 39.50 | 43.40 | 43.56 | 37.90 | bior3.1 | 56.84 | 32.81 | 39.24 | 42.77 | 40.60 | 37.49 |
| db19 | 58.47 | 34.23 | 40.54 | 44.28 | 43.25 | 37.26 | bior3.3 | 56.78 | 33.19 | 38.90 | 42.52 | 41.50 | **36.72** |
| db20 | 58.27 | 34.30 | 39.78 | 43.16 | **43.42** | 38.03 | bior4.4 | 54.13 | 30.24 | 36.76 | 40.84 | 39.69 | 32.88 |
| db21 | 57.39 | 33.77 | 39.66 | 41.94 | 43.01 | 37.19 | bior5.5 | 55.58 | 31.59 | 38.85 | 42.15 | 39.98 | 34.32 |
| db22 | 57.59 | 34.38 | 39.53 | 43.13 | 43.47 | 37.49 | bior6.8 | 56.84 | 32.63 | 39.82 | **43.40** | 41.83 | 35.56 |
| db23 | 58.22 | 35.89 | 40.55 | 43.49 | 43.36 | 37.13 | rbio1.1 | 27.94 | 11.58 | 9.83 | 18.31 | 15.13 | 11.35 |
| db24 | **57.90** | 35.58 | **40.64** | 42.61 | 43.28 | 37.29 | rbio2.2 | 53.64 | 31.46 | 35.34 | 39.66 | 37.88 | 31.54 |
| db25 | 57.39 | 35.47 | 40.17 | 42.92 | 43.65 | 37.57 | rbio3.1 | 39.60 | **17.96** | 13.35 | 21.94 | 18.29 | 16.14 |
| db26 | 57.99 | 35.16 | 40.73 | 43.09 | 43.21 | 37.71 | rbio3.3 | 55.31 | 29.48 | 32.05 | 42.12 | 35.63 | 32.70 |
| | | | | | | | rbio4.4 | 56.81 | 33.20 | 38.78 | 42.71 | 41.61 | 34.39 |
| | | | | | | | dmey | 56.85 | 33.07 | 39.27 | 43.37 | 41.97 | 36.34 |



Fig 2.    PSNR value of Haar and Daubechies wavelet transforms.



Fig 3.    PSNR of Symlet, Coiflet, biorthogonal, reverse biorthogonal, and Meyer wavelets.

| Original Image | Histogram - Original Image | Processed Image | Histogram - Processed Image |
|---|---|---|---|



(a)     Daubechies 24, PSNR: 32.54, MSSIM: 0.938, BPRI: 0.125, BMPRI: 57.90, size: 1125 x 1327 pixels



(b)     Reverse biorthogonal 3.1, PSNR: 28.96, MSSIM: 0.855, BPRI: 0.045, BMPRI: 17.96, size: 781 x 1001 pixels



(c)     Daubechies 24, PSNR: 24.32, MSSIM: 0.771, BPRI: 0.077, BMPRI: 40.64, size: 479 x 624 pixels



(d)     Biorthogonal 6.8, PSNR: 27.77, MSSIM: 0.927, BPRI: 0.091, BMPRI: 43.4, size: 505 x 672 pixels



(e)     Daubechies 20, PSNR: 29.29, MSSIM: 0.887, BPRI: 0.085, BMPRI: 43.42, size: 467 x 636 pixels



(f)     Biorthogonal 3.3, PSNR: 30.59, MSSIM: 0.821, BPRI: 0.072, BMPRI: 26.72, size: 472 x 635 pixels

Fig 4.     (a-f) Original and Processed image along with their histograms processed through selected wavelet transforms.

The highest PSNR in each of the six images is marked with bold letters. These correspond to Daubechies 24, reverse biorthogonal 3.1, Daubechies 24, biorthogonal 6.8, Daubechies 20, and biorthogonal 3.3 for successive images of Fig. 4(a)-(f). The specific wavelet is mentioned against each image in Fig. 4. The corresponding value of MSSIM, BPRI, and BMPRI are also highlighted with bold letters. The histograms of processed images are given on the right side in Fig. 4. It has been observed that the response to Haar wavelet, Symlet 1 wavelet, and the biorthogonal1.1 wavelet is same. Similarly, the performance of Daubechies 2, and Daubechies 3 are exactly same as Symlet 2 and Symlet 3, respectively. The

mean PSNR of six images is plotted in Fig. 2 and Fig. 3. The standard deviation of these images is also plotted along with the mean values in Fig. 2 and Fig. 3. The graphs in Fig. 2 corresponds to the Haar wavelet, and Daubechies wavelets. The mean value of the remaining wavelets is given Fig. 3. From Fig. 2, it is clear that variations among the values of Daubechies wavelets are within 1 dB value, but there is a definite increasing trend of PSNR as the value of N is increased from Daubechies 2 to Daubechies 20, and then it decreases slightly afterwards. This is understandable as with more number of taps, the amount of information increases until N becomes equal to 20, and then it slightly reduce as the undesired signal is possibly added in the extracted information.

The mean value of MSSIM is given in Table III. It is clear that the similarity index increases at the higher values of N. The mean values of BPRI, and BMPRI are given in Table IV, and Table V, respectively. The peak index values of BPRI and BMPRI for each image is encircled. It is clear that Daubechies 6 out performs in four out of six images. In the remaining two images, the value of Daubechies 6 closely follows two other wavelets Symlet 7, and Coiflet 2. The BPRI index of image (c) is same for Daubechies 6 and Symlet 7. The maximum value of BMPRI is not consistent, as the peak values of six images correspond to six different wavelets. The possible reason is that BMPRI compares the original image with a noisy image that has been generated by the aggregation of several noisy version of the original image. The image content has a larger role in generating the reference image, and therefore the response to various wavelets does not generate consistent results.

## V. Conclusions

This paper reviews the performance of seven discrete wavelet transforms in reducing the effect of speckle noise of the US medical images. The complete analysis is based on pre-processing, filtering, and post-processing. The pre-processing involves the application of Wiener filter, followed by filtering using seven discrete wavelet transforms. The selected wavelet transforms are Haar, Daubechies, Symlet, Coiflet, biorthogonal, reverse biorthogonal, and discrete Meyer wavelets. The post-processing involves the application of Median filter. The simulation results are based on US images of kidney that have previously been diagnosed for stones. The performance of various wavelets is compared using four metrics. Two are based on full-reference (FR) image quality assessment (IQA), and the remaining two are based on no-reference (NR) IQA. The two FR-IQA techniques are PSNR, and MSSIM, while the two NR-IQA techniques are blind pseudo-reference image (BPRI), and blind multiple PRI (BMPRI) metrics. The images with the highest PSNR are identified, and the other metrics of these images are marked. The output of the selected images with the highest PSNR, along with their histograms, are reproduced for comparison. The qualitative and quantitative analysis clearly show significant improvement in the processed images. In this paper only the first level of wavelet decomposition is considered. An extension of this work by considering multilevel wavelet decomposition has strong potential for better results. It would be interesting to see the effect of using one wavelet function for first level and another wavelet function for the second level.

### References

[1] C. D. Scales et. al., "Prevalence of Kidney Stones in the United States," European Urology, 2012, 62, pp. 160-165

[2] A. K. Jain, Fundamentals of Digital Image Processing, Prentice Hall, 1989

[3] C. Kao, X. Pan, E. Hiller, et. Al., "A Bayesian Approach for Edge Detection in Medical Ultrasound Images," IEEE Trans. On Nuclear Science, 1998, 45, (6), pp. 3089–3096.

[4] P. Perona, and J. Malik, "Scale-space and Edge Detection Using Anisotropic Diffusion," IEEE Trans. Pattern Analysis and Machine Intelligence, 1990, 12, (7), pp. 629-639

[5] S. Aja-Fernandez, and C. lberola-Lopez, "On the Estimation of the Coefficient of Variation for Anisotropic Diffusion Speckle Filtering," IEEE Trans. Image Processing, 2006, 15 (9), pp. 2694-2701

[6] K. Krissian, C. Westin, R. Kikinis, et. al., "Oriented Speckle Reducing Anisotropic Diffusion," IEEE Trans. Image Processing, 2007, 16 (5), pp. 1412-1424

[7] G. Vegas-Sanchez-Farrero, S. Aja-Fernandez, S. Martin-Fernandez, et. al., "Probabilistic-driven Oriented Speckle Reducing Anisotropic Diffusion with Application to Cardiac Ultrasound Images," Proc. Int. Conf. Beijing, China, Sep. 2010, pp. 518-525

[8] S. Finn, M. Glavin, and E. Jones, "Echocardiographic speckle reduction comparison," IEEE Trans. Ultrasonics, Ferroelectrics and Frequency Control, 2011, 58, (1), pp. 82-101

[9] S. G. Mallat, "Multifrequency Channel Decompositions of Images and Wavelet Models," IEEE Trans. Acoustics, Speech, and Signal Processing, 1989, 37, (12), pp. 2091-2110

[10] M. Antonini, M. Barlaud, and P. Mathieu, et. al., "Image Coding using Wavelet Transform," IEEE Trans. Image Processing, 1992, 1 (2), pp. 205-220

[11] G. Cincotti, G. Loi, and M. Pappalardo, "Frequency Decomposition and Compounding of Ultrasound Medical Images with Wavelet Packets," IEEE Trans. Medical Imaging, 2001, 20, (8), pp. 764-771

[12] M. I. H. Bhuiyan, M. O. Ahmad, and M. N. S. Swamy, "Spatially Adaptive Thresholding in Wavelet Domain for Despeckling of Ultrasound Images," IET Image Processing, 2009, 3 (3), pp. 147-162

[13] J. Tian, and L. Chen, "Image Despeckling using Non-parametric Statistical Model of Wavelet Coefficients," Biomedical Signal Processing and Control, 2011, 6, (4), pp.432-437

[14] D. Gupta, R. S. Anand, and B. Tyago, "Enhancement of Medical Ultrasound Images using Non-Linear Filtering Based on Rational-dilation Wavelet Transform," Proc. World Congress on Engineering and Computer Science, San Francisco, USA, Oct 2012

[15] S. Mukhopadhyay, and J. K. Manda, "Wavelet based Denoising of Medical Images using Sub-band Adaptive Thresholding through Genetic Algorithm," Procedia Technology, 2013, 10, pp. 680-689

[16] X. Fu, Y. Wang, L. Chen, et. al., "Quantum-inspired Hybrid Medical Ultrasound Images Despeckling Method," Electronic Letters, 2015, 51, (4), 2015, pp. 321-323

[17] H. Wen, and W. Qi, "Enhancement and Denoising Method of Medial Ultrasound Image Based on Wavelet Analysis and Fuzzy Theory," Proc. Int. Conf. Measuring Technology and Mechatronics Automation, 2015, China, pp. 448-452

[18] A. K. Yadav, A. P. Kumar, and S. K. Dhakad, "De-noising of Ultrasound Image using Discrete Wavelet Transform by Symlet Wavelet and Filters," Proc. Int. Conf. Advances in Computing, Communications and Informatics (ICACCI), 2015, India, pp. 1204-1208

[19] F. Adamo, F. Andria, A. M. Attivissimo, et. al., "A Comparative Study of Mother Wavelet Selection on Ultrasound Imaging Denoising," Measurement, 2013, 46, pp. 2447-2456

[20] I. Elamvazuthi, M. L. Zain, and K. M. Begam, "Despeckling of Ultrasound Images of Bone Fracture using Multiple Filtering Algorithms," Mathematical and Computer Modelling, 2013, 57, pp.152-168

[21] S. Sidhik, "Comparitive Study of Birge-Massart Strategy and Unimodal Thresholding for Image Compression using Wavelet Transform," Optik, 2015, 126, pp.5952-5955

[22] G. Wimmer, T. Tamaki, et. al., "Directional Wavelet Based Features for Colonic Polyp Classification," Medical Image Analysis, 2016, 31, pp.16-36

[23] Z. Wang, A. C. Bovik, H. R. Sheikh, et. al., "Image Quality Assessment: From Error Visibility to Structural Similarity," IEEE Trans. Image Processing, 2004, 13, (4), pp. 600-612

[24] X. Min, et. al., "Blind Quality Assessment based on Pseudo-Reference Image", IEEE Trans. Multimedia, 2018, 20, (8), pp. 2049-2062

[25] X. Min, et. al., "Blind Image Quality Estimation via Distortion Aggravation," IEEE Trans. Broadcasting, 2018, 64, (2), pp. 508-517

[26] Z. Wang, "The SSIM Index for Image Quality Assessment," http://www.cns.nvu.edu/~lcv/ssim/, 2014

[27] I. Daubechies, "Ten Lectures on Wavelets," CBMS-NSF Conference Series in Applied Mathematics SIAM Ed., 1992

[28] Ultrasound images, http://www.ultrasound-images.com

# BHA-160: Constructional Design of Hash Function based on NP-hard Problem

Ali AlShahrani

Faculty of Computing Studies
Arab Open University, Riyadh, Kingdom Saudi Arabia

*Abstract*—**Secure hash function is used to protect the integrity of the message transferred on the unsecured network. Changes on the bits of the sender's message are recognized by the message digest produced by the hash function. Hash function is mainly concerned with data integrity, where the data receiver needs to verify whether the message has been altered by eavesdropping by checking the hash value appended with the message. To achieve this purpose, we have to use a secure hash function that is able to calculate the hash value of any message. In this paper, we introduce an alternative hash function based on NP-hard problem. The chosen NP-hard problem is known as Braid Conjugacy problem. This problem has proved to be secure against cryptanalysis attacks.**

*Keywords*—*Hash function; integrity message; cryptanalysis; attack; NP-hard problem*

## I. INTRODUCTION

Hash function is the core of any cryptosystem. It is used for message integrity or for authenticating the data exchanging process between the connected parties. The design of a secure hash function consists of a special one-way function that receives any variable length input and produces a fixed length output. A one-way function is defined as a function that can simply take the input message and compute (generate) the corresponding hash value, but, it is computationally infeasible to recover the original message using the hash value. A hash function is called ideal if the hash value h cannot be distinguished from the values given by a random oracle [1]. Apart from hash functions, some cryptosystems are dependent on mathematically hard problems. An example of a mathematical hard problem is the braid theory. Generally, Braid Groups had been widely used as a tool to create various cryptographic primitives. There are a few of them such as a public key cryptosystem, key exchange, authentication and digital signature [2] [3]. Creating an ideal hash function using braid groups is connected to the general question of finding a function to map the braid groups to the sequence of {0,1}. The result of the secured hash function must be random enough and reveal no information about the argument of the hash function. The objective of this paper, therefore, is to create a secured hash function based on braid group's theory. The mechanism used in the core of this function is the braid multiplication, by which we multiply a pre-defined braid by the braid generated from message transformation (transformation of the message's content to a braid form). However, the importance of this research is related to the capability of our designed hash function to be attached to any cryptosystem for message integrity purposes with a high level of security. The rest of the paper is organized in six sections. The related works are discussed in Section 2. The proposed hash function is presented in Section 3. In Section 4, the algorithm performance is analyzed. The discussions and conclusion are presented in Sections 5 and 6, respectively.

## II. LITERATURE REVIEW

Hash functions have been applied for many security applications and protocols such as PGP, SSL, SSH, IPsec, TLS and S/MIME [4]. In order to provide these applications with a high level of security, we have to design a secure hash function against existing attacks. Let us now discuss the following scenario to understand the usage of a hash function: Alice wants to send a message $m$ to Bob. Alice needs to use the hash function $F_h$ to calculate the hash value $h$ of her message such that $F_h(m) = h$, and appends the hash value $h$ with the message. On the other hand, Bob (the receiver) needs to recalculate $h$ using the same hash function. By comparing the two hash values, Bob can judge whether the message has been altered or not. The message considered "Altered" if $F_h(m)_{\text{Alice}} \neq F_h(m)_{\text{Bob}}$.

The strength of any hash function can be measured by the complexity of its calculation and operation [5]. Recently, cryptosystems aimed to use some mathematical NP-hard problems in order to increase the complexity of their structure against the attackers. A problem is assigned to the NP (Nondeterministic Polynomial time) - hard problem class, if it is solvable in polynomial time by a nondeterministic oracle machine. Therefore, if we built a hash function based on a NP-hard problem, we will certain that the attackers cannot attack this function since it is based on a "hard-to solve" mathematical problem.

### A. Hash Function

A hash function $F_h$, is a transformation that takes an arbitrary size input $m$, and returns with a string of a fixed size, which is called the hash value $h$ (where $h = F_h(m)$) [6].

A cryptographically secure hash function should have the basic requirements in its design, which are:

- $F_h$ can be applied to an input of data of any size.

- $F_h$ produces a fixed-length of output.

- $F_h$ ($m$) is relatively easy to compute for any given $m$.

- $F_h(m)$ is one-way.

- $F_h(m)$ is collision-free.

MD2, MD5 and SHA [7] are good examples of well-known hash functions. In 1989, Ron Rivest introduced the MD2 Message Digest Algorithm that takes as input, a message of arbitrary length and produces as output, a 128-bit message digest by appending some redundancy to the message, and then iteratively applies 32 bytes to 16 bytes compression function. Researches done by [8] and [9] proved that the MD2 is not a one-way function, therefore, it is not collision-free and they also showed that it does not reach the ideal security level of $2^{128}$. However, the use of MD2 for new applications is discouraged. Similarly, MD5 takes as input, a message of arbitrary length and produces a 128-bit message digest; however, it is aimed at 32-bit machines instead of 8-bit machines in MD2. The algorithm consists of four distinct rounds with a similar structure, but each uses a different primitive logical function. According to the research done by [10] and [11], MD5 is not secure to be used in security applications since it is not collision-free. Therefore, MD5 in no longer recommended for new applications where collision-resistance is required. MD2 and MD5 are meant for digital signature applications where a large message has to be "compressed" in a secure manner. They are classified in a bit-operations based hash function category, since they depend on crossing, shifting and addition to the message's bits. The Secure Hash Algorithm (SHA-1) is another example of hash algorithms. It is one of the most widely used hash functions in the world. Indeed, four more variants have since been issued with increased output range and a slightly different design: SHA-224, SHA-256, SHA-384 and SHA-512 (sometimes they are collectively referred as SHA-2). However, SHA-1 takes a message with a maximum less than $2^{64}$ as an input producing a 160-bit message digest. The overall process of SHA-1 consists of five steps, starting from appending some of the padding bits to make the message congruent to a 448 modulo 512, ending with a 160-bit message digest. Through these steps, the message length must be appended to the message as well as XOR operations being applied to the message's bits. Research done by Chinese researchers showed that SHA-1 has been broken [12]. They presented new collision search attacks on the hash function SHA-1 and showed that the collision of SHA-1 can be found with a complexity of less than $2^{69}$.

### B. Braid Group

Number equations consecutively: Equation numbers, within parentheses, are to position flush right, as in (1), using a right the second category of our hash function's classification is the hash function based on the NP-hard problem. In this category, the heart of the hash function depends on a mathematical nondeterministic polynomial-time hard problem. Braid groups had been widely used as a tool to create various cryptographically primitives. There are a few of them, such as a public key cryptosystem, key exchange, authentication and digital signature. However, Conjugacy Problems are NP-hard problems in braid group theory. We say that braid $a$ and $b$ are conjugate if we have $a = s \, b \, s^{-1}$ for some braid $s$. Conjugacy Search Problem is one of the conjugacy problems in braid theory. This problem lies, that for some braids $(a,b) \in B_n$ X $B_n$

(where $B_n$ is braid group) such that $x$ and $y$ are conjugate, find $r \in B_n$ such that $b = rar^{-1}$. Any braid can be decomposed as a product of simple braids. One type of simple braid is the Artin generators $\sigma_i$, these have a single crossing between $i$-th and $(i+1)$-st strand as in Fig 1. Besides, the $n$-braid group $B_n$ can be presented by the Artin generators $\sigma_1, ..., \sigma_{n-1}$ and relations $\sigma_i \sigma_j = \sigma_j \sigma_i$ for $|i - j| > 1$ and $\sigma_i \sigma_j \sigma_i = \sigma_j \sigma_i \sigma_j$ for $|i - j| = 1$.



Fig. 1. Artin Generator $\sigma_i$

Many operations can be applied on two braids. For example, braid multiplication is the most used operation over braid. The multiplication of braids $a$ by $b$ where $(a,b) \in B_n$ results in a new braid which is unique. The process of ascertaining the original two braids (braid $a$ and $b$), given the resulted braid after the multiplication, is known to be a hard-problem. The multiplication of two braids is carried out by placing the braid $a$ under the braid $b$. As we previously mentioned, many cryptosystem's primitives have been built on braid group theory, but no hash function based on braid has been implemented yet. However, many researches are done in the braid group, and most of these researches showed the strength of this theory against attacks.

### III. PROPOSED BHA-160 HASH FUNCTUION ARCHITIECTURE

Currently, most of the existing hash functions are focusing on scrambling and shifting of the bits in the input blocks. With the intention of randomizing the bits of the input blocks, usually they are using the exclusive OR (XOR) operation and some additions in their implementation. For our work, we proposed a new approach of hash function architecture. In our opinion, hash function is not just scrambling or shifting the bits, but should also include the mathematical hard problems. We have found that the braid group's theory is the best way to do this, as it provides mathematical hard problems and also some advantages in computational aspects. The proposed structure consists of an initial vector called initial braid and blocks of text (represented as braid) to be the inputted into the hash function. We apply a braid operation (multiplication) on the braid groups to concentrate two different braids that then produce a completely new unique braid. By repeating this process, we will get a random braid that cannot be traced back to get the initial value of the hash function. This condition is able to fulfill the important properties of a secured hash function. The architecture of the proposed hash function will follow the steps as follows as in Fig. 2:

- Generate a random braid $B_{IV}$, to be as an initial vector of the hash function.

- Generate another braid by manipulating the bits from the text blocks.

- Do a multiplication operation on the initial braid and the braid generated beforehand.

- Repeat the iteration until the last text blocks.

### A. BHA-160 Processing Stages

The algorithm of BHA-160, takes as input a message of arbitrary length thereby producing as output, a 160-bit message digest. The input is processed in 192 bit blocks. The combined braid, as illustrated in the architecture, is achieved from the braid multiplication and will be processed in order to reduce the size of the digest to 160-bit. The overall processing of a message to produce a digest consists of four stages. The stages are:

- Stage 1: Append Padding Bits

The 192 bit block is padded to make sure the length is always in the desired length. The padding process is done by taking the first 8-bit block from the input message (which is less from the desired length) and then cyclic left shift the bits of the block by 2 bits as shown in Fig. 3.

After appending the padding bits, we will XOR every 8-bit in 192-bit block. The result of this stage is 12 8-bit blocks. Fig. 4 presents the input setup of BHA-160 in the first stage.



Fig. 2. Architecture of the Proposed Hash Function (BHA-160).



Fig. 3. Padding Process for BHA-160.



Fig. 4. Input Setup of BHA-160.

- Stage 2: Convert to Artin Generators, σ

This process is the beginning of mapping the input into a braid representation. As we can see in Fig. 5, B[i] represents the braid index or, in other words, the location where crossing occurs in braid groups. By mapping the input into a braid representation, we need to calculate the value of the crossing. With the number of strands $n = 128$, we convert the first 7 bits from binary to positive decimal. The 8th bit indicates the sign of the number that can be negative or positive. The positive sign indicates a positive crossing and the negative sign indicates a negative crossing.



Fig. 5. Relation between Blocks of Bits with Artin Generators.

- Stage 3: Braids Multiplication

The inputs of this stage are two braids with 12-byte size (12 crossing). The first braid represents the plain text block after the transformation (transforming the plain text block to Artin representation $B_i$. The second braid will be the initial value of $B_{IV}$ that is represented as braids $D_i$ with 24-byte ($B_{IV}$ that will be used for one time only, in the beginning of this stage). However, the initial value of $D_i$ will be reduced to 12-byte size to be multiplied by $B_i$. The braid reduction occurs by XORing $D_{2i-1}$ and $D_{2i}$ for all values of $1 \leq i \leq 12$. Therefore, the resulted braid, after reduction, is a 12-byte braid size which is represented as $D'_i$. The combined braid, resulted from multiplying the two braids $D'_i$ and $B_i$ as shown in Fig. 6 will therefore be in the size ranging from 0-24 crossing, since there is a possibility for zero crossing. However, the combined braid

then will be multiplied by the next input block after reduction to 12-byte instead of using the same value of $B_{IV}$.

- Stage 4: Message Digest Reduction and Production

This is the final stage where we produce the message digest of the corresponding plain text. However, this stage will be executed when we reach the last input block. The output will be in the size of 192-bit, meanwhile, we are looking for 160-bit size. Therefore, we will reduce the output size to 160-bit by keeping the first 160-bit and ignoring the remaining bits as portrait in Fig. 7.



X: Braid Multiplication

Fig. 6.    Braid Multiplication.



Fig. 7.    Message Digest Reduction.

## IV.   PERFORMANCE ANALYSIS

The performance of BHA-160 is examined against well-known hash functions, including: MD2, MD5, SHA-1, SHA-256, SHA-512. The experiment at each point is the mean of 10 measurements. The experiments of all hash functions are implemented on Core i7-4500U of a CPU of 2.4GHz, running Java 6 under Windows 7. The performance results are presented in Fig. 8.



Fig. 8.    Performance of Standard Hash Functions Against BHA-160.

The performance result shows that our BHA-160 is in the middle class, where it could outperform MD2 and it is almost close to the performance of other functions. However, we realized that a tradeoff between security and performance exists. Manipulating braids includes performing complex mathematical operations. In addition, the key size of BHA-160 is relatively larger than the key size used by most of the hash functions included in this study.

## V.   DISCUSSION

Two important parameters should be discussed, they are: the security and the performance of the proposed architecture. In terms of security, the 8-bit block of plain text will produce an Artin representation of a string in a 128 braid ($2^7$=128), which is big enough for security purposes, since the advice size for braid to be used for cryptography purposes is an 80 strings braid. Mathematically, the braid theory proved to be secure since it is virtually impossible to retrieve one of the multiplied braids after a braid multiplication operation. In terms of performance, a block of 192-bit will require one braid multiplication of two 128-strings braids, and 24 XOR operations. This can be considered a minimal operation that needs to be applied on every 192-bit plain text block.

## VI.   CONCLUSION

In conclusion, there is a presentation of less security level by bit-operations harsh functions that have a dependence on bits of XORing message as compared to the hash functions that are based on problems that are NP-hard. However, the proposed hash function is proved that it is secure due to the fact that it is based on mathematical problems that are hard to solve hence it is worth to be evaluated. The proposed hash function's internal stages tend to depend on the mapping of bits into braid representation as well as multiplying the resulted braids to each other. These stages form the core of the entire architecture of the hash function that is proposed and they could be able to fulfill the significant features of a hash function that is secured.

REFERENCES

[1] P. Hofmann and B. Schneier, Attacks on Cryptographic Hashes in Internet Protocols, facs.org. November 2005. [Online]. Available: https://tools.ietf.org/html/rfc4270 [Accessed May 20, 2019].

[2] K. Ko, S. Lee, J. Cheon, J. Han, J. Kang and C. Park, New Public-Key cryptosystem using braid groups, In advance in Cryptology: Crypto 2000.

[3] I. Al-Siaq, Public Key Cryptosystems based on Numerical Methods, Global Journal of Pure and Applied Mathematics, Vol.13, No.7 pp (2017) pp 3105-3112.

[4] R. Dobai, J. Korenek, L. Sekanina, Evolutionary design of hash function pairs for network filters, Applied Soft Computing, Volume 56, July 2017, Pages 173-181.

[5] G. Yu, Y. Zhao, C. Lu, J. Wang, HashGO: hashing gene ontology for protein function prediction, Computational Biology and Chemistry, Volume 71, December 2017, Pages 264-273.

[6] Y. Cui, J. Jiang, Z. Lai, Z. Hu, W. Wong, Supervised discrete discriminant hashing for image retrieval, Pattern Recognition, Volume 78, June 2018, Pages 79-90.

[7] W. Stalling, Cryptography and network security: principles and practices, Prentice Hall 2$^{nd}$ ED. 1999.

[8] N. Rogier and Chauvaud, The Compression Function of MD2 is not Collision Free, Selected Areas in Cryptography '95, 1995.

[9] F. Muller, The MD2 Hash Function is not One-Way, Advanced in Cryptology-Asia Crypt '2004, 2004.

[10] V. Klima, Finding MD5 Collision-a toy for a Notebook, Cryptology ePrint Archive, Report 2005/075.

[11] X. Wang, D. Feng, X. Lai, H. Yu, Collision for Hash Functions MD4, MD5, HAVAL-128 and RIPEMD, Cryptology ePrint Archive, Report 2004/199.

[12] X. Wang, Y. Yin, H. Yu, Finding Collisions in the Full SHA-1, Crypto 2005.

# Variable Reduction-based Prediction through Modified Genetic Algorithm

Allemar Jhone P. Delima[1], Ariel M. Sison[2], Ruji P. Medina[3]

Graduate Programs, Technological Institute of the Philippines, Quezon City, Philippines[1, 3]

Emilio Aguinaldo College, Manila, Philippines[2]

*Abstract*—Due to the massive influence in the use of prediction models in different sectors of society, many researchers have employed hybrid algorithms to increase the accuracy level of the prediction model. The literature suggests that the use of Genetic Algorithms (GAs) can sufficiently improve the performance of other prediction models; thus, this study. This paper introduced a new avenue of prediction integrating GA with the novel Inversed Bi-segmented Average Crossover (IBAX) operator paired with rank-based selection function to the KNN algorithm. The 70% of data from 597 records of student-respondents in the evaluation of the faculty instructional performance from the four State Universities and Colleges (SUC) in Caraga Region, Philippines were used as training set while the 30% was used for testing. The simulation result showed that the use of the proposed prediction model with the integration of the modified GA outperformed the KNN prediction model where GA with average crossover and roulette wheel selection function was used. The KNN where k value is three (3) was identified to be the optimal model for prediction with the 95.53% prediction accuracy compared to KNN with 1, 5, and 7 k values.

*Keywords*—*Enhanced prediction model; IBAX operator; modified genetic algorithm; prediction accuracy enhancement*

## I. INTRODUCTION

Data Mining (DM) is the process of extracting implicit information or knowledge from databases [1]-[3], that is drawn from the field of statistics [4] which uses mathematical and machine learning techniques and algorithms [5]. Knowledge Discovery in Databases (KDD) which is coined to data mining [6], represents the generally observed process in knowledge discovery where knowledge is the result of the data-driven discovery while data mining being the observed step in the process for efficiently automated discovery, employs diverse approaches of DM analysis [7].

The field of data mining has become standard practice in various disciplines such as business, finance, and marketing allowing to inadvertently impact social sciences and humanities in general [8]. The range of its application has also reached other sectors such as education [9], [10] and healthcare [11], [12]. DM is promising for researches applied in engineering, biomedical sciences, medical systems, web, sports, and shared market because of the accessibility to various vast datasets [13], [14].

There are several widely accepted major functions in data mining found in the literature such as association, classification, clustering, estimation, and prediction [15].

Prediction, as one of the optimal data mining approach, was defined by [16] as "a powerful tool in the process of planning that can provide the decision maker with a prediction about the future events according to using experiences and applying statistical, mathematical, or computational methods." It is commonly used in educational data mining (EDM) [17]-[19], crime mining [20]-[22], business and finance [23], [24], health [25], [26], and more.

Data preprocessing is one of the essential methods that are useful in data mining. It has led to the enhancement of the quality of data and improved the precision and accuracy level of a prediction model [22], [27]. Data reduction, as an important data preprocessing technique in DM, is performed through the selection and removal of the unneeded attributes in the dataset [28]. Reducing the training set or variables and retaining the most representative data is advisable. The goal is also to obtain nearly the same outcome or data-driven output [29]-[31]. Minimizing the size of the dataset aids in maximized accuracy [32], [28]. One of the widely used data reduction methods is the Genetic Algorithm [33] which was introduced by J.H. Holland. The average crossover, which is one of the crossover operators of GA, is modified in this study.

Due to the massive influence in the use of prediction methods in diverse fields such as weather and natural calamity, stock markets, telecommunication, transport organization, energy, economy and other sectors [16], researchers have employed models integrating algorithms for prediction as well as hybridizing algorithms and combining different techniques to elevate the accuracy level of a prediction model. To name some, the study of [34] employed a hybrid feature selection method integrating Weight by Relief and GA to select the best features in the dataset for myocardial infraction prediction using J48. An accuracy of 82.67% was depicted after applying the model to the imbalanced dataset. Also, a study of [35] used the K-Means segmentation technique and C4.5 algorithm to build a prediction model for customer loyalty in a multimedia service provider. The integration of K-Means and C4.5 algorithm have yielded an increase of 79.33% accuracy prediction from the identified 69.23% accuracy with the C4.5 algorithm alone.

Lastly, the prediction model of [36] used the K-Nearest Neighbor (KNN) algorithm to predict standard levels of OTOP's (One Tambon One Product) wood handicrafts product. Results showed that the model obtained the best prediction at the accuracy of 87.73%. The KNN algorithm is susceptible to noise and sensitive to irrelevant features [37]. Even though the prediction rate using KNN is already

acceptable, but with the advent of combining genetic algorithm for variable reduction to address the problem of KNN, an increase of accuracy through the hybridization is hoped to be established.

With the advent of combining two or more models, an increase of prediction accuracy is evident [38] such that of [34] that obtained 82.67% and [35] with 79.33% after integrating hybridization than employing prediction with one algorithm alone.

Therefore, the quest of this study is not only to modify the genetic algorithm and introduce a new avenue of crossover mating scheme but also to increase the accuracy of the prediction model of [36] who used KNN algorithm through the integration of the modified genetic algorithm for feature selection and variable minimization before prediction. The rest of the paper is arranged as follows: Section II discusses the literature review of genetic algorithm and other prediction models. Section III includes the design and methodology used in the study. Section IV discusses the results and discussions while Section V highlights the conclusion and recommendation.

## II. Literature Review

### A. Genetic Algorithm

The genetic algorithm is one of the many evolutionary algorithms anchored on the biological adaptation in the quest for global optimization. GA is deliberately one of the famous technique used in the search for the optimal solution for problems with a large search space. GA produces and controls some individuals by assigning optimal operators on its three fundamental operations namely the selection, crossover, and mutation functions. In this study, a modified genetic algorithm with the integration of novel Inversed Bi-segmented Average Crossover (IBAX) is used. This novel crossover is a modified version of the traditional average crossover of GA.

### B. Genetic Algorithm-based Prediction Models

The literature suggests that the genetic algorithm can efficiently increase the performance of other prediction models [39], [40]. The most significant benefit of the genetic algorithm is its ability to avoid being confined in local optima, and the use of GA or a hybrid GA gives the chance to select the best appropriate objective functions freely [41].

A recent study improved the accuracy of the self-organizing map (SOM), a type of unsupervised ANN, in predicting robotic manipulation failures for force-sensitive tasks using a genetic algorithm. The proposed hybrid GA-SOM model exhibited an increased accuracy prediction and improved the predictive capability of the SOM algorithm when used alone [42]. Moreover, the use of evolutionary technique like the genetic algorithm in enhancing ANN was observed along with SVM-Linear (L), SVM-Polynomial (P), SVM-Radial Basis Function (RBF), and CART in predicting the shape of carbon black reinforced rubbers. With the advent of the genetic algorithm, the prediction accuracy of each model has increased, and the most accurate model was obtained using GA-ANN hybrid model with those obtained using the GA-CART, GA-SVML, GA-SVM-P, and GA-SVM-RBF [43].

TABLE I. Indexed GA-Based Prediction Models

| Algorithms/ Procedure | Authors/ Year | Purpose | Significant Results |
|---|---|---|---|
| Genetic Algorithm-based Self Organizing Map (GA-SOM) prediction model | Parisi & RaviChandran, (2018) | To enhance the performance of SOM using GA in predicting robotic failures | The hybrid GA-SOM yielded 91.95% prediction accuracy compared with the 84.96% prediction of the standalone SOM algorithm |
| Hybrid GA-ANN, GA-CART, GA-SVML, GA-SVMP, GA-SVM-RBF prediction model | Martinez et al., (2018) | To use GA, ANN, SVM, and CART to characterize rubber blends. | The GA-ANN model exhibited the finest classification accuracy of 75.75% improving the 74.80% accuracy attained without GA. |
| Genetic Algorithm-based Back Propagation (GA-BP) neural network prediction model | Zheng, Qian, Liu, & Liu, (2018) | Hybrid GA-BP neural network was used to model skid resistance of epoxy asphalt mixture | The optimized GA-BP neural network hybrid model was able to give an effective and accurate forecast of long term skid resistance with 99% accuracy. |
| Combination of Genetic Algorithm, Levenberg-Marquardt algorithm, and Back Propagation neural network as a prediction model | Zhou et al., (2018) | Application of GA-LM-BP Neural network in fault prediction of drying furnace equipment. | The GA-LM-BP hybrid prediction model obtained the decision coefficient $R^2$ of 0.97511 which is higher than the BP and GA-BP models. |
| The use of Genetic algorithm in least squares-support vector machine (LS-SVM), Back Propagation Neural Network (BPNN), and Random Forest (RF) | Liu et al., (2018) | Analyze the origin of extra virgin olive oils. | Simulation results showed that GA-LS-SVM model obtained 96.25% prediction accuracy and a prediction of 86.25% for GA-BPNN while a prediction accuracy of 82.5% for GA-RF was identified. |
| Genetic Algorithm-based Random Forest (RF) prediction model | Kumar & Sahoo, (2017) | To propose a hybrid GA-RF prediction model for cardiovascular disease diagnosis | Hybrid GA-RF prediction model outperformed the principal component analysis-based random forest (PCA-RF), Relief F-based random forest (Relief-F-RF), sequential forward floating search-based random forest (SFFS-RF), and sequential backward floating search-based random forest (SFBS-RF) having 93.2%, 84.8%, 85.4%, 79.1%, and 85.8% prediction accuracy, respectively. |
| Genetic Algorithm-based Artificial Neural Network (GA-ANN) prediction model | Armaghani et al., (2016) | To enhance the prediction rate of ANN in predicting AoP from blasting operation in granite quarry site. | GA-ANN model obtained 0.965 coefficient of determination, variance account for (VAF) value of 96.380 and RMSE of 0.049 than the ANN with those statistical function values of 0.857, 84.257, and 0.117 respectively. |

Another study used the hybrid GA-BP neural network in predicting long-term skid resistance of epoxy asphalt mixture. The GA-BP model produced a great accuracy result when tested using the training set, validation set, and test set [44]. Meanwhile, the application of genetic algorithm, Levenberg-Marquardt (LM) algorithm, and backpropagation neural network were observed in fault prediction of drying furnace equipment. The hybrid GA-LM-BP model showed an increased prediction accuracy compared to both BP neural network and GA-BP neural network models [45].

Further, the hybrid genetic algorithm-based least squares-support vector machine (GA-LS-SVM), genetic algorithm-based back propagation neural network (GA-BPNN), and genetic algorithm-based random forest (GA-RF) were employed in identifying the topographical origin of extra-virgin olive oils. The simulation results showed that GA-LS-SVM obtained the highest prediction accuracy for features selection methods compared to GA-BPNN and GA-RF models [46]. To further prove the superiority of GA as variable minimization algorithm, the genetic algorithm was used to perform feature selection where the extracted features are taken as an input to random forest (RF) classifier in accomplishing cardiovascular diagnostic problem. The outcome shows that the GA-RF model obtained the highest prediction accuracy rate when compared to other feature selection algorithms [47].

Lastly, an artificial neural network (ANN) and genetic algorithm-based ANN (GA-ANN) were proposed and evaluated to predict air overpressure from blasting operation in a granite quarry site in Penang, Malaysia. Simulation results proved the superiority of GA-ANN model in predicting air overpressure than using ANN algorithm alone [39]. The indexed GA-based prediction models are shown in Table I.

## III. Methodology

### A. Modified Genetic Algorithm for Variable Reduction

To achieve the purpose of the study, the average crossover which is one of the crossover operators in the genetic algorithm as shown in Fig. 1, is modified. The modified crossover will be called Inversed Bi-segmented Average Crossover (IBAX) as depicted in Fig. 2. The use of rank-based selection function was observed in the simulation process.



Fig. 1.　Average Crossover with Roulette Wheel Selection Function.



Fig. 2.　Inversed Bi-Segmented Average Crossover with Rank-Based Selection Function.

For the IBAX operator to be realized, the following steps must be executed:

Step 1: Take the parents from the selection pool.

Step 2: Count the number of genes found in the chromosomes. Identify if the dataset is in odd or even numbers.

Step 3: Segment the chromosomes (x and y) by dividing the total number of genes in the chromosomes into two and make sure that both first and second segments must contain an equal number of genes in an even count.

Step 4: On the first segment, create offspring Z for each gene by inversely pairing the first gene from chromosome X to the last gene on chromosome Y. Repeat until the last gene of the chromosome X and the first gene of the chromosome Y have inversely mated and have produced an offspring using the formula:

$$z = [x + y] / 2 \qquad (1)$$

Step 5: Execute the same process on the second segment until genes from all segments have produced offspring. In case of odd datasets, the last genes of the chromosomes will not be combined in the second segment and will automatically be mated with each other to produce offspring.

### B. K-Nearest Neighbor Algorithm

Another recognized data mining algorithm for classification and prediction introduced by Fix and Hodges is the k-Nearest Neighbor (k-NN). This method adopts instance-based learning for prediction. The famous classifier is known as a non-parametric algorithm since it does not produce assumptions on the input data distribution; therefore, it is widely used in various applications [48], [49]. K-Nearest Neighbor (KNN) algorithm is simple and can be implemented through the following steps:

Step 1: Assign k values of the nearest neighbor of an instance in the algorithm.

Step 2: Perform the Euclidian distance calculation of each instance.

Step 3: Choose K neighboring attributes that have the lowest Euclidian distance.

Fig. 3.    Conceptual Framework of the Study.

A prediction model using K-Nearest Neighbor (KNN) algorithm was utilized in the study of [36] along with the many studies found in the literature.

### C. Enhanced KNN Prediction Model

The study evaluated the accuracy level of [36] prediction model when integrated with GA having AX operator and with the modified GA with IBAX operator having 1, 3, 5, and 7 k values. The Waikato Environment for Knowledge Analysis (WEKA) version 3.8.2 was instrumental in the simulation of KNN prediction model. The simulation results of both existing and enhanced prediction models were compared to check the improvement rate of the accuracy level of the prediction model. The conceptual framework of the study is presented in Fig. 3.

### D. Datasets

The datasets used in this study were the 597 records of student-respondents in the evaluation of the faculty instructional performance from the four State Universities and Colleges (SUC) in Caraga Region, Philippines. The thirty (30) variables that represent the faculty instructional performance (IP) having divided into six (6) parts viz., methodology, classroom management, student discipline, assessment of learning, student-teacher relationship, and peer relationship are reduced before the prediction to aid maximized accuracy. The 70% of the data were used as the training set while the remaining 30% were used for testing.

### E. Prediction Evaluation

An optimal model is selected once the model with the highest prediction rate is identified granted that the model has the lowest root mean squared error and mean absolute error values. Countless forecasting and prediction models found in the literature are evaluated using the various forecast error statistical tools. The following tools listed below will be used along with Precision, Recall, and F-Measure:

Root Mean Squared Error (RMSE)

$$R.M.S.E. = \sqrt{\sum_{t=T+1}^{T+h}(\hat{y}_t - y_t)^2 / h} \qquad (2)$$

Mean Absolute Error (MAE)

$$M.A.E. = \sum_{t=T+1}^{T+h} |\hat{y}_t - y_t| / h \qquad (3)$$

## IV.   RESULTS AND DISCUSSION

### A. Variable Minimization using GA with AX and IBAX Operators

The simulation on the genetic algorithm was done for ten generations utilizing the existing traditional average crossover and roulette wheel selection function. To generate new offspring from the two chromosomes (IP and Y), the average crossover was used where the average of the two chromosomes/parents was calculated. The new fitness values are then calculated based on the new offspring produced after the crossover function. Variables having the lowest fitness value were removed from the dataset. The sample simulation on the genetic algorithm having the original AX operator and roulette wheel selection function is presented in Table II.

First Generation: Variable C2 is removed from the chromosome since it obtained the lowest fitness value of 171396 as evident in Table II.

On the other hand, the simulation on the genetic algorithm with the novel Inversed Bi-segmented Average Crossover (IBAX) operator and rank-based selection function was done on the same dataset and number of generations.

First Generation: Variable C2 was removed from the list of variables after applying the rank-based selection. The variable C2 obtained the lowest fitness value in the rank-based selection. Hence, it does not have any chance to be selected. Moreover, after applying the inversed bi-segmented average crossover (IBAX) operator and obtained the fitness value of the offspring, variable C3 was removed from the chromosomes since it obtained the lowest fitness value of 224676 that will not warrant for the next generation. Thus, in the first generation, there were two variables removed from the list as shown in Table III.

Prior to prediction, the variables were minimized using GA with AX operator having roulette wheel selection function and GA with IBAX operator having rank-based selection function performed for ten generations.

TABLE II.    GENERATION 1 USING AN AX OPERATOR WITH RWS FUNCTION

| IP | X | Fitness | Rank | Pool Y | Pool IP | Off-spring | Fitness | Decision |
|----|---|---------|------|--------|---------|------------|---------|----------|
| M1 | 546 | 298116 | 22 | 552 | M5 | 549 | 301401 | |
| M2 | 565 | 319225 | 30 | 565 | M2 | 565 | 319225 | |
| M3 | 558 | 311364 | 27 | 548 | SD1 | 553 | 305809 | |
| M4 | 559 | 312481 | 28 | 546 | SD5 | 552.5 | 305256.3 | |
| M5 | 552 | 304704 | 24 | 474 | C3 | 513 | 263169 | |
| C1 | 490 | 240100 | 3 | 546 | A3 | 518 | 268324 | |
| C2 | 354 | 125316 | 1 | 474 | C3 | 414 | 171396 | Remove |
| C3 | 474 | 224676 | 2 | 556 | ST1 | 515 | 265225 | |
| C4 | 542 | 293764 | 18 | 490 | C1 | 516 | 266256 | |
| C5 | 528 | 278784 | 12 | 531 | ST3 | 529.5 | 280370.3 | |
| SD1 | 548 | 300304 | 23 | 500 | A2 | 524 | 274576 | |
| SD2 | 512 | 262144 | 5 | 542 | C4 | 527 | 277729 | |
| SD3 | 565 | 319225 | 29 | 546 | A3 | 555.5 | 308580.3 | |
| SD4 | 556 | 309136 | 26 | 558 | M3 | 557 | 310249 | |
| SD5 | 546 | 298116 | 21 | 528 | C5 | 537 | 288369 | |
| A1 | 513 | 263169 | 6 | 534 | ST2 | 523.5 | 274052.3 | |
| A2 | 500 | 250000 | 4 | 526 | P3 | 513 | 263169 | |
| A3 | 546 | 298116 | 20 | 513 | A1 | 529.5 | 280370.3 | |
| A4 | 518 | 268324 | 8 | 565 | M2 | 541.5 | 293222.3 | |
| A5 | 516 | 266256 | 7 | 516 | A5 | 516 | 266256 | |
| ST1 | 556 | 309136 | 25 | 556 | ST1 | 556 | 309136 | |
| ST2 | 534 | 285156 | 16 | 559 | M4 | 546.5 | 298662.3 | |
| ST3 | 531 | 281961 | 14 | 546 | M1 | 538.5 | 289982.3 | |
| ST4 | 541 | 292681 | 17 | 556 | SD4 | 548.5 | 300852.3 | |
| ST5 | 527 | 277729 | 11 | 552 | M5 | 539.5 | 291060.3 | |
| P1 | 531 | 281961 | 13 | 565 | SD3 | 548 | 300304 | |
| P2 | 533 | 284089 | 15 | 541 | ST4 | 537 | 288369 | |
| P3 | 526 | 276676 | 10 | 518 | A4 | 522 | 272484 | |
| P4 | 526 | 276676 | 9 | 565 | SD3 | 545.5 | 297570.3 | |
| P5 | 544 | 295936 | 19 | 565 | SD3 | 554.5 | 307470.3 | |

The variable minimization process using the genetic algorithm with AX operator and RWS function has depicted a decrease after the ten generations. From the 30 variables, it was minimized to 17 with a total reduction of 43%. Meanwhile, the variable minimization process using the genetic algorithm with the proposed novel mating scheme called inversed bi-segmented average crossover operator, and rank-based selection function has depicted a noticeable decrease after the ten generations. From the 30 variables, the numbers were minimized to 10 variables after the generations. A total of 66.66% of variables were removed as depicted in Table IV.

The simulation result showed that the modified genetic algorithm with a new crossover mating scheme outperformed the average crossover of genetic algorithm in reducing variables prior to prediction. Since dropping one or more variables helps reduce dimensionality, predictions using the dataset having 17 and 10 variables were conducted using the KNN algorithm.

Meanwhile, in the extent of fitness function, the proposed IBAX operator of the genetic algorithm has increased and outperformed the rate of the fitness functions generated using the genetic algorithm with the existing AX operator. The variables who obtained the lowest fitness function in each generation for ten generations were removed. It is evident in Fig. 4 and Table V that the fitness functions that were removed using the new crossover operator is higher compared to the fitness functions that were removed using the existing average crossover. This denotes that the modified genetic algorithm has managed to increase the fitness function of the variables compared to the genetic algorithm with traditional AX operator.

TABLE III.    GENERATION 1 USING IBAX WITH THE RANK-BASED SELECTION FUNCTION

| IP | X | Fitness | Rank-based Rank | Rank-based New Fitness | IBAX Parent 1 | IBAX Parent 2 | IBAX Offspring | Fitness |
|----|---|---------|------|-------------|----------|----------|-----------|---------|
| M2 | 565 | 319225 | 30 | 3629986.8 | 565 | 541 | 553 | 305809 |
| SD3 | 565 | 319225 | 29 | 3504670.8 | 565 | 542 | 553.5 | 306362.3 |
| M4 | 559 | 312481 | 28 | 3379354.8 | 559 | 544 | 551.5 | 304152.3 |
| M3 | 558 | 311364 | 27 | 3254038.8 | 558 | 546 | 552 | 304704 |
| SD4 | 556 | 309136 | 26 | 3128722.8 | 556 | 546 | 551 | 303601 |
| ST1 | 556 | 309136 | 25 | 3003406.8 | 556 | 546 | 551 | 303601 |
| M5 | 552 | 304704 | 24 | 2878090.8 | 552 | 548 | 550 | 302500 |
| SD1 | 548 | 300304 | 23 | 2752774.8 | 548 | 552 | 550 | 302500 |
| M1 | 546 | 298116 | 22 | 2627458.8 | 546 | 556 | 551 | 303601 |
| SD5 | 546 | 298116 | 21 | 2502142.8 | 546 | 556 | 551 | 303601 |
| A3 | 546 | 298116 | 20 | 2376826.8 | 546 | 558 | 552 | 304704 |
| P5 | 544 | 295936 | 19 | 2251510.8 | 544 | 559 | 551.5 | 304152.3 |
| C4 | 542 | 293764 | 18 | 2126194.8 | 542 | 565 | 553.5 | 306362.3 |
| ST4 | 541 | 292681 | 17 | 2000878.8 | 541 | 565 | 553 | 305809 |
| ST2 | 534 | 285156 | 16 | 1875562.8 | 534 | 490 | 512 | 262144 |
| P2 | 533 | 284089 | 15 | 1750246.8 | 533 | 500 | 516.5 | 266772.3 |
| ST3 | 531 | 281961 | 14 | 1624930.8 | 531 | 512 | 521.5 | 271962.3 |
| P1 | 531 | 281961 | 13 | 1499614.8 | 531 | 513 | 522 | 272484 |
| C5 | 528 | 278784 | 12 | 1374298.8 | 528 | 516 | 522 | 272484 |
| ST5 | 527 | 277729 | 11 | 1248982.8 | 527 | 518 | 522.5 | 273006.3 |
| P3 | 526 | 276676 | 10 | 1123666.8 | 526 | 526 | 526 | 276676 |
| P4 | 526 | 276676 | 9 | 998350.8 | 526 | 526 | 526 | 276676 |
| A4 | 518 | 268324 | 8 | 873034.8 | 518 | 527 | 522.5 | 273006.3 |
| A5 | 516 | 266256 | 7 | 747718.8 | 516 | 528 | 522 | 272484 |
| A1 | 513 | 263169 | 6 | 622402.8 | 513 | 531 | 522 | 272484 |
| SD2 | 512 | 262144 | 5 | 497086.8 | 512 | 531 | 521.5 | 271962.3 |
| A2 | 500 | 250000 | 4 | 371770.8 | 500 | 533 | 516.5 | 266772.3 |
| C1 | 490 | 240100 | 3 | 246454.8 | 490 | 534 | 512 | 262144 |
| C3 | 474 | 224676 | 2 | 121138.8 | 474 | 474 | 474 | 224676 |
| C2 | 354 | 125316 | 1 | -4177.2 | | | | |

TABLE IV.    VARIABLE MINIMIZATION SIMULATION RESULT FOR GENETIC ALGORITHMS WITH AX AND IBAX OPERATORS

| Basic GA with AX Operator | | | | | Proposed GA with IBAX operator | | | | |
| Number of Generations | Number of Variables Left | Number of Variables Removed | Variables Removed | Percentage | Number of Generations | Number of Variables Left | Number of Variables Removed | Variables Removed | Percentage |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 30 | 1 | C2 | 3.33% | 1 | 30 | 2 | C3, C2 | 6.66% |
| 2 | 29 | 2 | M5, A2 | 6.66% | 2 | 28 | 2 | ST2, C1 | 6.66% |
| 3 | 27 | 1 | C3 | 3.33% | 3 | 26 | 2 | ST4, A2 | 6.66% |
| 4 | 26 | 2 | C4, A5 | 6.66% | 4 | 24 | 2 | P2, A1 | 6.66% |
| 5 | 24 | 1 | C1 | 3.33% | 5 | 22 | 2 | ST3, A4 | 6.66% |
| 6 | 23 | 1 | P3 | 3.33% | 6 | 20 | 2 | C5, ST5 | 6.66% |
| 7 | 22 | 1 | A1 | 3.33% | 7 | 18 | 2 | SD5, SD2 | 6.66% |
| 8 | 21 | 1 | SD1 | 3.33% | 8 | 16 | 2 | M1, A5 | 6.66% |
| 9 | 20 | 1 | SD2 | 3.33% | 9 | 14 | 2 | A3, P3 | 6.66% |
| 10 | 19 | 2 | C5, A3 | 6.66% | 10 | 12 | 2 | M5, P4 | 6.66% |
| 10 | 17 | - | - | - | 10 | 10 | - | - | - |
| Total Percentage of Variables Removed | | | | 43.33% | Total Percentage of Variables Removed | | | | 66.66% |

# Fitness Function Evaluation



Fig. 4.    Comparison of the Fitness Function of the Removed Variables in Every Generation.

## B. Prediction Model Accuracy Evaluation

To evaluate the accuracy level of KNN as a prediction model, thirty percent (30%) of the data were used for testing while seventy percent (70%) were used as the training set. Table VI shows the comparison of results when GA with AX operator and roulette wheel selection function is used, and GA with IBAX operator with rank-based selection function are integrated prior to the prediction using KNN. The predictive capability of the KNN algorithm was also tested without the variable reduction stage and obtained a 90.50% prediction accuracy rate with a k value of 1.

The results showed that the prediction model gained an increase in the accuracy when integrated with genetic algorithm especially with the modified GA. The optimal model for predicting the instructional performance of the faculty in the four SUCs in the Caraga Region, Philippines is the KNN with a k value of 3. The model obtained a high 95.53% prediction accuracy. Meanwhile, the second best model that

has 94.97% accuracy is the KNN with k=5 where its MAE and RMSE values are 0.08 and 0.20, respectively.

TABLE V.    VALUE OF THE FITNESS FUNCTIONS REMOVED IN EVERY GENERATION

| Number of Generations | AX Operator | IBAX Operator |
|---|---|---|
| 1 | 171396 | 224676 |
| 2 | 263169 | 262144 |
| 3 | 265225 | 270920.3 |
| 4 | 266256 | 273529 |
| 5 | 268324 | 275100.3 |
| 6 | 272484 | 278256.3 |
| 7 | 274052.3 | 279841 |
| 8 | 274576 | 281961 |
| 9 | 277729 | 287296 |
| 10 | 280370.3 | 290521 |

TABLE VI.    INDEXED KNN AND GA-BASED KNN PREDICTION MODELS

| Model | K Value | Accuracy % | MAE | RMSE | Precision | Recall | F- Measure |
|---|---|---|---|---|---|---|---|
| KNN algorithm alone | 1 | 90.5028% | 0.0978 | 0.3055 | 0.907 | 0.905 | 0.906 |
| | 3 | 87.1508% | 0.1379 | 0.3241 | 0.872 | 0.872 | 0.872 |
| | 5 | 84.3575% | 0.1587 | 0.3311 | 0.840 | 0.844 | 0.841 |
| | 7 | 84.9162% | 0.1649 | 0.321 | 0.846 | 0.849 | 0.846 |
| GA-KNN with original AX and Roulette Wheel Selection | 1 | 93.8547% | 0.069 | 0.2434 | 0.938 | 0.939 | 0.938 |
| | 3 | 89.9441% | 0.1196 | 0.261 | 0.899 | 0.899 | 0.899 |
| | 5 | 89.9441% | 0.1293 | 0.2557 | 0.899 | 0.899 | 0.899 |
| | 7 | 89.3855% | 0.1377 | 0.2652 | 0.893 | 0.894 | 0.894 |
| GA-KNN with new IBAX and Rank Based Selection | 1 | 94.4134% | 0.0586 | 0.2364 | 0.944 | 0.944 | 0.943 |
| | **3** | **95.5307%** | **0.0662** | **0.1963** | **0.956** | **0.955** | **0.955** |
| | 5 | 94.9721% | 0.0828 | 0.2067 | 0.953 | 0.950 | 0.948 |
| | 7 | 94.9721% | 0.0956 | 0.2175 | 0.953 | 0.950 | 0.948 |

## V.    CONCLUSION AND RECOMMENDATION

With the integration of the genetic algorithm, the prediction model using the KNN algorithm has increased its prediction accuracy. The modified genetic algorithm with a new crossover mating scheme called Inversed Bi-segmented Average Crossover (IBAX) showed a considerably high prediction percentage than the genetic algorithm with average crossover having the roulette wheel as the selection function. Along with the GA-based prediction models found in the literature, the enhancement on the KNN as prediction model integrated with the modified genetic algorithm was a success and is added to the body of knowledge. Future researchers may consider using the modified GA-based KNN on different datasets as a prediction model.

## REFERENCES

[1] Savaliya, A. Bhatia, and J. Bhatia, "Application of Data Mining Techniques in IoT: A Short Review," Int. J. Sci. Res. Sci. Eng. Technol., vol. 4, no. 2, pp. 218–223, 2018.

[2] C. Suresh, K. T. Reddy, and N. Sweta, "A Hybrid Approach for Detecting Suspicious Accounts in Money Laundering Using Data Mining Techniques," Int. J. Inf. Technol. Comput. Sci., vol. 8, no. 5, pp. 37–43, 2016.

[3] K. Rajalakshmi, S. S. Dhenakaran, and N. Roobini, "Comparative Analysis of K-Means Algorithm in Disease Prediction," Int. J. Sci. Eng. Technol. Res., vol. 4, no. 7, pp. 2697–2699, 2015.

[4] I. A. Khan and J. T. Choi, "An application of educational data mining (EDM) technique for scholarship prediction," Int. J. Softw. Eng. its Appl., vol. 8, no. 12, pp. 31–42, 2014.

[5] E. Sugiyarti, K. A. Jasmi, B. Basiron, M. Huda, K. Shankar, and A. Maseleno, "Decision support system of scholarship grantee selection using data mining," Int. J. Pure Appl. Math., vol. 119, no. 15, pp. 2239–2249, 2018.

[6] E. Susnea, "Using data mining techniques in higher education," in The 4th International COnference on Virtual Learning ICVL 2009, 2009, vol. 1, no. 1, pp. 371–375.

[7] E. Petrova, P. Pauwels, K. Svidt, and R. L. Jensen, "Advances in Informatics and Computing in Civil and Construction Engineering," in 35th International Council for Research and Innovation in Building Construction W78 2018 Conference, 2019, pp. 19–26.

[8] R. Kitchin, "Big Data, new epistemologies and paradigm shifts," Big Data Soc., vol. 1, no. 1, pp. 1–12, 2014.

[9] B. M. M. Alom and M. Courtney, "Educational Data Mining: A Case Study Perspectives from Primary to University Education in Australia," Int. J. Inf. Technol. Comput. Sci., vol. 10, no. 2, pp. 1–9, 2018.

[10] A. Ikhwan et al., "A Novelty of Data Mining for Promoting Education Based on FP-Growth Algorithm," Int. J. Civ. Eng. Technol., vol. 9, no. 7, pp. 1660–1669, 2018.

[11] M. M. Malik, S. Abdallah, and M. Ala'raj, "Data mining and predictive analytics applications for the delivery of healthcare services: a systematic literature review," Ann. Oper. Res., vol. 270, no. 1–2, pp. 287–312, 2018.

[12] R. Wadhawan, "Prediction of Coronary Heart Disease Using Apriori algorithm with Data Mining Classification," Int. J. Res. Sci. Technol., vol. 3, no. 1, pp. 1–15, 2018.

[13] S. A. Aljawarneh, O. Bayat, and M. Essaaidi, "Introduction to the special section on new trends in data mining, games engineering and database systems," Comput. Electr. Eng., vol. 66, pp. 420–422, 2018.

[14] P. Kaur, M. Singh, and G. S. Josan, "Classification and Prediction Based Data Mining Algorithms to Predict Slow Learners in Education Sector," Procedia Comput. Sci., vol. 57, pp. 500–508, 2015.

[15] M. Brilliant, DwiHandoko, and Sriyanto, "Implementation of Data Mining Using Association Rules for Transactional Data Analysis," 3rd Int. Conf. Inf. Technol. Bus., pp. 177–180, 2017.

[16] M. J. Rezaee, M. Jozmaleki, and M. Valipour, "Integrating dynamic fuzzy C-means , data envelopment analysis and artificial neural network to online prediction performance of companies in stock exchange," Physica A, vol. 489, pp. 78–93, 2018.

[17] M. Zaffar, M. Ahmed, K. S. Savita, and S. Sajjad, "A Study of Feature Selection Algorithms for Predicting Students Academic Performance," Int. J. Adv. Comput. Sci. Appl., vol. 9, no. 5, pp. 541–549, 2018.

[18] E. Fernandes, M. Holanda, M. Victorino, V. Borges, R. Carvalho, and G. Van Erven, "Educational data mining: Predictive analysis of academic performance of public school students in the capital of Brazil," J. Bus. Res., vol. 94, no. August 2017, pp. 335–343, 2019.

[19] R. Ahuja, A. Jha, R. Maurya, and R. Srivastava, Analysis of Educational Data Mining, vol. 741. Springer Singapore, 2019.

[20] S. Prabakaran and S. Mitra, "Survey of Analysis of Crime Detection Techniques Using Data Mining and Machine Learning," J. Phys. Conf. Ser., vol. 1000, no. 1, pp. 1–10, 2018.

[21] O. Vaidya, S. Mitra, R. Kumbhar, S. Chavan, and R. Patil, "Comprehensive Comparative Analysis of Methods for Crime," Int. Res. J. Eng. Technol., pp. 715–718, 2018.

[22] P. Vrushali, M. Trupti, G. Pratiksha, and G. Arti, "Crime Rate Prediction using KNN," Int. J. Recent Innov. Trends Comput. Commun., vol. 6, no. 1, pp. 124–127, 2018.

[23] V. Ravi Jain, M. Gupta, and R. Mohan Singh, "Analysis and Prediction of Individual Stock Prices of Financial Sector Companies in NIFTY50," Int. J. Inf. Eng. Electron. Bus., vol. 10, no. 2, pp. 33–41, 2018.

[24] P. Carmona, F. Climent, and A. Momparler, "Predicting failure in the U.S. banking sector: An extreme gradient boosting approach," Int. Rev. Econ. Financ., pp. 1–54, 2018.

[25] A. Kiruthika, P. Deepika, S. Sasikala, and S. Saranya, "Predicting Ailment of Thyroid Using Classification and Recital Indicators," Int. J. Sci. Res. Comput. Sci. Eng. Inf. Technol., vol. 3, no. 3, pp. 1481–1485, 2018.

[26] K. Lakshmi, D. I. Ahmed, and G. Siva Kumar, "A Smart Clinical Decision Support System to Predict diabetes Disease Using Classification Techniques," 2018 Ijsrset, vol. 4, no. 1, 2018.

[27] S. García, J. Luengo, and F. Herrera, "Tutorial on practical tips of the most influential data preprocessing algorithms in data mining," Knowledge-Based Syst., vol. 98, pp. 1–29, 2016.

[28] A. Baldominos, P. Isasi, and U. C. I. I. I. De Madrid, "Feature Set Optimization for Physical Activity Recognition Using Genetic Algorithms," Proc. Companion Publ. 2015 Genet. Evol. Comput. Conf. - GECCO Companion '15, pp. 1311–1318, 2015.

[29] C.-F. J. Kuo, C.-H. Lin, and M.-H. Lee, "Analyze the energy consumption characteristics and affecting factors of Taiwan's convenience stores-using the big data mining approach," Energy Build., vol. 168, pp. 120–136, 2018.

[30] I. M. El-hasnony, H. M. El Bakry, and A. A. Saleh, "Comparative Study among Data Reduction Techniques over Classification Accuracy," Int. J. Comput. Appl., vol. 122, no. 2, pp. 9–15, 2015.

[31] F. Herrera and S. Garc, "Prototype Selection for Nearest Neighbor Classification Taxonomy and Empirical Study," IEEE Trans. Pattern Anal. Mach. Intell., vol. 34, no. 3, pp. 417–435, 2012.

[32] Y. Cheng, K. Chen, H. Sun, Y. Zhang, and F. Tao, "Data and knowledge mining with big data towards smart production," J. Ind. Inf. Integr., vol. 9, pp. 1–13, 2018.

[33] M. Mafarja, I. Aljarah, A. A. Heidari, A. I. Hammouri, H. Faris, and A. M. Al-zoubi, "Evolutionary Population Dynamics and Grasshopper Optimization Approaches for Feature Selection Problems," Knowledge-Based Syst., vol. 154, no. 7, pp. 25–45, 2017.

[34] A. Daraei and H. Hamidi, "An efficient predictive model for myocardial infarction using cost-sensitive J48 model," Iran. J. Public Health, vol. 46, no. 5, pp. 682–692, 2017.

[35] S. Moedjiono, Y. R. Isak, and A. Kusdaryono, "Customer Loyalty Prediction In Multimedia Service Provider Company With K-Means Segmentation And C4 . 5 Algorithm," in 2016 International Conference on Informatics and Computing (ICIC), 2016, pp. 1–6.

[36] J. Tarapitakwong, B. Chartrungruang, and N. Tantranont, "A Classification Model for Predicting Standard Levels of OTOP ' s Wood Handicraft Products by Using the K-Nearest Neighbor," Int. J. Comput. Internet Manag., vol. 25, no. 2, pp. 135–141, 2017.

[37] D. García-gil, J. Luengo, S. García, and F. Herrera, "Enabling Smart Data : Noise filtering in Big Data classification," Inf. Sci. J., vol. 479, pp. 135–152, 2019.

[38] U. O. Cagas, A. J. P. Delima, and T. L. Toledo, "PreFIC : Predictability of Faculty Instructional Performance through Hybrid Prediction Model," Int. J. Innov. Technol. Explor. Eng., vol. 8, no. 7, pp. 22–25, 2019.

[39] D. J. Armaghani, M. Hasanipanah, A. Mahdiyar, M. Z. A. Majid, H. B. Amnieh, and M. M. D. Tahir, "Airblast prediction through a hybrid genetic algorithm-ANN model," Neural Comput. Appl., vol. 29, no. 9, 2018.

[40] V. Rashidian and M. Hassanlourad, "Predicting the Shear Behavior of Cemented and Uncemented Carbonate Sands Using a Genetic Algorithm-Based Artificial Neural Network," Geotech. Geol. Eng., vol. 31, no. 4, pp. 1231–1248, 2013.

[41] L. D. Chambers, Practical handbook of genetic algorithms: complex coding systems. CRC Press, 2010.

[42] L. Parisi and N. RaviChandran, "Genetic Algorithms and Unsupervised Machine Learning for Predicting Robotic Manipulation Failures for Force-Sensitive Tasks," in 2018 4th International Conference on Control, Automation and Robotics (ICCAR), 2018, pp. 22–25.

[43] R. F. Martinez, P. Jimbert, J. Ibarretxe, and M. Iturrondobeitia, "Use of support vector machines, neural networks and genetic algorithms to characterize rubber blends by means of the classification of the carbon black particles used as reinforcing agent," Soft Comput., pp. 1–10, 2018.

[44] D. Zheng, Z. Qian, Y. Liu, and C. Liu, "Prediction and sensitivity analysis of long-term skid resistance of epoxy asphalt mixture based on GA-BP neural network," Constr. Build. Mater., vol. 158, pp. 614–623, 2018.

[45] W. Zhou, D. Liu, and T. Hong, "Application of GA-LM-BP Neural Network in Fault Prediction of Drying Furnace Equipment," in MATEC Web of Conferences, 2018, vol. 232, pp. 1–5.

[46] W. Liu et al., "Discrimination of geographical origin of extra virgin olive oils using terahertz spectroscopy combined with chemometrics," Food Chem., vol. 251, pp. 86–92, 2018.

[47] S. Kumar and G. Sahoo, "A random forest classifier based on genetic algorithm for cardiovascular diseases diagnosis," Int. J. Eng., vol. 30, no. 11, pp. 1723–1729, 2017.

[48] A. Tharwat, H. Mahdi, M. Elhoseny, and A. E. Hassanien, "Recognizing human activity in mobile crowdsensing environment using optimized k-NN algorithm," Expert Syst. Appl., vol. 107, pp. 32–44, 2018.

[49] M. Huang, R. Lin, S. Huang, and T. Xing, "A novel approach for precipitation forecast via improved K-nearest neighbor algorithm," Adv. Eng. Informatics, vol. 33, pp. 89–95, 2017.

# IoT Testing-as-a-Service: A New Dimension of Automation

Babur Hayat Malik[1], Myda Khalid[2], Maliha Maryam[3], M. Nauman Ali[4], Sheraz Yousaf[5], Mudassar Mehmood[6], Hammad Saleem[7]

Department of CS & IT
University of Lahore, Chenab Campus, Gujrat Pakistan

*Abstract*—**Internet of Things (IoT) systems has become a global trend enhancing the capabilities smart computing era involving a variety of distributed end-devices and multi- scalable applications. The collaborative nature of IoT systems connected through the Internet increases the heterogeneity of coming data streams that need to be processed for correct decision making in a real-time environment. The processing of huge data streams for remotely distributed IoT systems create loops for data breaches and open new challenges for security and scalability of system testing. Thus, the testing of IoT systems is becoming the necessity, requires automated testing framework due to the amount of IoT devices and processing of data events is prone to error by traditional software testing. An automated IoT testing service based framework is purposed in this paper, to test the distributed IoT systems by reducing cost and scalability issues of software testing. The infrastructure of IoT systems demands a large number of platforms be developed which requires systematic testing approach. Therefore, the purposed automated IoT testing as a service model performs distributed interoperability testing, oneM2M based conformance testing, security testing of distributed systems and validating semantics/syntactic testing of IoT devices in a systematic approach. Lastly, to provide more strength to the work we discussed and analyze existing IoT testing models to evaluate our proposed model.**

*Keywords*—*Testing automation; IoT; interoperability testing; conformance testing; security testing; semantic testing*

## I. INTRODUCTION

The connectivity of people, objects and devices termed as Internet of Things (IoT). IoT uses sensors devices let physical objects and virtual world to conjoin their environment via the Internet [1]. The connectivity of virtual and physical over the Internet opens up new dimensions such as smart cities, smart homes, etc. the growth in IoT services and devices give rise to humongous data streams and a single device will create four times in its five-year duration and this amount will lead to more than 600 Zeta-byte data and number of IoT devices increase to 29 billion in the year 2020 [2]. IoT environment provides an opportunity for high-scalable end-devices with constrained computing and storage along with cloud integration to maintain latency sensitive systems in IoT [3].

Due to highly complex distributed computing structure, lack of communication frameworks and multiple protocols developing these type of systems is a tedious task because such complexity in the structure of IoT systems is vulnerable to unauthorized access and external attacks [4], [5]. Usage

test-case and test suites for ensuring the security and interoperability of IoT system is challenging but testing IoT systems with the integration of diverse technologies and handling big data streams makes it more difficult [6]. Therefore, there is a need to implement a testing framework to ensure both conformance and interoperability along with the security of IoT systems [4]. Although, a number of researches on exploring the security of IoT has been done in the past. Only rare searches highlight the security, correctness, and completeness of both hardware and software attached remote IoT applications [7]. Mostly, the IoT system (software and hardware) has been overlooked from past years.

In this paper, we have presented a model to automate IoT testing in a real-world scenario. Firstly, we go through the previous work done on IoT interoperability and conformance testing approaches to strategies potential IoT framework for testing. Our proposed work is a model-based approach to testing as a service which is interconnected to the Internet of things (IoT) systems. A distributed cloud service based approach over the network is being adapted to facilitate the IoT system with automated testing as a service. Basically, we are extending a testing service model and using a holistic approach to debug network-related features and perform efficient testing of remote IoT systems. This model analyzes technicalities of testing IoT devices and present solutions by incorporating interoperability, conformance testing along with the security validation of IoT devices and also performing semantic and syntactic testing.

This article is organized as Section II. Background work and present the overview of automated testing. In Section III, various testing methodologies in IoT are discussed. Section IV describes our purposed model of automated IoT testing as-a-service then Section V discusses the strengths of this model. Lastly, Section VI concludes the article.

## II. BACKGROUND

The model-based testing is studied thoroughly in literature [8], [9]. Whereas, mostly IoT connected approaches are premeditated for mobile applications. Though, other related work emphasizes more on the liability of IoT based systems testing. Authors in [10] developed an approach named as IoT testing as a service for creating automatic test cases with the use of several patterns. Cloud consumers and cloud providers are also provided with a testing service known as Testing as a Service (TaaS) [11]. Work done on TaaS is more related to web services and cloud computing. Zech et al. in [12]

presented an approach for creating test cases with help of risk analysis of cloud computing for ensuring security.

The work in [13], described standardized interoperability testing implementation and other affecting issues like cost, scalability, and coordination. These are few issues which arise during the IoT products development and methodologies being applied such as testing methods used by the telecommunication industry aren't flexible enough to command the IoT systems. Few of them aren't able to interact with high-level protocols because of their small sizes. Interoperability issue occurs at the semantic level because the streams of data that passes on are needed to be checked at semantic and syntactic levels for the purpose of data correction if any of data is a flaw. In the past, IoT testing was used for handling this issue of interoperability from the creation and execution of test cases and testbeds for real-world IoT devices.

System security concepts are described in this paper [14], along with language-based technique which describes the process of dealing with BOF. It also describes the data structures and techniques which are used for code and memory analysis (Control Flow Graph (CFG) and Dependence Graph (DG)), (Points-To and TFA), and Memory Safety tools (Address Sanitizer and SAFE Code). Buffer overflow contains a large quantity of data to let go of the upper bound of the buffer which means overwriting of data on another. It usually occurs in the form of heap or stack. Dependence Graph (DG) is used for modeling of data and instructions reliance in the program. Tainted flow analysis is another technique for poking paths of information that moves through inputs to sensitive operations.

In [15], the RM model is being used for knowing the IoT domain along with other models for finding out IoT concepts and constraints. It also works as a base of RA. RM particularly consists of a first level IoT concepts description known as Domain Model. It also consists of an Information Model that deeply describes the processing of IoT information. RA offers key Functional Groups (FG) which is required by IoT architecture through its functional view. FG explains applications functionalities which are made on the peak of IoT infrastructure. It also offers IoT –aware demonstration which is accomplished during process execution.

### A. Overview of Internet of Things (IoT) Automation Testing

Test automation in IoT is used for execution management and compares actual result with the predicted one. It also helps out in enhancing the speed of unit testing, API testing, and GUI testing. Though it also executes regression tests and its extremely economical and with the technology shift industry depends for testers with automated skills [16]. Apart from the regression test, compatibility tests are also run by it which improves productivity level and make sure that customer is provided with quality software. Although it enhances the efficiency of tests few drawbacks also exist over here. These drawbacks involve unable to enhance test potency and identifying errors. Therefore, one major drawback involves scripts automation. Mohd Ehmer Khan discussed such software testing tools which are used for testing software like performance, reliability, etc. As these, all are categorized

according to their main working and in different types [17]. Researcher Manjit Kaur did a comparison between different automation tools just like QTP, TC [18]. The tools QTP Pro and QA are compared on the basis of various characteristics of cost, time and scripts creation, etc. QTP is basically more efficient in regard to those applications which requires more security whereas test complete efficiently work for those applications which require less security. Author Harpreet Kaur presented a comparative analysis of different tools like Selenium, etc. and identified their performance on the basis of cost, application support, etc. QTP [19] is compared with selenium and TC and deeply analyzed and compared according to each possible factor and considered best among all [20].

## III. TESTING METHODOLOGIES IN INTERNET OF THINGS

### A. Interoperability Testing in Internet of Things (IoT)

For the assurance of network interoperability testing standard bodies of communication (ETSI and Bluetooth SIG) benchmarks some rules and processes which includes plug-test events and conformance testing. Plug-test events involve a meeting of organizations who implement technologies, each party test and check their systems against others. Such events cost high overdue of organizations and also attended by IoT communities and research centers, it also requires one developer and tester which is economically not preferable without sponsors for open-source communities [10], [21]. To test the network interoperability in IoT systems, an external IoT system (both software and hardware) is integrated by the third party service providers with the minimum code already written to apply initial functions to set the system in a stable state for interoperability testing that can be handled by the abstraction layer like resetting device or configuring network. Test cases to test the interoperability of the system are presented at the top of the abstraction layer. The challenged faced in executing this method requires both third party and already integrated communication systems to be present in the similar location which can be controlled by implementing transparent network bridges in distributed test system scenario [10]. The communication systems involving System Under Test (SUTs) need to be connected with the end point of the bridges in order to transfer communication to and from the third party systems placed in a different locality. In addition to this, to create a network bridge both the wireless transceivers and endpoint systems must be using the same distributed messaging service (e.g. IEEE 802.15.4).

*1) Interoperability testing models of IoT*: Various types of configuration testing are modeled to deal with the diverseness of IoT test and deployment systems controlled by testbed alliance are discussed as [13]:

- *Simple Conformance Testing:* Appropriate for testing the conformance of only one IUT at a time. It can only check the functionality of the IoT devices.

- *Simple Conformance and Interoperability Testing:* This model is suitable for both conformance and interoperability testing of only single new IUT with a standard testbed.

- *Simple Conformance and Compound Interoperability Testing:* This mode is appropriate for conformance and interoperability testing of a new IUT with a number of testbeds in the system when reference implementation is unclear.

- *Compound Conformance Testing:* This model performs conformance and interoperability testing to analyze the cooperative and collaborative behavior of several IUTs without testbeds.

- *Compound Conformance and Compound Interoperability Testing:* This model performs conformance and interoperability tests on several IUTs numbers of testbeds.

### B. Conformance Testing in Internet of Things (IoT)

Conformance testing measures and ensures the implementation of a specific standard to the required level. Generally, the conformance testing model comprises of two parts, one carries out the Implementation Under Test (IUT), is System Under Test (SUT) and other is Means of Testing (MOT) which includes coordination, logging and reporting activities handled by a minimum of one tester depends on IUT's architecture and interface [13].

*1) Architecture of conformance testing:* There are various creation elements are used for conformance testing such as Implementation Conformance Statement (ICS), Implementation eXtra Information for Testing (IXIT), Testing Description Language (TDL), Executable Test Suite (ETS), Abstract Test Suite (ATS), and Equipment under Test (EUT) from [22]. Product Functionalities and abilities used for the purpose of checking and provision of interoperability signals are used by ICS. Extra important metadata is provided by IXIT for testing purpose. Test cases are described by formal language known as TDL. Test suites are described by another formal language known as Tree and Tabular Combined Notation version 3 (TTCN-3) and it is done by ETSI and few others like SDOs i.e. 3GPP and oneM2M. Group of test cases which shows test completion and described in a normal language like TTCN-3 is all done by ATS whereas, ETS uses TTCN compiler and its totally irreplaceable. Conformance testing isn't only used for testing of normal behaviors, rather also used for testing of extraordinary behaviors. It also enables the tester to perform broader functional testing. It doesn't completely ensure the interoperability of the system with other systems because the standard test might leave some space for configuration and conformance purposes.

### C. Security Testing in Internet of Things (IoT)

The security requirements of IoT system are of extreme importance as functional requirements due to the vulnerable of security functions. However, analytics of IoT systems and security tests summarize factors of IoT systems which resist security issues such as usage of other systems, system security threats and vulnerability, and security function's exploitation. Some Model-Based Testing (MBT) Standards like M2M (Machine to Machine) and its extension oneM2M identified some security vulnerabilities and requirements level that needs

to be satisfied before developing the system [23]. Three testing strategies need to be implemented in groups or discretely to check the validation of these requirements. These are as follows:

- *Security Functional testing (compliance with agreed standards/ specification):* It analyzes the system against the required functional specifications in order to ensure that the implementation of security functions is implemented in an approved manner.

- *Vulnerability testing (pattern driven):* It analyzes security attacks and risk-based vulnerabilities. It is based on security patterns used to initialize the security testing, then the targeted test patterns are used to apply appropriate test cases for possible security threats identified.

- *Security robustness testing (behavioral fuzzing):* Measuring invalid messages created by test cases in order to deal unpredicted behavior of the system for the security threat and attacks on large scale IoT systems.

Also, Model-Based Testing (MBT) approaches have used with their shown their benefits and usefulness for security testing of large-scale IoT systems undergoing particular standards defining guidelines and solutions for these security elements of the system [23].

### D. Semantic Testing in Internet of Things (IoT)

Testing IoT is a level based approach in which conformance testing and interoperability testing performs protocol level testing, security testing focus on vulnerabilities in a system. While the basic purpose of semantic testing is to test the semantic accuracy of IoT data streams in accordance with the pre-defined standards [13].

The implementation of semantic testing in IoT paradigm is most challenging because of the heterogeneity of IoT devices and semantics testing performs validation in the semantic description at various targeted levels like testing lexical and syntactic validation and then logical and semantic validation. Some reference ontologies have already been defined including oneM2M ontology, W3C-SSN ontology, etc. After defining these ontologies, the next step is the conformance test against these reference ontologies to achieve semantic interoperability. Such diversity in concepts and relations of semantics models could make application of semantic interoperability more complex [24]. Some ongoing research projects like H2020 Fiesta-IoT, provide a unique cloud platform for conducting a test on semantic technologies using semantic IoT testbeds. These cloud platforms give access to semantic data of various testbeds such as smart cities, smart homes, etc., through uniquely identified access points. For analyzing the correctness of semantic and syntactic validation, data regarding particular ontology is selected from the semantic database and used for experiments against defined standard semantic. The database will reject the data in order to keep it clean and accurate if the data does fulfill all the requirements and semantics description reporting all the errors will be provided to ontology developers to model these errors while improvement phase. To complete semantic interoperability of IoT systems, achieving semantic testing is

required. The inclusion of lexical, syntactical, semantic correctness and test feature are crucial for attaining semantic testing [25].

## IV. PROPOSED TESTING-AS-A-SERVICE MODEL

We discussed various IoT testing methodologies in the previous section. Both interoperability and conformance are traditional testings in IoT. Also, security testing and semantic testing models are used in IoT which is a major part of IoT testing. Therefore, in this section, we integrate these testing methods used and formulate a model shown in Fig. 1, for IoT testing-as-a-service.



Fig 1. Proposed IoT testing-as-a-service Model.

### A. Distributed Interoperability Testing for Remote Devices

The application of Interoperability testing on our model is based on a remotely distributed test system architect for automating both interoperability and conformance testing. Previously discussed network interoperability testing applies a suite of test cases but here we presented the extension of distributed interoperability testing with using distributed test plugs which will help developers and third party service providers with quick test case response from different locations. ETSI designed a Constrained Application Protocol known as CoAP for such plug-tests [26]. CoAP specifies a group of test requirements for interoperability testing and each requirement defines CoAP properties, after finalizing these requirements a test case is derived for each of them. From details of these test cases expected system behavior of CoAP protocol is analyzed [27]. Therefore passive testing methodology is appropriate for such resource constrained and operational architecture of IoT, which does not allow overheads in networks. Furthermore, to test the implementation of passive testing a message (Pass, Inconclusive or Fail) is released if a packet is captured by packet sniffer between client and server shown in Fig. 2 as:



Fig 2. Architecture of CoAP Interoperability Testing.

Distributed CoAP test plugs involve two different configurations. The basic system includes TS (Test System) and SUT involving two IUTs as CoAP server and client shown in Fig. 3. However, using passive testing technique might cause capturing packets by sniffers while exchanging packets between IUTs shown in Fig. 4. Thus, distributed CoAP interoperability testing in such environment uses a UDP gateway in between CoAP server and CoAP client to replicate a lossy medium.



Fig 3. Basic CoAP Testing Model.



Fig 4. CoAP Passive Testing Model.

Based on the testing of some test requirements for CoAP, working model in both consistent and packet lost scenarios. The distributed CoAP interoperability testing verifies the correctness of client and server interaction involving HTTP method (GET, POST, PUT and DELETE) by analyzing each request/response by the client have correct message code [28]. On the other hand, the server sends a piggy-backed reply upon receiving a request from client such as if the request is confirmed send ACK (acknowledgment), if there is delay in getting a request the server first sends an empty ACK message then upon receiving the request it sends a confirmed response and it will send not-confirmed response for non-confirmable requests. There are some major options selected on the basis of basic transactions such as Token option analyzing any delay in request and response timing, each client request is assigned a token to synchronize the response, URI schemes are used to identify and locate CoAP resources. Here we are using URI query option in which requested resource is allotted and a correct response with accurate message code/type is sent by the server against the client's request [29].

### B. Conformance Testing based on oneM2M

M2M testing framework is developed for lab-based conformance testing. In this type of testing developers and vendors have to go to labs for conformance testing purpose. Due to large SMEs (Small to medium-sized enterprises), an appropriate testing method is required whereas individual developers working generates less number of IoT devices and this M2M isn't sufficient for this purpose. Identifying low-cost IoT testing processes isn't any big deal as IoT testing has to manage many communicative standards and protocols [13], [30] creation and coverage of different protocols, logging, etc. Isn't an easy task for SMEs and developers. The automated testing attribute is used in the web-based remote testing framework. It also resolves all the occurring problems during testing whereas the main purpose of this testing is to allocate main conformance logic and provision of APIs. This provision of APIs helps testers in configuration selection according to their needs.

New IoT protocols can be chosen by web-based testing instead of enabling unknown third parties to include their protocols to the core system. For initiative test case communication triggering of the device is necessary by M2M. For helping various network protocols M2M ensures flexibility with the usage of network protocol through the UpperTester performs the previous action. UpperTester is a software which is used for converting test pointer to a message that is perceivable by IUT. IUT's ability decides the implementation of UpperTester either inwards or outwards of IUT.

Testing configuration information is required by a tester who is about to test IoT device such that selection of test cases, protocols, and devices that performs web interface. An actuated message consisting of test cases and configuration data is sent by test system to UpperTester on basis of inputs entered. On the basis of the provided guidelines, one M2M action is performed by the test system as tester passes on the message to the test system. M2M function consists of creating, retrieve, update, delete and notification. One of these M2M function is guided to the IoT device by UpperTester [13].

An agreement is required between UpperTester and IoT device for doing test procedures mentioned above. The specific operation is applied by IoT device on the basis of mentioned test cases in actuated messages when UpperTester provides test case data. After verification of conforming standard messages, test system develops findings of IoT device's conformance testing [31].

The second step of the testing model is the provision of support for managing communicative variables and also automated assistance for developing conformance tests verdicts. IoT uses different kinds of protocols as their integration is quite necessary for a testing framework. An automated IoT testing feature is being developed by us for usability and test distribution. That feature in the framework is described as follows:

*1) Protocol adapter:* Ascendable testing is done by various protocols of various domains and it is done by IoT devices. IoT devices need scalable testing using various protocols for different domains of application. Normal data integrity is done by TCP (Transmission Control Protocol) and request-response time is also dependent on TCP. Whereas, publish-subscribe is only used when real-time communication is required by the environment. Scalability by external code is based on IoT TaaS.

*2) Automated device testing:* Different applied transactions in the target device are actuated for testing IoT devices. Required action performed by developers is done by using stimulus. IoT TaaS defined various transaction and message types. It has helped IoT developers to perform automated testing by simply inserting a code.

### C. Securing IoT using Distributed Systems Analysis

Securing Internet of Things (SIoT) model comprises of two parts: application independent and dependent part and executed on Top level of LLVM compiler [32]. The SIoT core is the application independent part for static analysis. The application dependent part is the SIoT instance involves libraries generated by the users for the implementation of the static analysis. The DDG graph is always created for a program which acts as a bridge between SIoT core and instance and works for each instance particularly as shown in Fig. 5.

*1) The architecture of the SIoT core:* SIoT uses LLVM IR a low-level language to process code files. Formed\ by bytecodes [14] (3-address instructions) use the various size of integers: bit vectors, floating point numbers, arrays, and labels. Using a group of files in this format it creates a DDG graph by undergoing a two-level process of merging and linking. In merging, multiple files are mere into one file reducing the naming conflict of files i.e. several files as the same name. this tool evaluates the network function each bytecode file which is required by merging phase to name bytecode files with the Send and Recv functions [33]. SIoT can identify these functions to add different tags to bytecode files and then merge them into a single file to ease the analysis. Whereas linking uses the recv function for the

creation of SEND and RECEIVE graph of all programs to generate DCFG. Using DCFG, Dist.-Dep graph generates DDG which can detect the vulnerabilities in the data flow.

*2) Buffer overflow instance*: For detecting the liabilities of BOF attacks we analyze memory and input dependencies. If data is functioned to unreliable input, then we highlight the vulnerability of the array. An unreliable input may be accessed by some malicious user with the sensors, serial ports or by the file system to use the memory access. Through DDG graph we can analyze the flow of information and detect the vulnerabilities in the program by providing DDG and input values to DistVulArrays. The LLVM pass checks various paths among memory access and unreliable inputs in the DDG. Once the analysis of program completed DistVulArrays gives these possible outputs such as true-positive rate and potential false-positive rate and number of malicious paths in the program in graph [14].



Fig 5. Architecture of SIoT.

### D. Validating Semantics Testing

F-Interop project is being created for the implementation of semantic interoperability in conformance testing. F-Interop helps test systems and SUT which are placed in far areas by providing a cloud-based platform. This platform has enabled developers to work from their residences instead of moving from place to place, in this way more applicable tests are being generated in a better way in regard to time and cost. High-level testing premises are discussed in upcoming paragraphs which are applied within EU H2020 F-Interop project 11.

There are various scripts of semantic conformance test described as follows: There is fundamental interaction between the tester and SUT (System Under Test). SUT sends semantic data which is then checked by the tester whereas at end of conformance testing tester provides a report regarding

completion of semantic data according to ontology acknowledgment. And if any issue occurs in that semantic data that issue is mentioned in the report. Semantic conformance testing chart is discussed in Fig. 6.



Fig 6. Basic Workflow of Semantic Testing.

There are two test scenarios of semantic interoperability. One carries SUT whereas a halfway tester is required for completion of the test in the second scenario. Here is technology agnostic that tells they aren't obliged of testing any specific platform or semantic attribute such as identifying if the semantic descriptor is generated in a proper way that it consists of M2M system's semantic data. The main purpose behind all this generation of such tests is applicable for each type of semantic data that obeys some specifications [34]. This enables test integration by applying specific standard along with all types of test which is limited to standard.

Semantic interoperability is considered as data interpretation from the system. In this premise, every portion generates a piece of semantic data processing which are as follows: semantic data and semantic query. Results of semantic processing parts are then compared if they are equal or not. Their equality shows their mutual understanding. SUT1 and SUT2 should have similar semantic queries for executing tests. At last, if a similar query is executed from similar data of both SUT1 and SUT2 it results that both SUT1 and SUT2 have the same level of understanding of data.

Interoperability performed at the data level. As discussed before interoperability is given on the basis of ontology. Hence, our purpose is to identify semantic data which is used on the basis of ontology used. Ontology is a combination of vocabulary and the relationship between vocabularies. In this test, two SUT's data submitted is checked if they have the same vocabulary which is discussed in the same ontology. If they share similar vocabulary, then it is implemented at a semantic data level because they are workable.

Transmitted semantic data (D1 and D2) produced from SUT1 and SUT2 have verified their conformance as it's a condition of test. Tester recovers D1 and D2 vocabulary and verifies if they share a similar vocabulary. If similar vocabulary is shared, then D1 and D2 are totally practical [35].

## V. DISCUSSION

The integration of cloud and Internet of Things platforms provides various services such as platform as a service, infrastructure as a service and software as a service. The distributed nature of IoT devices also requires such kind of testing service for analyzing and reconfiguring IoT application during development. In this papers, we purposed ad extension of plug-test with existing IoT testing methodologies, the framework of automation testing as a service has four phases: interoperability testing is performed using CoAP protocol to verifies the correctness of client and server interaction and analyze the request/response of target message type [28]. Conformance testing based on oneM2M use test plugs to test system specifications using test case on IUTs. Validating semantic testing used different ontologies to validate the semantic/syntactic correctness of the particular document. Furthermore, the addition of security testing in the model identifies the vulnerabilities in the system and provides a solution to increase system reliability. Therefore, this framework could allow developers to easily implement automation testing as a service to enhance correctness, reliability, and interoperability of IoT application being developed.

## VI. CONCLUSION

The testing model presented in this paper is a service-based approach of IoT system testing which enables automated testing for distributed IoT systems by providing constraints on cost, scalability, and complexity of IoT applications. Firstly, we analyze automation testing in IoT and then, in accordance with this we presented an insight into existing methodologies of IoT testing with its design and implementations. Furthermore, we extended an existing testing concept and introduced a novelty framework to generalize testing-service in remote IoT systems. Automation IoT testing as a service model architecture incorporating four IoT testing methodologies: distributed interoperability testing, conformance testing based on oneM2M, security testing distributed systems and semantic/syntactic testing in a systematic approach. Our model creates a distributed plug-test to enables network interoperability testing without delaying data transfer from one SUT to another irrespective of location constraints. As future work, we will extend this work in order to design automation testing suites to enables the development team to analyze and enhance the security of IoT devices.

### REFERENCES

[1] L. Ling, M. Loper, Y. Ozkaya, A. Yasar, and E. Yigitoglu, "Machine to machine trust in the IoT era," in 18th International Workshop on Trust in Agent Societies, Singapore, 2016.

[2] "Cisco global cloud index: Forecast and methodology, 20152020," Tech. Rep., 2016.

[3] E. Yigitoglu, M. Mohamed, L. Liu, and H. Ludwig, "Foggy: A Framework for Continuous Automated IoT Application Deployment in Fog Computing," 6th International Conference on AI & Mobile Services, IEEE, 2017.

[4] I. Schieferdecker, S. Kretzschmann, A. Rennoch and M. Wagner, "IoT-Testware- an Eclipse Project," International Conference on Software Quality, Reliability and Security, IEEE, 2017.

[5] M. Ammar, G. Russello and B. Crispo, "Internet of Things: A survey on the security of IoT frameworks," Journal of Information Security and Applications, vol. 38, pp. 8-27, 2018.

[6] M. Leotta, F. Ricca, D. Clerissi, D. Ancona, G. Delzanno, M. Ribaudo and L. Franceschini, "Towards an Acceptance Testing Approach for Internet of Things Systems," ICWE, Springer, 2018.

[7] Z. B. Celik, P. McDaniel and G. Tan, "Soteria: Automated IoT Safety and Security Analysis," USENIX Annual Technical Conference, 2018.

[8] M. Utting, B. Legeard, F. Bouquet, E. Fourneret, F. Peureux and A. Vernotte, "Chapter 2 - recent advances in model-based testing," Advances in Computers 101, pp. 53–120, 2016.

[9] M. Utting, A. Pretschner and B. Legeard, "A taxonomy of model-based testing approaches," STVR vol. 22, issue 5, pp. 297–312, 2012.

[10] P. Rosenkranz, M. Wahlisch, E. Baccelli and L. Ortmann, "A Distributed Test System Architecture for Open-source IoT Software," Proceedings of Workshop on IoT challenges in Mobile and Industrial Systems, pp.43-48, 2015.

[11] L. M. Riungu, O. Taipale and K. Smolander "Research issues for software testing in the cloud," 2nd International Conference CloudCom, IEEE, pp. 557–564, 2010.

[12] P. Zech, M. Felderer and R. Breu, "Towards a model-based security testing approach of cloud computing environments," 6th International Conference SERE-C, 2012.

[13] H. Kim, A. Ahmad, J. Hwang, H. Baqa, F. Gall, M. A. R. Ortega and J. Song, " IoT-TaaS: Towards a Prospective IoT Testing Framework," in IEEE Access, vol. 6, 2018.

[14] F. A. Teixeira, F. M. Q. Pereira, H. Wong, J. M. S. Nogueira and L. B. Oliveira, "SIoT: Securing Internet of Things through Distributed System Analysis," Future Generation Computer Systems, vol. 92, 2019.

[15] S. De, F. Carrez, E. Reetz, R. Tonjes and W. Wong, "Test-Enabled Architecture for IoT Service Creation and Provisioning," Springer, 2013.

[16] K. Saravanan and E. P. C. Prasad, "Open Source Software Test Automation Tools: A Competitive Necessity," International Journal of Management and Development, vol. 3, issue 6, pp. 103-110, 2016.

[17] M. E. Khan, "Different forms of Software testing techniques for finding errors," International Journal of Computer Science Issues, vol. 7, issue 3, pp. 11-16, 2010.

[18] M. Kaur, and R. Kumari, "Comparative study of automated testing tools: Test complete and quick test pro," International Journal of Computer Applications, vol. 24, issue 1, pp. 1-7, 2011.

[19] Mercury Quick Test Professional tutorial, version 8.0. Mercury Interactive Corporation, Documentation, 2004.

[20] Automated Testing: Process, Automated Testing Tutorials: What is Process, Benefits and Tools Selection. Guru99. Retrieved (2015), from http://www.guru99.com/automation-testing.html.

[21] L. Atzori, A. Iera, and G. Morabito, "The Internet of Things: A survey," Computer Networks, vol. 54, 2010.

[22] S. Moseley, S. Randall, and A. Wiles, ``Experience within ETSI of the combined roles of conformance testing and interoperability testing,'' in Proc. 3rd Conf. Standardization Innovation Information Technology, pp. 177-189, 2003.

[23] A. Abbas, G. Baldini, P Cousin, S. N. Matheu, A. Skarmeta, E. Fourneret and B. Legeard, "Large Scale IoT Security Testing, Benchmarking and Certification," Cognitive Hyperconnected Digital Transformation, Chapter: 7, 2017, pp.189-220.

[24] M. Bermudez-Edo, T. Elsaleh, I. P. Barnaghi, and K. Taylor, "A lightweight Semantic model for the Internet of Things," in Proc. Int. IEEE Conf. Ubiquitous Intell. Comput., Adv., Trusted Comput., Scalable Comput. Commun., Cloud Big Data Comput., Internet People, Smart World Congr, 2016.

[25] Linked Data_Connect Distributed Data Across theWeb. (Online), 2017 Available: http://linkeddata.org/.

[26] C. Lerche, K. Hartke, M. and Kovatsch, "Industry adoption of the internet of things: A constrained application protocol survey," In Proceedings of the 7th International Workshop on Service Oriented Architectures in Converging Networked Environments, 2012.

[27] D. Lee, A. N. Netravali, K. K. Sabnani, B. Sugla, and A. John, "Passive testing and applications to network management," International Conference on Network Protocols, IEEE, pp. 113-122, 1997.

[28] A. Ahmad, F. Bouquet, E. Fourneret, F. L. Gall and B. Legeard, "Model-Based Testing as a Service for IoT Platforms," ISoLA, 2016.

[29] N. Chen, C. Viho, A. Baire, X. Huang and J. Zha, "Ensuring Interoperability for Internet of Things: Experience with CoAP Protocol Testing," Journal for Control, Measurement, Electronics, Computing and Communications, vol. 54, issue 4, 2017.

[30] B. Ahlgren, M. Hidell, and E. C. H. Ngai, "Internet of Things for smart cities: Interoperability and open data," IEEE Internet Computer, vol. 20, issue 6, pp. 52-56, 2016.

[31] Functional Architecture, document oneM2M TS-0001-V2.10.0, 2016.

[32] C. Lattner and V. S. Adve, "LLVM: A compilation framework for lifelong program analysis & transformation," In: CGO, IEEE, 2004.

[33] C. Cowan, F. Wagle, C. Pu, S. Beattie and J. Walpole, "Buffer overflows: attacks and defenses for the vulnerability of the decade," DISCEX, DARPA, 2000.

[34] "Web of Tthings (WoT) thing description specification." https://github.com/_w3c/wot-thing-description.

[35] S. K. Datta, C. Bonnet, H. Baqa, M. Zhao and F. L Gall, "Approach for Semantic Interoperability Testing in Internet of Things," GIoTS, 2018.

# Assessment of Technology Transfer from Grid power to Photovoltaic: An Experimental Case Study for Pakistan

Umer Farooq[1], Habib Ullah Manzoor[2], Aamir
Mehmood[3], Awais Iqbal[4], Rida Younis[5], Amina Iqbal[6]
Faculty of Electrical & Mechanical Engineering, University
of Engineering and Technology
Lahore (FSD-Campus)-38000
Pakistan

Fan Yang[7], Muhammad Arshad Shehzad Hassan[8],
Nouman Faiz[9]
State Key Laboratory of Power Transmission Equipment,
and System Security and New Technology
School of Electrical Engineering, Chongqing University
Chongqing 400044, China

*Abstract*—**Pakistan is located on the world map where enough solar irradiance value strikes the ground that can be harnessed to vanish the existing blackout problems of the country. Government is focusing towards renewable integration, especially solar photovoltaic (PV) technology. This work is focused to assess the techno-economic viability of different PV technologies with aim of recommending the most optimum type for domestic sector in high solar irradiance region of the country. For this purpose, standalone PV systems are installed using monocrystalline (m-Si), polycrystalline (p-Si), and amorphous crystalline (a-Si) modules on the rooftop at 31.4 °N latitude position. The performance of PV modules is evaluated based on, average output power, normalized power output efficiency, module conversion efficiency, and performance ratio. Results elaborated that m-Si module is the optimum type for the application with 23.01% average normalized power output efficiency. Economics of the system has also been evaluated in terms of the price of power value produced by PV modules with respect to the consumption of that power value from grid source in base case. Integration of such type of domestic PV systems are a need of time to make the future sustainable.**

*Keywords—Solar energy; photovoltaic technologies; module efficiency; power demand satisfaction; economics*

## I. INTRODUCTION

Energy is the main driving force behind the dynamics of the world. The whole world is running after energy sources to make the future brighter and sustainable. Since the evolution of modern civilization on earth, electricity production is largely dependent on fossil fuels. With scientific technology advancement in 21st century, innovative techniques and inexhaustible energy sources are making their way to meet new standards. In this context, renewable energy sources (RESs) are the focus as they are everlasting and are not associated with threat of being extinct, as fossil fuels are. Significant amount of energy security, climate change mitigation, and certain economic benefits can be achieved by the speedy deployment of RESs [1-3].

Among renewables, solar (photovoltaic) PV technology is one of the most auspicious and emerging one in the whole world. Pakistan is considered at naturally blessed location on

world map for solar applications due to its geographical location, climatic situation, and high solar insulation value [4]. Even Pakistan is facing severe energy crisis presently despite of the reported potential of 455.3 GWh electric power generation through harnessing solar irradiance energy using off-grid type PV systems only [5].

Solar energy harnessing through PV system is dependent on the type and performance of PV module. The performance of PV modules varies depending on the geography and climatic conditions of the location [6]. In PV technology, electricity can be generated directly from solar energy but with low conversion efficiency, and needs much improvement in conversion technology [7]. It requires a lot of efforts to achieve maximum energy from the PV panels. The basic need to test PV modules at outdoor real time conditions is to estimate the output power of PV modules under variant environmental conditions [8-9]. Parameters of PV modules estimated on the standard testing conditions (irradiance 1000 $W/m^2$, module temperature 25°C, and AM 1.5) are not equal to the real time operating conditions due to the variations in environmental parameters [10].

Amin et al., [11] conducted an experimental study for the evaluation of the performance of four PV modules (m-Si (monocrystalline), p-Si (polycrystalline), a-Si (amorphous crystalline), and copper indium diselenide (CIS)) in Malaysia for three consecutive days. Results showed that CIS module has higher performance ratio (PR) value while m-Si module has high module efficiency among all the modules tested. Midtgard et al., [12] conducted the experimentation for the performance evaluation of three different types of solar panels (m-Si, p-Si, and a-Si) in the climate of Norway. It was found that m-Si was better in terms of average output power and module efficiency compared to p-Si and a-Si modules.

The efficiency of PV modules depends on various environment associated factors such as dust accumulation, wind speed, wind direction, temperature, humidity etc., [13-14]. Catelani et al., investigated the effect of dust and pollution on PV module performance using statistical approach [15]. From the statistical hypothesis test, the maximum power values showed decrement with the dust accumulation on panel

surface. Mekhilef et al., investigated the effect of dust, humidity, and air velocity simultaneously [16]. They concluded that effect of each parameter should not be studied separately during estimating the cell efficiency by ignoring the other factors. P. Ferrada et al., investigated the performance of PV modules in the coastal zone of Antofagasta, Northern Chile for about 16 months [17]. They concluded that PR value of PV modules decreases due to dust accumulation. They also analyzed that difference of energy yield between the different technologies became higher for summer and lower for winter season. Generally, performance of PV modules decreases when temperatures of surroundings increases above the nominal operating temperature range [18]. Due to increase in air temperature, PV modules undergo severe degradation, consequently they cannot achieve maximum performance, therefore immediate cooling must be provided [19-20].

The best approach is to manage the light collecting capability and optimizing the temperature of solar PV cell. Jiang et al., performed experiments on three solar cells in a laboratory with sun simulator and test chamber and concluded with decrease in output efficiency value of up to 26% when dust deposition density increased to 22 $g/m^2$ [21]. Bashir et al., reported the performance data for different PV modules (m-Si, p-Si, and a-Si) for the month of January in the climate of Taxila and found that m-Si module was more efficient with demonstrated average module efficiency of 13.01% [22]. It was also found that a-Si module showed average PR 1.03 and performed better than other two modules at that site. Furthermore, module efficiency and PR showed decreasing trend with the increase in solar irradiance and module temperature after nominal operating temperature point.

The present experimental study was focused to assess the impact of technology transfer from grid source to renewables source i.e. PV. For this purpose, three types of PV technologies were tested. It was aimed to identify the optimum PV technology to recommend for domestic applications in middle planes of the country.

Three standalone systems using commercially available PV technologies (m-Si, p-Si, and a-Si) were installed at Faisalabad city (Latitude 31.4 °N, longitude 73.1 °E) of Pakistan, a city with 3rd largest population count. The experimental data was recorded in form of incident solar irradiance, output voltage, current, and module temperature values. Using these recorded values, techno-economic performance of different PV technologies was assessed based on output power, normalized power output efficiency (NPOE), module conversion efficiency, performance ratio, and economic impact in terms of domestic power demand satisfaction and consequent reduction in electricity bill.

The reminder of this paper is organized as follows: the methodology is presented in Section II. Section III briefly explains the results and discussions. The economic impact of power source replacement and conclusions are explained in Section IV and Section V.

## II. METHODOLOGY

The standalone PV systems were installed using three commercially available different PV technologies to satisfy the curiosity that what will be the impact of transferring the power source from grid to PV renewable and which one PV technology can satisfy the domestic power demand optimally in the middle planes of Pakistan.

### A. PV System

The standalone PV system is a type of off-grid system that generates direct current (DC) electric power. This harnessed solar power can be used directly through DC appliances, or indirectly in conventional AC (alternating current) systems after conversion.

The main component of a PV system configuration is its PV module. Considered types from available PV modules for current study are: m-Si, p-Si, and a-Si. The m-Si modules are composed of silicon bars that are converted into wafers by cutting form single crystal of silicon for each cell. The m-Si panels are usually blackish in color. In p-Si solar panels, more than one crystal of silicon is heated to make wafers. Single cell of p-Si module consists of mixture of silicon crystals. These panels are generally blue in color. The a-Si modules fall into the category of non-crystalline thin film silicon solar cells. They typically hold less efficiency compared to m-Si and p-Si solar panels [23].

The specifications and characteristic parameters of all three types of PV modules, along with rated and measured values are tabulated in Table I. Rated values are those provided by manufacturer, while measured values are recorded at outdoor environment in real time conditions.

### B. Experimental Procedure and Setup

Three standalone PV systems were installed at the rooftop of Electrical Engineering Department, University of Engineering and Technology (Faisalabad Campus). Each PV system was configured using a certain type of PV module from selected three technologies. The PV modules were tilted at latitude angle position of the location with the aim of working year around, as shown in Fig. 1.

The measurements were taken for thirty days in the months of Feb-March 2018 during sunshine hours starting from 7:00 a.m. to 6:00 p.m. on hourly basis. The PV modules were connected to digital multimeters to measure corresponding module output voltage and current values. For measuring global solar irradiance values incident on PV modules, light intensity meter with, 2000 $W/m^2$ measuring capacity, 0.1 $W/m^2$ resolution, 400~1000 nm spectral response, and ±10 $W/m^2$) accuracy value, was used. The solar irradiance values were measured consecutively after every hour. The temperature of all the modules was recorded using infrared (IR) laser meter (with specifications: ±2 $^{o}$C accuracy, -32~380 $^{o}$C measuring range, 0.1 $^{o}$C resolution). Digital multimeters were attached with each PV module, which was connected to electronic load to vary the output of the PV modules from zero to maximum. Complete schematic configuration of installed PV system with all measuring devices is portrayed in Fig. 2.

The characteristic parameters of PV modules based upon which PV module can be characterized and differentiated from each other are: maximum power, fill factor, normalized output efficiency, module efficiency, and performance ratio. These

properties are calculated using base measured values (like V, I, E), and applying the following mathematical relations:

Maximum power: $P_{max} = V_{max} \times I_{max}$ (1)

Fill Factor: $FF = \dfrac{(V_{max} \times I_{max})}{(V_{oc} \times I_{sc})}$ (2)

NPOE: $\eta_p = \left(\dfrac{P_{meas}}{P_{max(STC)}}\right) \times 100$ (3)

Module efficiency: $\eta_{mod} = \left(\dfrac{P_{meas}}{E \times A_a}\right) \times 100$ (4)

Performance ratio: $PR = \dfrac{\left(P_{meas}/P_{max(STC)}\right)}{(E \times 1000)}$ (5)

Comparison cannot be done considering module output power values that are different for all the modules. So, to compare the performance of each module with other, measured output power is normalized to its output value calculated at standard temperature conditions (STC) by using Eq. (3). The module efficiency also depends on various other parameters like active area of the module which is actual area on which solar radiations strike. The PR is considered to check the performance of PV modules based on actual and theoretical energy outputs [24].



Fig. 1. PV Modules Installed at Experimentation Site.



Fig. 2. Schematic Diagram of Experimental Setup.

TABLE I. PV MODULES SPECIFICATIONS, THEIR MEASURED AND RATED OUTPUT PARAMETRIC VALUES

| Parameter Name | m-Si | p-Si | a-Si |
|---|---|---|---|
| **PV Module Dimensions** | | | |
| PV module dimensions (mm × mm) | 1500×700 | 1500×700 | 1220×620 |
| Cell dimensions (mm × mm) | 155×155 | 155×150 | 1200×600 |
| No. of cells in series | 4×9 | 4×9 | 1 |
| Cell area (m$^2$) | 0.8649 | 0.837 | 0.72 |
| **Rated PV Module Specifications** | | | |
| $P_{max}$ (W) | 150 | 150 | 65 |
| $I_{max}$ (A) | 8.75 | 7.9 | 0.94 |
| $V_{max}$ (V) | 18.5 | 19 | 32.5 |
| $I_{sc}$ (A) | 9.83 | 8.5 | 0.97 |
| $V_{oc}$ (V) | 22.5 | 22 | 56.1 |
| **Measure Values of PV Module Outcomes** | | | |
| Avg. module temp. (ºC) | 36.70 | 37.47 | 37.18 |
| Avg. module current (A) | 2.15 | 1.82 | 0.44 |
| Avg. module voltage (V) | 14.52 | 14.08 | 9.83 |
| Avg. module power (W) | 34.52 | 28.62 | 4.57 |
| Avg. power output efficiency (%) | 23.01 | 19.08 | 7.03 |
| Avg. module efficiency (%) | 8.29 | 6.30 | 1.16 |
| Avg. performance ratio | 0.46 | 0.36 | 0.14 |

## III. Results and Discussions

### A. Meteorological Parameters

The PV off-grid type standalone system was designed and installed using three different PV technologies (m-Si, p-Si, a-Si). The performance of PV technology [25] is directly dependent on incident solar irradiance value. Incident solar irradiance intensity value is governed by meteorological parameters like ambient temperature, wind speed, and humidity mainly.

During the study, the daily average solar irradiance value recorded was above 500 W/m$^2$, with minimum 99.3 W/m$^2$ and maximum 683.3 W/m$^2$ values, as depicted in Fig. 3. These daily average irradiance values were in accordance with daily environmental temperature fluctuations except two days (rainy days). Because some hour measurements of these two days could not be recorded due to rainy weather. Fig. 4 shows the variations in hourly average solar irradiance for measurement days. The hourly average solar irradiance value increases linearly from morning hours with 101.07 W/m$^2$ to 1055.45 W/m$^2$ value at 12:00 p.m. hours, then decreases in the afternoon hours upto 18.49 W/m$^2$. The maximum hourly average incident solar irradiance value was recorded between 12:00-1:00 p.m. (as shown in Fig. 4).

Further the performance of three PV technologies is evaluated using following approaches to choose the optimum one:

- Module temperature analysis
- Average and normalized power output efficiency analysis and power delivered to load
- Module efficiency analysis
- Performance ratio analysis

### B. Module Temperature Analysis

The PV module surface temperature depends upon incident solar irradiance intensity. The variations in hourly average module's temperature are shown in Fig. 5.

Analysis of Fig. 5 shows that increasing or decreasing module's temperature trend is in accordance with incident solar irradiance values (evident from the comparison of Fig. 5 bars against Fig. 4). The lowest hourly average module's surface temperature is recorded for m-Si, p-Si and a-Si were equivalent to 15.12 °C, 15.09 °C, and 16.05 °C at 7:00 a.m. in the morning respectively, and the maximum modules temperature values were 50.81 °C, 50.35 °C, and 52.02 °C at 12.00 p.m. at noon respectively against highest hourly average incident solar irradiance values (as evident from Fig. 4). The Difference in surface temperature of different PV modules is very small, and that might be due to variations in glass type of glazing layer of module, as there was no shading effect on any of the module's surfaces.



Fig. 3.    Daily Avg. Solar Irradiance Values for Experimentation Days.



Fig. 4.    Hourly Avg. Incident Solar Irradiance Against Time Hour for Experimentation Days.

Fig. 5.   Hourly Avg. Module Temperature Values.

## C. *Normalized Power Output Efficiency (NPOE) and Power Delivered to Load*

Output power ($P_{mea}$) is the main parameter to check the performance of any power generation source, to be installed. The average daily and hourly based output power ($P_{mea}$) in watts of installed standalone PV generation sources is elaborated in Fig. 6 and Fig. 7, respectively. Output power in relation with maximum possible rated power from modules at STC is elaborated in Fig. 8 and Fig. 9 in form of normalized output power (Eq. (3)) against incident solar irradiance values and time, respectively.

Analysis of Fig. 6 shows that m-Si type PV module produced daily output power in 27-49 W range, while p-Si and a-Si type PV modules generated $P_{mea}$ in 21-35 W and 3-7 W ranges respectively, ignoring the output of 2 rainy days. The m-Si module generated the highest output power value of 48.8 W against 654.01 W/m² incident solar irradiance and 26 °C environmental temperature values while p-Si produced

maximum $P_{mea}$ of 34.1 W against 651.6 W/m² and 28 °C environmental temperature values. a-Si module generated in almost constant behavior with very small variations against measured range of incident solar irradiance values.

Analysis of Fig. 6 to 9 elaborates that increase in incident solar irradiance value corresponds to the increase in output power ($P_{mea}$) and normalized output power values of all three types of PV modules, with maximum values generated by m-Si type module followed by p-Si and a-Si type PV modules in descending order. Daily individual analysis is also evident of the fact that output power value is in accordance with variations in solar irradiance value with maximum power generated values during 12:00-1:00 p.m. hour (evident from Fig. 7 and Fig. 9). Analysis of output power ($P_{mea}$) and normalized output power curves conclude that m-Si PV technology is optimum type with 23.01% NPOE value, followed by p-Si and a-Si PV technologies with 19.08% and 7.03% NPOE values, respectively.



Fig. 6.   Daily based Avg. Output Power of Experimentation Days.



Fig. 7.   Hourly based Avg. Power Output of Experimentation Days.

Fig. 8.    Hourly based Avg. NPOE of PV Modules Against Avg. Incident Solar Irradiance Values.



Fig. 9.    Hourly based Avg. NPOE of PV Modules for Experimentation Days.

### D. Module Efficiency Analysis

The PV module efficiency is an indicator of overall panel's quality, tells about the percentage of incident sun light spectrum being converted into useable electric power (Eq. (4)). The module efficiency is directly related to $P_{mea}$, while is in inverse relation with product of module surface area and solar irradiance value incident on it. Average module efficiency of three types of PV technologies against incident irradiance intensity values, and time of the days is described through Fig. 10 and Fig. 11, respectively.

Analysis of the Fig. 10 and Fig. 11 elaborates that module efficiency at real-time conditions could never be the same as reported by the manufacturers, because that value is reported under STC. Analysis of Fig. 10 and Fig. 11 supports the fact that module efficiency is in inverse relation with incident solar irradiance values, but after certain intensity value. There is a critical relation between incident solar irradiance value and surface area of the module. As every PV technology type works efficiently upto the certain module surface area and temperature in relation with incident solar irradiance value. The further increase in incident irradiance value against certain area value leads to decrease in overall module efficiency value, as is happening here. In the start of efficiency curves, module efficiency increases for m-Si and p-Si types against increasing irradiance value upto certain value at early day hours. Then further increase in irradiance value caused decrease in module efficiency for both types (m-Si and p-Si) of modules. Overall analysis of Fig. 10 and Fig. 11 elaborates that both m-Si and p-Si PV module types give high efficiency at early (i.e. morning)

and late (i.e. afternoon) hours of the day with comparatively low incident solar irradiance values.

Comparative analysis of m-Si and p-Si modules shows that m-Si type PV technology gives higher efficiency values at lower incident solar irradiance values, while p-Si type modules gave higher efficiency value at relatively high incident solar irradiance values (as shown in Fig. 10). The a-Si type module shows almost constant behavior against the increase or decrease in incident solar irradiance values. Module efficiency-based analysis of three types of PV technologies to choose the optimum one among them concludes that m-Si type PV technology is the best option with 8.29% real time measured module efficiency, higher compared to that of p-Si and a-Si type PV modules with 6.3% and 1.16% module efficiency values, respectively.

### E. PR Analysis

The PR is description of real time quality factor of the PV plant/ module that tells about the relationship between actual and theoretically expected output of the PV plant (calculated using Eq. (5)). The PR value helps in estimating the final value of power would be available for exporting to grid after all loss's deduction. The hourly average (Avg.) PR against module temperature, and time is plotted in Fig. 12 and Fig. 13 respectively. Analysis of Fig. 12 and Fig. 13 elaborates that PV module types show the same trend in case of PR parameter as was for module efficiency (Fig. 10 and Fig. 11). The PR has decreasing trend in value against the increase in incident solar irradiance value. That's why m-Si and p-Si have relatively high

PR values in early (morning) and late (afternoon) hours of the day with low incident solar irradiance values, compared to the hours (noon time) with high incident solar irradiance values. Comparative analysis of three PV module types concludes that

m-Si type PV modules is the optimum choice with 46.38% PR value (0.46), higher than that of p-Si and a-Si type PV modules generating power with 36.43% (0.36) and 13.82% (0.14) PR values.



Fig. 10. Avg. Module Efficiency Against Avg. Incident Solar Irradiance.



Fig. 11. Hourly Avg. Module Efficiency Values Against Time Hours During Experimentation Days.



Fig. 12. Hourly Avg. PR Against Hourly Avg. Module Temperature Values.



Fig. 13. Hourly Avg. PR Values for Time Hours During Experimentation Days.

## IV. ECONOMIC IMPACT OF POWER SOURCE REPLACEMENT

Economy is the major concern in making any technology or project viable. The current study elaborates the performance of three types of PV technologies to recommend the optimum choice for installing in domestic sector with the purpose of satisfying domestic load demand. Economic analysis of the experimentally tested systems is done in terms of the price of the consumed power from grid in base case. Analysis reveals that the installation of 1 kW PV system using m-Si technology will generate 2.53 kWh avg. normalized power per day which is 22.2% more than the power value that could be produced through p-Si PV technology-based system (equivalent to 2.07 kWh avg. normalized power per day value). While 1 kW PV system installed using a-Si technology will generate 1.09 kWh avg. normalized power per day. In domestic sector of the country region with 3$^{rd}$ largest population [26] count (Faisalabad, Pakistan), where average power consumption per house is around 100-200 kWh per month [27]: installation of such 1 kW m-Si technology based PV system would result in 38-76% reduction in electricity bills, while p-Si technology based 1 kW PV system would result in 31-62% reduction and a-Si technology based 1 kW PV system would lead to 16-33% reduced monthly electricity bills.

## V. CONCLUSION

A standalone PV system setup was installed at the rooftop of Faisalabad, Pakistan. Three different types of PV technologies (m-Si, p-Si, a-Si) were tested to identify the impact of these renewable technologies on domestic power demand satisfaction. Within the scope of this work, the hourly output power, NPOE, module efficiency, and PR of three modules were measured, and variations in these parameters with variant solar irradiance and surface temperature of PV modules conditions were investigated. Outcomes have elaborated that the output power, and corresponding module efficiency and PR factors, increases with the increase in incident solar irradiance value and surface area of the module, but upto certain limit value. After that limit value, efficiency values start decreasing because module temperature has crossed its nominal operating temperature value.

With the aim of addressing the curiosity that how much impact a certain type of PV technology will impose on domestic sector load demand at testing location, a comparative analysis of three types of PV module has been carried out. Analysis concluded that m-Si PV technology is the optimum choice with 23.01% avg. NPOE, 8.29% avg. module efficiency, 46.68% PR values, and 2.53 kWh normalized power per day value which is 22.2% more than the power generated value compared to p-Si type. While p-Si PV technology performs with 19.08% avg. NPOE, 6.3% avg. module efficiency, 36.43% PR values, and 2.07 kWh normalized generated power per day. a-Si PV technology, although showed almost uniform behavior against variant incident solar irradiance values, but is the least efficient type with 7.03% avg. NPOE, 1.16% avg. module efficiency, 13.82% PR values, and 1.09 kWh normalized generated power per day.

Economic evaluation of PV technologies elaborates that m-Si technology will result in maximum possible reduction in electricity bills equivalent to 38-76%. Integration of renewable energy in domestic sector and installation of such systems would put an impressive environmental and economic impact towards power sector sustainability and reducing shortfall.

### REFERENCES

[1] Muhammad Arshad Shehzad Hassan, Minyou Chen, Houfei Lin, Mohammed Hassan Ahmed, Muhammad Zeeshan Khan, Gohar Rehman Chughtai, "Optimization Modeling for Dynamic Price Based Demand Response in Microgrids," Journal of Cleaner Production, vol. 222, pp. 231-241, 2019.

[2] Muhammad Arshad Shehzad Hassan, Minyou Chen, Qiang Li, M. Ali Mehmood, Tingli Cheng, Bo Li, "Microgrid Control and Protection State of the Art: A Comprehensive Overview," Journal of Electrical Systems, Vol. 14, Issue 2, PP. 148-164, 2018.

[3] Ghafoor, A., et al., Current Satus and Overview of Renewable Energy Potential in Pakistan for Continuous Energy Sustainability. Renewable and Sustainable Energy Reviews, 60 (2016), pp. 1332-1342.

[4] Mehmood, A., F.A. Shaikh, and A. Waqas. Modeling of the Solar Photovoltaic Systems to Fulfill the Energy Demand of the Domestic Sector of Pakistan using RETSCREEN Software. Proceedings (IEEE), in Green Energy for Sustainable Development (ICUE), 2014 International Conference and Utility Exhibition on, 2014.

[5] Harijan, K., M.A. Uqaili, and U.K. Mirza, Assessment of Solar PV Power Generation Potential in Pakistan. Journal of Clean Energy Technologies, 3 (2015), 1, pp. 54-56.

[6] Tahri, F., A. Tahri, and T. Oozeki, Performance Evaluation of Grid-Connected Photovoltaic Systems Based on Two Photovoltaic Module Technologies Under Tropical Climate Conditions. Energy Conversion and Management, 165 (2018), pp. 244-252.

[7] Huld, T., et al., Mapping the Performance of PV Modules, Effects of Module Type and Data Averaging. Solar Energy, 84 (2010), 2, pp. 324-338.

[8] Kawagoe, K., Y. Hishikawa, and N. Yamada, Outdoor Direct STC Performance Measurement of PV Modules Based on a Sun-Shading Technique. IEEE Journal of Photovoltaics, 7 (2017), 6, pp. 1725-1730.

[9] Ramli, M.S., S.S.A. Wahid, and K.K. Hassan. A Comparison of Renewable Energy Technologies Using Two Simulation Softwares: HOMER and RETScreen. Proccedings, AIP Conference, 2017, AIP Publishing.

[10] Kroposki, B., et al. A Comparison of Photovoltaic Module Performance Evaluation Methodologies for Energy Ratings. Proceedings (IEEE), in Photovoltaic Energy Conversion, Conference Record of the Twenty Fourth. IEEE Photovoltaic Specialists Conference-1994, 1994 IEEE First World Conference on.

[11] Amin, N., C.W. Lung, and K. Sopian, A Practical Field Study of Various Solar Cells on Their Performance in Malaysia. Renewable Energy, 34 (2009), 8, pp. 1939-1946.

[12] Midtgard, O.-M., et al., A Qualitative Examination of Performance and Energy Yield of Photovoltaic Modules in Southern Norway. Renewable Energy, 35 (2010), 6, pp. 1266-1274.

[13] Rahman, M., M. Hasanuzzaman, and N. Rahim, Effects of Various Parameters on PV-Module Power and Efficiency. Energy Conversion and Management, 103 (2015), pp. 348-358.

[14] Singh, J.P., et al., Comparison of Glass/Glass and Glass/Backsheet PV Modules using Bifacial Silicon Solar Cells. IEEE Journal of Photovoltaics, 5 (2015), 3, pp. 783-791.

[15] Catelani, M., et al. Characterization of Photovoltaic Panels: The effects of Dust. Proceedings (IEEE), International Energy Conference and Exhibition (ENERGYCON), 2012.

[16] Mekhilef, S., R. Saidur, and M. Kamalisarvestani, Effect of Dust, Humidity and Air Velocity on Efficiency of Photovoltaic Cells. Renewable and Sustainable Energy Reviews, 16 (2012), 5, pp. 2920-2925.

[17] Ferrada, P., et al., Performance Analysis of Photovoltaic Systems of Two Different Technologies in a Coastal Desert Climate Zone of Chile. Solar Energy, 114 (2015), pp. 356-363.

[18] Skoplaki, E. and J.A. Palyvos, On the Temperature Dependence of Photovoltaic Module Electrical Performance: A Review of Efficiency/Power Correlations. Solar energy, 83 (2009), 5, pp. 614-624.

[19] Li, X., et al., Outdoor Performance and Stability Under Elevated Temperatures and Long - Term Light Soaking of Triple - Layer Mesoporous Perovskite Photovoltaics. Energy Technology, 3 (2015), 6, pp. 551-555.

[20] Tomar, V., et al., Thermal Modeling and Experimental Evaluation of Five Different Photovoltaic Modules Integrated on Prototype Test Cells with and without Water Flow. Energy Conversion and Management, 165 (2018), pp. 219-235.

[21] Jiang, H., L. Lu, and K. Sun, Experimental Investigation of the Impact of Airborne Dust Deposition on the Performance of Solar Photovoltaic (PV) Modules. Atmospheric Environment, 45 (2011), 25, pp. 4299-4304.

[22] Bashir, M.A., et al., An Experimental Investigation of Performance of Photovoltaic Modules in Pakistan. Thermal Science, 19 (2015), Supplement 2, pp. 525-534.

[23] Green, M.A., et al., Solar Cell Efficiency Tables (Version 45). Progress in Photovoltaics: Research and Applications, 23 (2015), 1, pp. 1-9.

[24] Cañete, C., J. Carretero, and M. Sidrach-de-Cardona, Energy Performance of Different Photovoltaic Module Technologies under Outdoor Conditions. Energy, 65 (2014), pp. 295-302.

[25] Said, Z. and A. Mehmood, Standalone Photovoltaic System Assessment for Major Cities of United Arab Emirates Based on Simulated Results. Journal of Cleaner Production, 142 (2017), pp. 2722-2729.

[26] PBS, Populaiton, Pakistan Bureau of Statics, Government of Pakistan. 2018.

[27] FESCO, Domestic electricty billing, Faisalabad Electric Supply Company (FESCO), WAPDA, Government of Pakistan. 2018.

# MHealth for Decision Making Support: A Case Study of EHealth in the Public Sector

Majed Kamel Al-Azzam[1]
Business Administration Department
Yarmouk University
Irbid, Jordan

Malik Bader Alazzam[2]
Software Engineering Department
Ajloun National University
Jordan

Majida Khalid al-Manasra[3]
Biomedical and Communication
Engineering, Summit International
Academy-Jordan

*Abstract*—**This paper seeks to explore factors that determine the acceptance of the MHealth application patients. The research relied on (UTAUT2) Unified Theory of Acceptance and Use of Technology to assess the level of acceptance of a new mobile health application by patients. The study involved conducting test surveys across medical hospitals in Jordan with the goal of collecting data from hospital visitors and their patients concerning their intention to use the new mobile health application. 98 questionnaires were collected and 44 valid responses drawn from them for onward data analysis. The UTAUT2 research model was the most appropriate one for conducting the evaluation on MHealth's user acceptance. Its results would support the government's goal of building m-health solutions that meet user needs. The model also enhances the roles of DSS in facilitating adoption of MHealth applications. This study provides a theoretical framework for pursuing future research work on the rates of adoption of m-health applications by patients.**

*Keywords—Mobile health application; UTAUT1; UTAUT2; trust factors*

## I. INTRODUCTION

Mobile technology is growing rapidly [1][2], prompting more healthcare organisations to consider mobile health technologies to be a feasible solution for monitoring patients' status. This development has formed the possible for converting healthcare system supply into a more available, reasonable, and active form [3][4][5]. The electronic health system (e-health system) used in the traditional setting are faulted for their reliance on wired connections and computers. As an improvement, mobile health uses wireless cellular communication to achieve mobility and flexibility [6][7]. It also has the advantage of portability and long-life battery power. Due to their mobility, portability, and flexibility of mobile technologies, they are preferred over traditional systems when seeking to enhance access to healthcare services. They also help to reduce the costs and time incurred in delivering healthcare [8]. These systems may also be used to create motivation amongst healthcare professionals to stick to professional habits by reshaping [9][10]. Mobile health system is gaining increased prominence as a favoured technology on health communication. Its importance goes beyond more communication by including aspects such as management [11], facilitation and delivery of health information through monitor, cell phones, wireless infrastructure, tablets, and sensors. Mobile healthcare technologies include many healthcare services and applications such as site-based health services, mobile telemedicine, and pervasive information access to healthcare systems, and patient monitoring [12][13] [14]. They bring great benefits to both patients and physicians.

In addition the potential financial and medical additions of mobile health service area, the implementation of the use of MHA opposite tasks and walls at the social, technological, culture, governmental and governmental points, particularly in developing nations [15] [16] after interviews with Patients and reviewers and health organizations there are several barriers. one of the key notes that we observed is the lack of the cognizance about mobile health applications (MHA) [17][14] , facilities and its benefits. Second, there are concerns with act of application. Third, patients do not neediness to change the routine of providing healthcare system, equal with the probable assistances of mobile health app [17] [18] [19]. Several patients sense that physicians are so busy and do not have period to use the app to follow up and remotely display them. Fourth, there is a unlimited mission with respect to the belief, extra than one patient do not belief any administration system.

Lastly, it is hard to accept mobile health app without there is community acceptance also official support. The acceptance of mobile health system depends on several factors beyond the technology's skills. These issues include the willingness of the patients, health professionals, and care contributors to adopt embrace new technology [20][21][22][23][24]. To some extent, it also depends on the level to which the management will provide the necessary support. It is vital to remember that any role of DSS is to help support the aims of the healthcare organizations.

When creating m-health system, an institution needs to have information on the potential for its users' acceptance. In this regard, they need to rely on proper research on the established factors that influence users to adopt or reject new mobile healthcare technology [25]. This is especially important in the modern healthcare setting, as most hospitals strive to encourage the public and their staff to embrace new interventions for health promotion as a way of reducing healthcare costs and enhancing the overall health standards.

The researchers from several disciplines have been reviewing computerized DSS for around 55 years [26][27][28]. DSS frameworks can be distributed into five types including: first–Model(driven DSS), second(Data-driven

DSS) [29][26], third Communications (driven DSS), forth Document (driven DSS), fifth Knowledge (driven DSS).

- Model driven DSS emphasize admission to and handling of economic, or simulation models. Model-driven DSS use incomplete factors delivered by decision makers to assistance decision makers in investigating a state. Academics motivated on model management and on improving further diverse styles of models for use in DSS such as optimization, and simulation models.

- data-driven DSS give emphasis to admission to manipulation of a time series of internal establishment data and at times external , real-time data. Simple file systems accessed by query and recovery tools provide the most basic level of functionality.

- Communications-driven DSS- use system and communications technologies to enable and communication. In these systems, communication technologies are the dominant architectural factor. Tools used contain groupware. In overall, groupware, bulletin boards, audio and videoconferencing.

- Document-driven DSS- usages handling technologies to deliver document retrieval and analysis. Vast document databases may hold scanned documents, images, sounds, and video. Cases of documents that might be accessed by a document-driven DSS are rules and techniques, product specifications, corporate historical documents.

- Knowledge-driven DSS- can propose actions to managers. These DSS are person-computer systems with intensive problem-solving expertise [30][31]. The expertise involves of knowledge about an exact domain, considerate of problems inside domain.

*A. Review of Decision Support System in Healthcare Sector*

There are two main areas for making primary decisions in the healthcare setting. The first lower level area has to do with diagnosis and treatment, patient management, finance management, inventory and record keeping The second higher level area is meant to position the healthcare facility in a competitive place relative to its competitors [32]. In this area, some of the functions handled include patient management and inventory management. This second higher level of decision making is mostly meant to benefit shareholders, while the lower level decisions is mainly targeted to improve the work of nurses and doctors. An example of a decision support system in the healthcare setting is the PRODIGY [33]. This DSS allows healthcare professionals to access evidence and knowledge on disease symptoms and conditions when delivering primary healthcare [34]. PRODIGY is an acronym for Prescribing Rationally with Decision support In General Practice Study. The system provides full text guidance, drug information, patients information leaflets, self-help contacts, and quick reference guides to pharmacists, nurses and patients. Other important information maintained by PRODIGY include finance information such as accounts receivable, track billing, accounts payable and payroll; patient insurance policy, insurance payment options, and insurance claims.

MedSphere also offers decision support systems that have various modules [35][10]. Their systems have modules that are integrated to deliver a comprehensive functionality. They capture important information on the patient's billing as they move from registration, diagnosis, prescription, admission and discharge. It also has modules for handling collection processing, Supply chain management, and other features. It gives instantaneous view on the financial state of the hospital at any given time [7]. DSS supports patient diagnosis, as it displays patients' background information and provides clinicians with knowledge on symptoms for various health problems and their recommended treatment processes [36]. As showed in Fig. 1, clinician can use DSS to determine proper drug usage, diagnose the health problem of a patient, and send reminders to other staff to administer drugs to hospitalised patients on time. Another example of such a DSS is known as Isabel. This system has a database of medical records for patients which are accessible over the web.



Fig. 1.   DSS Progress Example [27].

## II. Literature Review

### A. UTAUT1

More research has been conducted on the technology acceptance models. Different models and theories propel these studies. The researchers of UTAUT model have identified and harmonized 8 different models and theories that form a comprehensive acceptance model. These theories are [11], [37]: Social Cognitive Theory-SCT-, Technology Acceptance Model –TAM-, Motivational Model -MM-, Theory Reason Action TRA-, and Theory of Planned Behavior- TPB, Innovation Diffusion Theory –IDT.

The unification of these studies make a summation of all the constructs from the eight models into four different determinants which help in predicting the usage, the intentions, and the specific moderators of all important relationships (Morris et al., 2003) [38]. Fig. 2 makes an illustration of the relationships inherent within the UTAUT [13][39]. This model is composed of 4 exogenous variables (EV) that are namely; performance expectancy (PE), effort expectancy (EE), facilitating conditions (FC) and social influence (SI) [39]. These exogenous variables are mainly used in technology intention to the behavior and usage. Among these exogenous variables [38][37], there are four moderators which include age, gender, experience and voluntariness.

### B. UTAUT2

The study has extended the unified theory and the technology model by examining how technology is used and accepted by consumers as shown in Fig. 3. The goals of UTAUT2 align with those of UTAUT1 and the concepts of HM, PV and HT[40]. According to this study by Venkatesh, the user's demographic characteristics became the moderator variables[41][1][40] They included experience, age, and gender; and how they affected technology use.

### C. Hedonis Motivation (HM)

It is described as the intrinsic motivation of a user of technology [42][43][13]. It is considered critical in constructing a model to determine use and acceptance of technology. HM can be compared to playfulness or the user's enjoyment of TAM as a factor that has intrinsic value.

### D. Price Value (PV)

People are known to choose the products and services that benefit them more than price value. As such [44][45], one can define price value as conscious trade-off that people make between the money costs of acquiring a new application and the perceived benefits that they would derive from the application.



Fig. 2. UTAUT1 Model.



Fig. 3. UTAUT 2 Model.

### E. Mobile Health Application

EHealth-related information can be collected, stored, and exchanged effectively through a range of tools provided through various information and communication technologies (ICTs)[3]. Healthcare can harness these technologies to enhance their level of safety, quality and cost performance. MHealth is considered as one of the key applications in healthcare delivery, as it combines several functions such as telemedicine, electronic prescription, test ordering, emergency information and digital imaging. These elements help the healthcare professionals to obtain medical evidence in a reliable manner. Such evidence supports the healthcare organizations in delivering their clinical mandate. MHealth application also enhances service delivery and organisational efficiency in the healthcare setting. Different healthcare stakeholders such as professionals, patients, general public and organizations can use the MHealth application[46]. Different studies indicate that MHealth applications bring great benefit to patients. One such benefit is that it raises the quality of care by providing easy access to important health data that a patient may need from different health service providers. Most MHealth applications are also modelled to support disease management programs that bring great benefit to patients.

Table I summarises related studies on the topic of user's intention to adopt new MHealth applications.

TABLE I. MHEALTH APPLICATION REVIEW

| Author | Description | Origin | Method, sample size | Model | Results |
|---|---|---|---|---|---|
| [19] | Examines the user adoption of a new tablet application that aims to provide support for cognitive stimulation for the elderly | Paris | , survey , 15 senior users | Cognitive Therapy | good acceptability of the app's games that continues and improves with time |
| [47] | Applying an Acceptance Model to assess the adoption of M-Health Services by health related users in UAE | UAE | Survey , 144 | (TAM) | model. PU, PEOU, TR and SE found directly influencing the intention to use M-Health system |
| [9] | Achieving privacy and security in MHealth applications | USA, EU | Review and Recommendations | Review and Recommendations | |
| [16] | Assessing the current state of the art in mobile health-related and clinical apps | USA, Europe, | brief survey of evaluation studies | evaluation, regulation and certification, quality | Interactions may require substantial effort. |
| [48] | | Developing Countries | CASE STUDY | capacity for improved access | RECOMANDATIONS |
| | consumer's acceptance of mobile technology | Egypt and Yemen | 302 survey | (TAM) | Positively Resistance to change, Technology anxiety factors |
| [12] | Creating a summary of 7 strategies for conducting evaluation and selection of health-related apps: | | Interview .1 | Case study | |
| [26] | Developing and performing user evaluation for a mobile DSS running on iOS known as OphthalDSS | Spain | Survey, 50 physicians answered | Quality, Ease of Use, Availability, Performance | Positively Quality, Ease of Use, Availability, Performance |

## III. RESEARCH METHODS AND RESULTS

### A. Pilot Study for MHealth Application by UTAUT2

Quantitative research methods were used for primary research in this study. This approach helps to generate contextual information on user acceptance for MHealth applications across Jordan. It will also give background information on the use of UTAUT2 in the assessment of user acceptance for new technology [34][49]. The research will also use a correlational study design to determine whether the conceptual model has any relationships that can be interpreted as independent variables and dependent variables [6]. The

section below describes the detail of the phases used in the study. 9 constructs from the technology acceptance model were measured using 32 items by following the guidelines UTAUT2. The study also involved the collection of the necessary demographics on the number of users that had embraced the new M-health application across Jordan. The demographics were intended for use in comparing different levels of user acceptance for MHealth applications in hospitals across Jordan [50]. The study recognised that the adoption of information technology is largely connected to most business activities and services.

TABLE II.    SAMPLES DEMOGRAPHIC

| Variable | Description | Frequency | Percentage |
|---|---|---|---|
| Gender | male | 28 | 63.6 % |
|  | Female | 16 | 36.4% |
| Age | 21-30 | 6 | 13.65 |
|  | 31-40 | 14 | 31.80% |
|  | 41-50 | 20 | 45.50% |
|  | 51-60 | 4 | 9.10% |

TABLE III.    VARIABLES

| Variable Group | No. Items | Item | Min | Max | Mean | Std. Deviation |
|---|---|---|---|---|---|---|
| Performance expectancy | 4 | I would find MHEALTH (useful) in my job | 2.00 | 5.00 | 4.0554 | .89853 |
|  |  | Using MHEALTH (increases) my chances of achieving things that are vital to me | 2.00 | 5.00 | 3.5455 | 1.22386 |
|  |  | Using MHEALTH aids me accomplish things more rapidly. | 1.00 | 5.00 | 3.8081 | .97145 |
|  |  | Using MHEALTH in my job | 2.00 | 5.00 | 3.6818 | 1.39340 |
| Effort Expectancy | 4 | Learning to operate MHEALTH app would be relaxed | 1.00 | 5.00 | 3.8081 | 1.30600 |
|  |  | My interaction with MHEALTH app is clear | 1.00 | 5.00 | 3.7708 | 1.42716 |
|  |  | I find MHEALTH app easy to use | 1.00 | 5.00 | 3.5455 | 1.40500 |
|  |  | It is relaxed for me to develop expert at using MHEALTH application | 1.00 | 5.00 | 3.7702 | 1.46828 |
| Social Influence | 3 | Persons who are significant to me reflect that I should use MHEALTH application | 1.00 | 5.00 | 3.2044 | 1.21052 |
|  |  | Persons who effect my behaviour think to use MHealth application | 1.00 | 5.00 | 3.2201 | 1.27920 |
|  |  | Persons whose views that x value select that I use MHEALTH application | 1.00 | 5.00 | 3.6601 | 1.16217 |
| Facilitating condition | 4 | I have the resources essential to use MHEALTH | 1.00 | 5.00 | 3.5909 | 1.05375 |
|  |  | I have the knowledge needed to use MHEALTH | 1.00 | 5.00 | 3.9545 | .89853 |
|  |  | MHEALTH is compatible with other tools I use | 1.00 | 5.00 | 3.5213 | 1.00755 |
|  |  | get help from experts when I have problems using MHEALTH application | 1.00 | 5.00 | 3.4440 | 1.29685 |
| Hedonic Motivation | 3 | MHealth application is so fun. | 1.00 | 5.00 | 3.0114 | 1.42413 |
|  |  | MHealth application is so enjoyable | 1.00 | 5.00 | 3.0124 | 1.31590 |
|  |  | MHealth application is so entertaining. | 1.00 | 5.00 | 2.8541 | 1.46311 |
| Price Value | 3 | MHEALTH application is judiciously priced | 1.00 | 4.00 | 3.4471 | 1.21677 |
|  |  | MHEALTH application is a not bad value for the money | 1.00 | 5.00 | 3.8854 | 1.46015 |
|  |  | existing price, MHealth application provides a respectable value | 1.00 | 5.00 | 3.6547 | 1.54863 |
| Habit | 4 | use of MHEALTH application has become a nice and habit | 1.00 | 5.00 | 3.4401 | 1.09801 |
|  |  | I am addicted to using MHealth app | 1.00 | 5.00 | 3.8801 | 1.13294 |
|  |  | I must use MHealth app | 1.00 | 5.00 | 3.3182 | 1.21052 |
|  |  | Using MHealth app has become natural to me | 2.00 | 5.00 | 4.0014 | 1.15470 |
| Behaviour | 3 | I plan to remain using MHealth app in the future | 1.00 | 5.00 | 3.2014 | 1.17422 |
|  |  | I will continuously try to use MHealth app in my daily life | 1.00 | 5.00 | 3.8824 | 1.34277 |
|  |  | I plan to continue to use MHealth app frequently | 1.00 | 5.00 | 3.7547 | 1.24924 |
| Intention to use | 3 | I frequently used MHEALTH app to understand health problem | 1.00 | 5.00 | 3.4401 | 1.24924 |
|  |  | I often use MHEALTH app to serve patient | 2.00 | 5.00 | 3.8792 | 1.21677 |
|  |  | I often use MHEALTH app to establish information about health issues problem | 2.00 | 5.00 | 3.8821 | 1.20317 |

TABLE IV.    CRONBACH ALPHA

| Variables | N of (Items) | N of delete (items) | Cronbach C Alpha A(CA) |
|---|---|---|---|
| Performance P expectancy E (PE) | 4 | 0 | 0.91 |
| Effort Expectancy(EE) | 4 | 0 | 0.901 |
| Social Influence(SI) | 3 | 0 | 0.875 |
| Facilitating condition(FC) | 4 | 0 | 0.876 |
| Hedonic Motivation(HM) | 3 | 0 | 0.952 |
| Price value(PV) | 3 | 0 | 0.932 |
| Habit(H) | 4 | 0 | 0.876 |
| Behaviour(B) | 3 | 0 | 0.941 |
| Intention to use(IU) | 3 | 0 | 0.917 |

The study identified 8 important success factors that researchers need to look out for when using the UTAUT2 model. These include Hedonic Motivation, Performance Expectancy, Social Influence, Effort Expectancy, Price Value, Habit, Facilitating Conditions, and Behavior. In light of this model, the researcher sought to use the Acceptance and Use of Technology approach to assess the perception that hospital visitors and patients have about the use of MHealth applications. More than 70 users participated in the survey that sought to examine their intention to use MHealth applications. The users were drawn from two major hospitals across Jordan[51][7][43][52]. The total number of valid responses collected from the survey was 44. This information is presented in Table II. Most of the valid responses were submitted by women (63.6%) rather than men (36.4%). Table III presented the min and max and mean of valid variables. And Table IV showed the Cronbach Alpha values.

### B. Reliability Test

A reliability test was conducted in this study to measure the validity or acceptability of the measures. Cronbach's alpha test was used to determine the degree of internal consistency in the data. It sought to ascertain that the data used in the study had an appropriateness value of above 0.7 (Sekaran 2003). The results in this test indicate that the data met the basic threshold for internal consistency and the other factors used in the study were reliable.

## IV. CONCLUSION AND LIMITATION

More developing countries continue to adopt m-health applications as a critical technology that would drive a positive change in their healthcare services delivery approach. It is required to know the questions that affect related health users to approve or reject mobile health system. This research presented a model for MHA acceptance founded on UTAUT2 and explored the rationale of technology adoption for users in the mobile health system. This model united cultural, social, technological, political, and structural sides. DSS is important in conducting medical diagnosis, as it documents information on health problems and patients' background information that can be used by clinicians to quickly identify a patient's specific ailments. Limitation of this study this study have been conducted in two Jordan hospitals only. Researchers only measured trust factors and UTAUT factors. This sets a foundation for future research on patients' adoption of MHealth applications.

## V. ACKNOWLEDGMENT

### REFERENCES

[1] M. B. Alazzam, Y. M. Al-sharo, and M. K. Al-, "Developing ( UTAUT 2 ) Model of Adoption Mobile Health Application in Jordan e-Government," J. Theor. Appl. Inf. Technol. 30th, vol. 96, no. 12, 2018.

[2] A. M. B. Al-azzam, Majed Kamel, "Smart City and Smart-Health Framework , Challenges and Opportunities," Int. J. Adv. Comput. Sci. Appl., vol. 10, no. 2, pp. 171–176, 2019.

[3] M. Technologies, S. D. Impact, and M. I. Usage, "Mobile Technologies and Services Development Impact on Mobile Internet Usage in Latvia « Mobile Technologies and Services Development Impact on Mobile Internet Usage in Mobile Technologies and Services Development Impact on Mobile Internet Usage in Latvia," 2013.

[4] A. Dashti, I. Benbasat, and A. Burton-jones, "Working Papers on Information Systems Trust , Felt Trust , and E-Government Adoption : A Theoretical Perspective," vol. 10, no. 2010.

[5] A. Spring, A. Sukkar, and A. World, "Mobile Crowdsourcing Technology Acceptance and use in the Crises Management of Arab Spring Societies," no. June, pp. 901–910, 2014.

[6] J. Lee and M. J. Rho, "Perception of Influencing Factors on Acceptance of Mobile Health Monitoring Service: A Comparison between Users and Non-users.," Healthc. Inform. Res., vol. 19, no. 3, pp. 167–76, Sep. 2013.

[7] L. Eljawad, R. Aljamaeen, M. K. Alsmadi, I. Al-marashdeh, and H. Abouelmagd, "Arabic Voice Recognition Using Fuzzy Logic and Neural Network 1," vol. 14, no. 3, pp. 651–662, 2019.

[8] H. Cicibas and T. Internet, "Current and Emerging mHealth Technologies," pp. 283–302, 2018.

[9] I. De, T. Díez, M. Lopez-coronado, and M. López-coronado, "Privacy and Security in Mobile Health Apps : A Review and Recommendations Privacy and Security in Mobile Health Apps : A Review and Recommendations," no. October 2017, 2014.

[10] Y. Mohammad Al-Sharo, G. Shakah, M. Sh Alkhaswneh, B. Zeyad Alju-Naeidi, and M. Bader Alazzam, "Classification of big data: machine learning problems and challenges in network intrusion prediction," Int. J. Eng. Technol., vol. 7, no. 4, pp. 3865–3869, 2018.

[11] M. Doheir, B. Hussin, A. Samad, H. Basari, and M. B. Alazzam, "Structural Design of Secure Transmission Module for Protecting Patient Data in Cloud-Based Healthcare Environment," Middle-East J. Sci. Res., vol. 23, no. 12, pp. 2961–2967, 2015.

[12] S. Mullen and S. Pagoto, "Evaluating and selecting mobile health apps: strategies for healthcare providers and healthcare organizations," no. May 2012, 2017.

[13] M. B. Alazzam, "Physicians' Acceptance of Electronic Health Records Exchange: An Extension of the with UTAUT2 Model Institutional Trust," Adv. Sci. Lett., vol. 21, pp. 3248–3252, Feb. 2015.

[14] S. Vongjaturapat, "Mobile Technology Acceptance for Library Information Service : A Theoretical Model," pp. 290–292, 2013.

[15] L. M. Telefons, "Mobile Technologies and Services Development Impact on Mobile Internet Usage in Latvia," vol. 1142, 2013.

[16] M. N. K. Boulos, A. C. Brewer, C. Karimkhani, D. B. Buller, and P. Robert, "Mobile medical and health apps : state of the art , concerns , regulatory control and certification," vol. 5, no. 3, pp. 1–23, 2014.

[17] C. Y. Tang, C. C. Lai, C. W. Law, M. C. Liew, and V. V. Phua, "Examining key determinants of mobile wallet adoption intention in Malaysia: an empirical study using the unified theory of acceptance and use of technology 2 model," Int. J. Model. Oper. Manag., vol. 4, no. 3, pp. 248–265, 2014.

[18] R. Breitschwerdt, R. Iedema, S. Robert, A. Bosse, and O. Thomas, "Health Information Technology in the International Context," Adv. Health Care Manag., vol. 12, pp. 171–187, 2012.

[19] M. Yasini and G. Marchand, "Adoption and Use of a Mobile Health Application in Older Adults for Cognitive Stimulation," 2016.

[20] L. Wang, "The Important of Enjoyment and Mobility for Continuance with Mobile Data Services," 2014.

[21] A. H. H. M. Mohamed, H. Tawfik, D. Al-Jumeily, and L. Norton, "MoHTAM: A Technology Acceptance Model for Mobile Health Applications," 2011 Dev. E-systems Eng., pp. 13–18, Dec. 2011.

[22] A. H. H. M. Mohamed, H. Tawfik, L. Norton, and D. Ai-jumeily, "Does e-Health technology design affect m-Health informatics acceptance ? A case study using a particular system would enhance his or her job," vol. 25, no. C, pp. 968–971, 2012.

[23] K. M. Unertl, K. B. Johnson, and N. M. Lorenzi, "Health information exchange technology on the front lines of healthcare: workflow factors and patterns of use.," J. Am. Med. Inform. Assoc., vol. 19, no. 3, pp. 392–400, 2012.

[24] T. R. Campion, A. M. Edwards, S. B. Johnson, and R. Kaushal, "Health information exchange system usage patterns in three communities: practice sites, users, patients, and data.," Int. J. Med. Inform., vol. 82, no. 9, pp. 810–20, Sep. 2013.

[25] M. Russell and J. M. Brittain, "Health informatics," vol. 36, no. Monograph, pp. 591–628, 2002.

[26] I. De et al., "mHealth App for iOS to Help in Diagnostic Decision in Ophthalmology to Primary Care Physicians mHealth App for iOS to Help in Diagnostic Decision in Ophthalmology to Primary Care Physicians," no. October, 2017.

[27] K. Rajalakshmi and S. C. Mohan, "Decision Support System in Healthcare Industry," vol. 26, no. 9, pp. 42–44, 2011.

[28] A. Wolfenden, "Factors Predicting Oncology Care Providers' Behavioral Intention To Adopt Clinical Decision Support Systems," no. January, 2012.

[29] "An Empirical Analysis of Citizens ' Acceptance Decisions of Electronic-Government Services : A Modification of the Unified Theory of Acceptance and use of Technology (UTAUT) Model to Include Trust as a basis for Investigation by Lawrence J . Awuah RICHA," 2012.

[30] Y. Arshad and A. R. Ahlan, "Understanding ITO decisions and implementations in Malaysia public healthcare sector: The evidence from a pilot case study," 2011 Int. Conf. Res. Innov. Inf. Syst., pp. 1–6, Nov. 2011.

[31] M. Sambasivan, P. Esmaeilzadeh, N. Kumar, and H. Nezakati, "Intention to adopt clinical decision support systems in a developing country: effect of physician's perceived professional autonomy, involvement and belief: a cross-sectional study.," BMC Med. Inform. Decis. Mak., vol. 12, no. 1, p. 142, Jan. 2012.

[32] M. B. Alazzam, "A Proposed Framework to Investigate the User Acceptance of Personal Health Records in Malaysia using UTAUT2 and PMT," Int. J. Adv. Comput. Sci. Appl., vol. 8, no. 3, pp. 386–392, 2017.

[33] M. B. Alazzam, "A Proposed Framework to Investigate the User Acceptance of Personal Health Records in A Proposed Framework to Investigate the User Acceptance of Personal Health Records in Malaysia using UTAUT2 and PMT," Int. J. Adv. Comput. Sci. Appl., no. March, 2017.

[34] M. R. Ramli, Z. A. Abas, M. I. Desa, Z. Z. Abidin, and M. B. Alazzam, "Enhanced convergence of Bat Algorithm based on dimensional and inertia weight factor," J. King Saud Univ. - Comput. Inf. Sci., 2018.

[35] L. Carter and V. Weerakkody, "E-government adoption: A cultural comparison," Inf. Syst. Front., vol. 10, no. 4, pp. 473–482, 2008.

[36] A. Karahoca, D. Karahoca, and M. Aksöz, "Examining intention to adopt to internet of things in healthcare technology products," Kybernetes, vol. 47, no. 4, pp. 742–770, 2018.

[37] M. B. Alazzam, A. Samad, H. Basari, and A. S. Sibghatullah, "Trust in stored data in EHRs acceptance of medical staff : using UTAUT2," Int. J. Appl. Eng. Res., vol. 11, no. 4, pp. 2737–2748, 2016.

[38] S. M.Alazzam, BASARI, "EHRs Acceptance in Jordan Hospitals By UTAUT2 Model: Preliminary Result," J. Theor. Appl. Inf. Technol., vol. 3178, no. 3, pp. 473–482, 2015.

[39] M. Rasmi, M. B. Alazzam, M. K. Alsmadi, A. Ibrahim, R. A. Alkhasawneh, and S. Alsmadi, "Healthcare professionals ' acceptance Electronic Health Records system : Critical literature review ( Jordan case study ) Healthcare professionals ' acceptance Electronic Health Records system : Critical literature review ( Jordan case study )," Int. J. Healthc. Manag., vol. 0, no. 0, pp. 1–13, 2018.

[40] A. S. MB.Alazzam, "Review of Studies With Utaut As Conceptual Framework," Eur. Sci. J., vol. 10, no. 3, pp. 249–258, 2015.

[41] S. Yang, "Understanding Undergraduate Students' Adoption of Mobile Learning Model: A Perspective of the Extended UTAUT2," J. Converg. Inf. Technol., vol. 8, no. 10, pp. 969–979, May 2013.

[42] E. L. Slade and M. Williams, "An extension of the UTAUT 2 in a healthcare context An extension of the UTAUT 2 in a healthcare context," 2013.

[43] M. B. Alazzam, A. B. D. Samad, H. Basari, and A. Samad, "PILOT STUDY OF EHRS ACCEPTANCE IN JORDAN HOSPITALS BY UTAUT2," vol. 85, no. 3, 2016.

[44] G. Rodrigues, J. Sarabdeen, and S. Balasubramanian, "Factors that Influence Consumer Adoption of E-government Services in the UAE: A UTAUT Model Perspective," J. Internet Commer., vol. 15, no. 1, pp. 18–39, 2016.

[45] F. Fares, A. Mashagba, and M. Othman, "Modified UTAUT Model to Study the Factors Affecting the Adoption of Mobile Banking in Jordan," Int. J. Sci. Basic Appl. Res., pp. 83–94, 2012.

[46] M. El-wajeeh, "Technology Acceptance Model for Mobile Health Systems 1, 2, 3," vol. 1, no. 1, pp. 21–33, 2014.

[47] C. Paper, "Technology Acceptance Model for the Use of M- Health Services among Health Related Users in UAE," no. December, 2015.

[48] P. N. Mechael, "The Case for mHealth in Developing Countries," pp. 103–118, 2009.

[49] M. B. Alazzam, "Theories and factors applied in investigating the user acceptance towards personal health records : Review study Theories and factors applied in investigating the user acceptance towards personal health records : Review study," Int. J. Healthc. Manag., vol. 0, no. 0, pp. 1–8, 2017.

[50] M. Ally and M. Gardiner, "The moderating influence of device characteristics and usage on user acceptance of Smart Mobile Devices," no. 2010, pp. 1–10, 2012.

[51] A. A. Alalwan, A. M. Baabdullah, N. P. Rana, K. Tamilmani, and Y. K. Dwivedi, "Examining adoption of mobile internet in Saudi Arabia: Extending TAM with perceived enjoyment, innovativeness and trust," Technol. Soc., vol. 55, pp. 100–110, 2018.

[52] M. B. Alazzam, "Factors Influencing Medical Professional Adoption of Electronic Health Record in Jordan Hospital," UTeM University, 2017.

# Sea Lion Optimization Algorithm

Raja Masadeh[1]

Computer Science Department, The World Islamic Sciences
and Education University, Amman, Jordan

Basel A. Mahafzah[2], Ahmad Sharieh[3]

Computer Science Department
The University of Jordan, Amman, Jordan

*Abstract*—**This paper suggests a new nature inspired metaheuristic optimization algorithm which is called Sea Lion Optimization (SLnO) algorithm. The SLnO algorithm imitates the hunting behavior of sea lions in nature. Moreover, it is inspired by sea lions' whiskers that are used in order to detect the prey. SLnO algorithm is tested with 23 well-known test functions (Benchmarks). Optimization results show that the SLnO algorithm is very competitive compared to Particle Swarm Optimization (PSO), Whale Optimization Algorithm (WOA), Grey Wolf Optimization (GWO), Sine Cosine Algorithm (SCA) and Dragonfly Algorithm (DA).**

*Keywords*—*Optimization; Metaheuristic optimization algorithms; Benchmarks; Sea Lion Optimization Algorithm (SLnO)*

## I. INTRODUCTION

Metaheuristic optimization algorithms are becoming more popular in application because they depend on simple concepts and easy to implement. They do not demand gradient information. They can bypass local optima and they can be applied in a wide range of issues covering various disciplines [1-5].

Metaheuristic optimization algorithms are introduced in order to solve optimization problems by imitating physical or biological phenomena [6-11]. Therefore, these algorithms are categorized into three classes; evolution- based, physics-based, and swarm-based methods [1, 2, 12-13]. Evolution-based techniques are inspired by the natural evolution' laws. The search operation begins by randomly generating population that is improved is through subsequent descent. Usually, these techniques are characterized by combining the best individuals to form the next individuals' generation. This leads the population over the generations. The most common algorithms of evolution-inspired are Genetic Algorithms (GA) [14], Evolution Strategy (ES) [15], Genetic Programming (GP) [16], Biogeography-Based Optimizer (BBO) [17] and Probability-Based Incremental Learning (PBIL) [18].

Physics-based methods mimic the physical principles in the world. Some of the most common techniques are Ray Optimization (RO) [19], Black Hole (BH) [20], Small-World Optimization Algorithm (SWOA) [21], Simulated Annealing (SA) [22], Big-Bang Big-Crunch (BBBC) [23], Gravitational Search Algorithm (GSA) [24], Charged System Search (CSS) [25] and Curved Space Optimization (CSO) [26].

Swarm- based methods are the third class of nature inspired techniques which imitate the social behavior of animals in nature. The most common technique is Particle Swarm Optimization (PSO) [27] which is mimics the bird flocking's social behavior. PSO employs number of particles which indicate to the candidate solutions that wing in the search space in order to detect the best solution that represent the optimal solution. Moreover, at the same time, they all track the best solution in their routes. Ant Colony Optimization (ACO) algorithm [28] is considered as another common swarm-based technique. ACO imitates the social behavior of ants in their colony. The most significant characteristic of ants is in finding the nearest route from the colony to the food's source; which is the major inspiration of this technique. New metaheuristic optimization algorithm is proposed by [29]. The proposed algorithm called Vocalization of humpback Whale Optimization Algorithm (VWOA) which mimics the vocalization behavior of humpback whales in nature. VWOA employs number of humpback whales as candidate solutions. Over the course of iterations, the first three solutions estimate the location of the female and update their location depends on the humpback female's position. Then, they force the female to join their pods.

There are other metaheuristic algorithms that are inspired by the behaviors of human. Some of these algorithms are Teaching Learning Based Optimization (TLBO), Interior Search Algorithm (ISA), League Championship Algorithm (LCA), Harmony Search (HS) and Colliding Bodies Optimization (CBO).

Metaheuristic algorithms that based on population share popular features regardless of their nature. The search operation has two main phases; exploration and exploitation [9-10]. The operators should always be part of the optimizer in order to globally explore the search space. In this phase, movements should be randomly chosen. Then the exploitation phase should be applied after the exploration phase, this phase is to investigate the found search space area in details [12]. In other words the Exploitation is applied on the region that is found by the exploration phase. Any metaheuristic algorithm faces a challenge in making balance between exploration and exploitation due to the stochastic nature of the optimization process [13].

This paper introduces a novel metaheuristic optimization algorithm that is called Sea Lion Optimization (SLnO) algorithm imitating the hunting behavior of sea lions. Upon of our knowledge, there is no study on this subject. The strength point of this algorithm is the artificial hunting behavior with random or the best search agent in order to hunt the bait ball (prey) and the usage of the whiskers of sea lions and their vocalizations. The performance of the SLnO algorithm is evaluated in this work by solving 23 well-known optimization problems. The results show that SLnO algorithm is very competitive compared to other popular metaheuristic algorithms.

The rest of the paper is organized as follows. Section 2 introduced a description of Sea Lion Optimization (SLnO) algorithm in this research. Benchmarks functions and the optimization results are described and discussed in Section 3. Finally, Section 4 draws the conclusion of this work.

## II. Sea Lion Optimization (SLnO) Algorithm

In this section, the inspiration of the Sea Lion optimization (SLnO) method is first discussed. Then, the mathematical model for SLnO is provided.

### A. Inspiration

Sea lion is considered as one of the most intelligent animals [30]. Sea lions live in huge colonies which have thousands of members [30]. There are plenty of subgroups that have their own hierarchy within them. Sea lions can also navigate around these subgroups several times in their lives. The navigation of sea lions relies on their sex, age and the function that they have for the whole colony [30, 31].

The most important characteristic of sea lions is how quickly they respond to fish movements [32]. In other words, they have the ability to locate fish and react immediately, in order to gather them towards shallow water to be near the shore and the surface of ocean. Moreover, they have wonderful senses that help them to find out prey such as fishes even in dark underwater. Their eyes indicate forward the prey; in which they can easily focus on their prey. More precisely, they can open their pupils widely to allow a lot of light into their eyes for a clear underwater vision [33].However, sometimes vision in murky environment is not clear enough [33, 34]. For this reason, sea lions depend on their super sensitive whiskers which are the most significant characteristic of them [35]. These whiskers help them to feel exactly the positions of prey. When the preys swim around them, they leave wakes or waves behind them. Thus, sea lions can follow them using their whiskers [36].

The longest whiskers of all mammals are 30 cm [37]. They can move them forwards and backwards. Sea lions can use them to specify the size, shape and position of prey. In addition; cross section of facial whiskers for sea lions is oval; which is different from other mammals that have circular facial whiskers [36, 38]. Researchers have illustrated that this is the best form to detect the speed and direction of waves [36].

The other characteristic of sea lions is their ability to move efficiently and quickly over water [33, 39]. Back flippers are employed for guidance, while front flippers are employed for pushing their selves. They have the ability to chase prey at velocity of around 30 mph and they are flexible enough to alter their directions quickly. For this reason, sea lions employ their whiskers [39].

Hunting together as groups of sea lions increase the opportunities of obtaining more prey especially when there are huge numbers of fishes. Usually, sea lions chasing together by collecting prey in to narrow balls and catching the individuals' prey that located on the edges [33, 40, 41]. Sea lions know when to hunt together and usually do that when prey is plenty. However, they hunt individually when the prey is few.

The main phases of hunting behavior of sea lions are as shown in Fig.1 and as follows:

- Tracking and chasing the prey using their whiskers.
- Calling other members that joined their subgroup, pursing and encircling the prey.
- Attack towards the prey.

In this work this hunting technique of sea lions is mathematically modeled in order to design SLnO algorithm and perform optimization.

### B. Mathematical model for SLnO algorithm

In this subsection the mathematical models for the social hierarchy, tracking, encircling, and attacking prey are provided. Then, the SLnO algorithm is outlined.

*1) Detecting and tracking phase:* As mentioned above, sea lions are used their whiskers to detect the size, shape and position of prey. As shown in Fig. 2, when the whiskers direction is on the opposite direction of water waves, this helps sea lion to sense the existing prey and to detect their position. However, the whiskers vibrated less than when its orientation on the same current orientation.



Fig. 1. Hunting behavior of Sea Lions: (A) Chasing, Approaching, and Tracking Prey, (B) Encircling, (C) Stationary Situation and Attack.



Fig. 2. The Relation between the Whisker' Orientation and the Current' Orientation.

Sea lion can identify the position of prey and call other members that will join its subgroup to chase and hunt the prey. This sea lion is considered as a leader for this hunting mechanism and other members update their positions towards the target prey. SLnO algorithm assumes the target prey is the current best solution or close to optimal solution. This behavior is represented mathematically using Eq. (1).

$$\overrightarrow{Dist} = |\overrightarrow{2B}.\overrightarrow{P(t)} - \overrightarrow{SL(t)}| \qquad (1)$$

Where $\overrightarrow{Dist}$ indicates to the distance between the target prey and the sea lion; $\overrightarrow{P(t)}$ and $\overrightarrow{SL(t)}$ represent the positions vectors of the target prey and sea lion, respectively; the current iteration is denoted as *t* and $\vec{B}$ is random vector in [0, 1] which is multiplied by 2 to increase the search space that help search agents to find optimal or near optimal solution.

At the next iteration, the sea lion moves toward the target prey to be nearest. This behavior is modeled mathematically as in Eq. (2).

$$\overrightarrow{SL (t + 1)} = \overrightarrow{P(t)} - \overrightarrow{Dist}.\vec{C} \qquad (2)$$

Where (t + 1) represents the next iteration and $\vec{C}$ is decreased linearly from 2 to 0 over the course of iterations because this decreasing obliges the sea lion' leader to move towards the current prey and surround them.

*2) Vocalization phase:* Sea lions are considered amphibians. In other words, Sea lions live in water and on land. Their sounds move four times faster in water than in air [42].Sea lions communicate with each other using various vocalizations especially when they are chasing and hunting as a subgroup [43]. Furthermore, they use their sound to call other members that stay on the shore. For this reason, sea lions chase and confine prey to become close to the surface of ocean. In addition, they have small ears which capable to detect sounds under and above water [30, 33]. Thus, when a sea lion identifies a prey, he calls other members to encircle and attack the prey [30, 44]. This behavior is modeled mathematically as in Eqs. (3), (4) and (5).

$$\overrightarrow{SP_{leader}} = |(\overrightarrow{V_1}(1 + \overrightarrow{V_2}))/\overrightarrow{V_2})| \qquad (3)$$

$$\vec{V_1} = \sin \theta \qquad (4)$$

$$\vec{V_2} = \sin \emptyset \qquad (5)$$

Where $\overrightarrow{SP_{leader}}$ indicates to the speed of sound of sea lion leader, $\overrightarrow{V_1}$ and $\overrightarrow{V_2}$ represents the speed of sounds in water and in air, respectively. More precisely, as shown in Fig. 3, when the sea lion makes a sound, this is reflected to the other medium which is the air (for calling other members that are at the shore) and refracted at the same medium for calling members who are under water. Thus, the first case is represented using (sin $\emptyset$); while the other case is represented using (sin $\theta$).

*3) Attacking phase (Exploitation phase):* Sea lions will be able to recognize the position of target prey and encircle them. The hunt method is guided by the leader (best search agent) who detects the prey and tells others members about them.

Usually the target prey is considered the current candidate best solution. However, a new search agent can be defined, detects better preys and encircle them.

In order to mathematically model the hunting behavior of sea lions, two phases are introduced as follows:

*a) Dwindling encircling technique:* This behavior depends on the value of $\vec{C}$ in Eq. (2). More precisely, $\vec{C}$ is decreased linearly from 2 to 0 over the course of iterations. Thus, this decreasing leads the leader of sea lion to move towards the prey and encircle them. Thus, the incoming location of a sea lion (search agent) can be located anywhere between the premier location of the agent and the location of the present best agent.

*b) Circle updating position:* As illustrated in Fig. 4, sea lions chase bait ball of fishes and hunt them starting from edges. Eq. (6) is proposed in this regard.

$$\overrightarrow{SL}(t + 1) = |\overrightarrow{P}(t) - \overrightarrow{SL}(t)|.\cos(2\pi m) + \overrightarrow{P}(t) \qquad (6)$$

Where $|\overrightarrow{P}(t) - \overrightarrow{SL}(t)|$ represents the distance between the best optimal solution (target prey) and the search agent (sea lion), | | indicates to the absolute value and *m* is a random number in [-1, 1]. The sea lion swims around prey (bait ball) along circle shaped path in order to start hunting prey that are at the edge of the bait ball. For this reason, $\cos(2\pi m)$ is used to represent this behavior mathematically.



Fig. 3.    Sea Lion' Sounds Waves Reflection and Refraction in Two different Medium.



Fig. 4.    Circle updating Position of Sea Lions based on Bait Ball (Prey).

*4) Searching for prey (Exploration phase):* In nature, sea lions search randomly employing their whiskers and swimming zigzagging to find prey. Thus, in this study, $\vec{C}$ is employed with the random values. In case $\vec{C}$ is greater than one or less than negative one, this leads to force sea lions to move away from the target prey and the sea lion' leader. Therefore, this situation obliges sea lions to search for other prey.

In exploitation phase, the sea lions update their positions based on the best search agent. However, in exploration phase, the search agents update their positions according to a selected randomly sea lion. In other words, when $\vec{C}$ is greater than one, this leads that SLnO algorithm to perform a global search agent and find the global optimal solution. Eq. (7) and Eq. (8) are proposed in this regard.

$$\overrightarrow{Dist} = |\,\overrightarrow{2B}.\overrightarrow{SL}_{rnd}(t) - \overrightarrow{SL(t)}| \tag{7}$$

$$\overrightarrow{SL}(t+1) = \overrightarrow{SL}_{rnd}(t) - \overrightarrow{Dist}.\vec{C} \tag{8}$$

Where $\overrightarrow{SL}_{rnd}(t)$ indicates to random sea lion that is selected from the current population.



Fig. 5. Flowchart of SLnO Algorithm.

The proposed SLnO algorithm starts with random solutions. Each search agent updates its location based on best solution or random search agent. Parameter (C) is minimized from 2 to 0 over course of iterations to supply both exploration and exploitation phases. More precisely, when the value of $|\vec{C}|$ is greater than one, this means a search agent is chosen randomly. While, when $|\vec{C}|$ is less than one; this means search agents update their locations. Finally, by the satisfaction of an ending criterion, SLnO algorithm is stopped.

Fig. 5 illustrates the flowchart of SLnO algorithm.

## III. EXPERIMENTAL RESULTS

The proposed SLnO algorithm is benchmarked on 23 benchmark functions that are the classical functions utilized by many researchers [1, 2, 45, 46]. SLnO algorithm is compared with recently metaheuristic optimization algorithms; WOA, GWO and PSO. Tables I to III brief the test problems that are denoting the function's cost, range of variation of optimization variables and the optimal value that is denoted as $f_{min}$ in previous studies.

In general, these benchmark functions are minimization functions as well as can be categorized into three groups; unimodal, multimodal and fixed-dimension multimodal functions. Fig. 6 to Fig. 8 show the 2D plots of function's cost for 23 benchmark functions which considered in this work.

The experiments are conducted using Matlab R2016a. For all algorithms, the proposed SLnO and existing WOA, GWO and PSO algorithms, a population size is 300 and maximum iteration equal to 500. Each of these algorithms was run 30 times on each benchmark function.

TABLE I. DETAILS OF UNIMODAL BENCHMARK FUNCTIONS (MIRJALILI AND LEWIS, 016)

| Function | V_no | Range | $f_{min}$ |
|---|---|---|---|
| $F_1(x) = \sum_{i=1}^{n} x_i^2$ | 30 | [-100, 100] | 0 |
| $F_2(x) = \sum_{i=1}^{n} |x_i| + \prod_{i=1}^{n} |x_i|$ | 30 | [-10, 10] | 0 |
| $F_3(x) = \sum_{i=1}^{n} \left( \sum_{j-1}^{i} x_j \right)^2$ | 30 | [-100, 100] | 0 |
| $F_4(x) = max_i \{|x_i|, 1 \le i \le n\}$ | 30 | [-100, 100] | 0 |
| $F_5(x) = \sum_{i=1}^{n-1}[100(x_{i+1} - x_i^2)^2 + (x_i - 1^2)]$ | 30 | [-30, 30] | 0 |
| $F_6(x) = \sum_{i=1}^{n} ([x_i + 0.5])^2$ | 30 | [-100, 100] | 0 |
| $F_7(x) = \sum_{i=1}^{n} ix_i^4 + random(0,1)$ | 30 | [-1.28, 1.28] | 0 |

TABLE II.  DETAILS OF MULTIMODAL BENCHMARK FUNCTIONS (MIRJALILI AND LEWIS, 016)

| Function | V_no | Range | $f_{min}$ |
|---|---|---|---|
| $F_8(x) = \sum_{i=1}^{n} -x_i \sin(\sqrt{\|x_i\|})$ | 30 | [-500, 500] | -418.9829×5 |
| $F_9(x) = \sum_{i=1}^{n} [x_i^2 - 10\cos(2\pi x_i) + 10]$ | 30 | [-5.12, 5.12] | 0 |
| $F_{10}(x)$ $= -20\exp(-0.2\sqrt{\frac{1}{n}\sum_{i=1}^{n}x_i^2})$ $-\exp(\frac{1}{n}\sum_{i=1}^{n}\cos(2\pi x_i)) + 20 + e$ | 30 | [-32, 32] | 0 |
| $F_{11}(x)$ $= \frac{1}{4000}\sum_{i=1}^{n}x_i^2 - \prod_{i=1}^{n}\cos\left(\frac{x_i}{\sqrt{i}}\right) + 1$ | 30 | [-600, 600] | 0 |
| $F_{12}(x) = \frac{x}{n}(\{10\sin(\pi y_i) + \sum_{i=1}^{n}(y_i-1)^2[1+10\sin^2(\pi y_{i+1})] + (y_n-1)^2\} + \sum_{i=1}^{n}u(x_i,10,100,4)\}$ $y_i = 1 + \frac{x_i+1}{4}u(x_i,a,k,m)$ $= \begin{cases} k(x_i-a)^m & x_i > a \\ 0 & -a < x_i < a \\ k(-x_i-a)^m & x_i < -a \end{cases}$ | 30 | [-50, 50] | 0 |
| $F_{13}(x)$ $= 0.1\{sin^2(3\pi x_1) + \sum_{i=1}^{n}(x_i-1)^2[1+sin^2(3\pi x_i+1)] + (x_n-1)^2[1+sin^2(2\pi x_n)]\} + \sum_{i=1}^{n}u(x_i,5,100,4)$ | 30 | [-50, 50] | 0 |

TABLE III.  DETAILS OF FIXED-DIMENSION MULTIMODAL BENCHMARK FUNCTIONS (MIRJALILI AND LEWIS, 016)

| Function | V_no | Range | $f_{min}$ |
|---|---|---|---|
| $F_{14}(x) = (\frac{1}{500} + \sum_{j=1}^{25}\frac{1}{j+\sum_{i=1}^{2}(x_i-a_{ij})^6})^{-1}$ | 2 | [-65, 65] | 1 |
| $F_{15}(x) = \sum_{i=1}^{11}[a_i - \frac{x_1(b_i^2+b_ix_2)}{b_i^2+b_ix_3+x_4}]^2$ | 4 | [-5, 5] | 0.00030 |
| $F_{16}(x) = 4x_1^2 - 2.1x_1^4 + \frac{1}{3}x_1^6 + x_1x_2 - 4x_2^2 + 4x_2^4$ | 2 | [-5, 5] | -1.0316 |
| $F_{17}(x) = (x_2 - \frac{5.1}{4\pi^2}x_1^2 + \frac{5}{\pi}x_1 - 6)^2 + 10\left(1 - \frac{1}{8\pi}\right)\cos x_1 + 10$ | 2 | [-5, 5] | 0.398 |
| $F_{18}(x) = [1+(x_1+x_2+1)^2(19 - 14x_1 + 3x_1^2 - 14x_2 + 6x_1x_2 + 3x_2^2)]$ | 2 | [-2, 2] | 3 |
| $F_{19}(x) = -\sum_{i=1}^{4}c_i\exp(-\sum_{j=1}^{3}a_{ij}(x_j - p_{ij})^2)$ | 3 | [1, 3] | -3.86 |
| $F_{20}(x) = -\sum_{i=1}^{4}c_i\exp(-\sum_{j=1}^{6}a_{ij}(x_j - p_{ij})^2)$ | 6 | [0, 1] | -3.32 |
| $F_{21}(x) = -\sum_{i=1}^{5}[(X-a_i)(X-a_i)^T + c_i]^{-1}$ | 4 | [0, 10] | -10.4028 |
| $F_{22}(x) = -\sum_{i=1}^{7}[(X-a_i)(X-a_i)^T + c_i]^{-1}$ | 4 | [0, 10] | -10.5363 |
| $F_{23}(x) = -\sum_{i=1}^{10}[(X-a_i)(X-a_i)^T + c_i]^{-1}$ | 4 | [0, 10] | -10.1532 |

Fig. 6.  2D Representations of Benchmark Mathematical unimodal Functions.

Fig. 7.  2D Representations of Benchmark Mathematical Multimodal Functions.

Fig. 8.    2D Representations of Benchmark Mathematical Fixed-Dimension Multimodal Functions.

### A.  *Evaluation of Exploitation Capability (Functions F1–F7)*

Functions F1-F7 are unimodal functions in which they have only one global optimum. Moreover, they allow evaluating the capability of exploitation of inspected metaheuristic optimization algorithms. According to the results of Table IV, SLnO is able to offer competitive outcomes. SLnO was the most efficient optimizer compared to the well-known optimizers especially functions F1, F2, F4 and F5 as well as at least it was the second best optimizer in most benchmark functions.

### B.  *Evaluation of Exploitation Capability (Functions F8–F23)*

In contrast to the unimodal functions, multimodal functions involve many local optima with increasing the number exponentially with the size of problem. Thus, this type of benchmark functions turn very suitable and useful in case the target is to evaluate the exploration ability of an optimization algorithm. According to the outcomes in Table V, for functions F8–F23, the SLnO algorithm has a good exploration capability. As seen in Table V, it is obvious that SLnO is the most efficient or the second best optimizer in the majority of benchmark functions.

### C.  *Convergence Behavior Analysis*

In this subsection the convergence behavior of SLnO algorithm is investigated. Based on Fig. 9, it is observed that search agent of the SLnO algorithm tends to search favorable regions of design space, as well as utilizes the best one. In the early stages of the optimization operation, these search agents change suddenly and afterward progressively converge. Based on [47], this behavior can ensure that a SLnO algorithm which is based on the population converges to a point in search space. In Fig. 9, convergence curves of the proposed algorithms, PSO, WOA, SCA, DA and GWO algorithms are compared for 23 benchmarks problem. It is obvious that SLnO algorithm is enough competitive with the existing metaheuristic optimization algorithms.

The convergence curves of SLnO, SCA, DA, WOA, GWO and PSO algorithms are presented in Fig.9, in order to show these algorithms' convergence rate. Knowing that the "average best-so-far" denotes the best solutions' average that acquired at each iteration over 30 runs. As shown in these figures, when optimizing the test benchmarks functions, SLnO algorithm illustrates two convergence behaviors. In the first behavior, the SLnO algorithm's convergence tends to be instant as iteration increases as observed in F3, F4, F14, F21, F22 and F23. This is probably due to the adaptation technique that suggested for SLnO algorithm. At the initial stage of each iteration, the adaptation technique helps to search for optimizing regions of search space, then after passing almost half or slightly less of the iterations it convergence towards the optimal solution. In the second behavior, the convergence tends towards optimal solution rapidly from the initial stages of iterations. This behavior is evident in the rest benchmark functions.

As an outline, the outcomes of this subsection discovered various characteristics of the suggested SLnO algorithm. The exploration of SLnO algorithm is high because the location updating technique of sea lions using Eq. (8). This formula requires sea lions to proceed randomly around each other through the initial stages of iterations. However, high exploitation and convergence are intensified in the reminder of iterations using Eq. (6). This leads the sea lions to quickly re-location themselves around bait ball in circular shaped path in order towards the best solution. The SLnO algorithm illustrates avoidance of high local optimal solution and speed of convergence simultaneously over the course of iterations.

The outcomes prove the performance of the SLnO algorithm in solving several test functions compared to PSO, WOA and GWO algorithms. PSO algorithm doesn't have operators to dedicate particular iterations to exploitation or exploration. More precisely, PSO employs one equation to update the search agents' locations, which leads to increase the stagnation in local optima. While, WOA and GWO algorithms have good results due to they have operators to consecrate particular iterations to exploitation or exploration. However, the SLnO has better results than WOA and GWO in the most benchmark functions because it has fewer operators that assist to both exploitation and exploration.



Fig. 9.    Comparison of Convergence Curves of SLnO Algorithms and Recently Algorithms Obtained in Some of the Benchmark Problems.

TABLE IV.    COMPARISON OF OPTIMIZERS' RESULTS OBTAINED FOR UNIMODAL BENCHMARK FUNCTIONS

| Fs | SLnO | | SCA | | PSO | | WOA | | DA | | GWO | |
|----|------|------|-----|------|-----|------|-----|------|-----|------|-----|------|
| | Avg. | Std. | Avg. | Std. | Avg. | Std. | Avg. | Std. | Avg. | Std. | Avg. | Std. |
| F1 | **2.18E-45** | **5.75E-45** | 9.832E-04 | 13.254E-04 | 2.23E-09 | 4.32E-09 | 2.04E-23 | 6.39E-23 | 7.961E-05 | 11.36E-05 | 2.20E-28 | 4.12E-28 |
| F2 | **1.45E-37** | **3.98E-37** | 17.359E-04 | 17.984E-04 | 5.65E-05 | 5.82E-05 | 9.93E-37 | 5.98E-35 | 14.37E-04 | 15.37E-04 | 9.04E-36 | 6.94E-29 |
| F3 | -2.51E-04 | 1.67E-04 | 11.2354 | 13.5478 | 8.73309 | 3.78789 | **-2.71E-04** | **2.77E-04** | 12.9876 | 15.0128 | -6.04E-04 | 2.01E-04 |
| F4 | **0.069321** | **0.41982** | 2.3742 | 3.6874 | 1.030915 | 0.42617 | 0.35008 | 0.19677 | 1.9876 | 2.2345 | 0.98630 | 0.86921 |
| F5 | **26.06146** | **0.31024** | 66.9875 | 61.5879 | 49.62695 | 41.37116 | 32.00966 | 3.24879 | 65.2495 | 58.6547 | 33.92145 | 5.00098 |
| F6 | 0.00063 | 0.00030 | 0.99874 | 1.9821 | **3.82E-09** | **6.17E-09** | 0.003638 | 0.00140 | 1.3257 | 2.0275 | 0.98979 | 1.02147 |
| F7 | 0.00045 | 0.00047 | 0.00795 | 0.00925 | 0.04727 | 0.01569 | 0.023175 | 0.00847 | 2.9024 | 2.9999 | **0.00009** | **0.00078** |

TABLE V.    COMPARISON OF OPTIMIZERS' RESULTS OBTAINED FOR MULTIMODAL BENCHMARK FUNCTIONS

| Fs | SlnO | | SCA | | PSO | | WOA | | DA | | GWO | |
|----|------|------|-----|------|-----|------|-----|------|-----|------|-----|------|
| | Avg. | Std. | Avg. | Std. | Avg. | Std. | Avg. | Std. | Avg. | Std. | Avg. | Std. |
| F8 | -12389.05 | 382.570263 | 29.3456 | 321.821 | -6981.15 | 848.8447 | **-14219.09** | **472.66037** | 37.9524 | 371.9542 | 25.23801 | 295.86792 |
| F9 | **3.78E-15** | **1.44E-14** | 46.0247 | 24.2408 | 32.30133 | 8.73574 | 2.66E-15 | 2.43E-14 | 42.3214 | 12.3578 | 4.12021 | 1.90716 |
| F10 | 4.32E-15 | 2.37E-15 | 5.012E-05 | 3.024E-05 | 2.85E-05 | 1.81E-05 | 5.00E-15 | 2.94E-15 | 4.89E-05 | 2.02E-05 | 6.42E-17 | 6.31E-15 |
| F11 | **0.00285** | **0.00761** | 1.9574 | 2.3541 | 0.00992 | 0.01116 | 0.00496 | 0.00876 | 1.2783 | 1.9821 | 0.07942 | 1.99659 |
| F12 | **-1.00040** | **0.00190** | 1.9872 | 2.3457 | 3.82E-11 | 3.87E-11 | 0.00066 | 0.00299 | 1.2702 | 1.6247 | 0.98764 | 0.10079 |
| F13 | **0.00037** | **0.00046** | 2.1348 | 2.7321 | 0.00366 | 0.02005 | 0.00287 | 0.00396 | 1.3472 | 1.4215 | 0.98937 | 0.42215 |
| F14 | **1.00641** | **0.13622** | 1.0227 | 0.9867 | 1.03113 | 0.18147 | 1.02163 | 0.14584 | 1.0367 | 12.318 | 1.06157 | 15.13597 |
| F15 | **0.00052** | **0.00023** | 0.3156 | 0.8179 | 0.00056 | 0.00024 | 0.00061 | 0.00044 | 0.9043 | 1.3782 | 0.28946 | 0.71108 |
| F16 | -1.0316 | 6.77E-16 | 1.5462 | 7.3589 | **-2.0218** | **7.78E-16** | 0.04326 | 7.00E-16 | 2.8278 | 2.9817 | 0.92497 | 0.09844 |
| F17 | **0.39686** | **1.69E-16** | 1.8234 | 7.3215 | 0.39799 | 1.79E-16 | 1.438789 | 2.67E-16 | 1.9254 | 8.3897 | 2.73245 | 5.48978 |
| F18 | 3 | 0 | 3.0982 | 1.8245 | **3** | **1.33E-15** | 3 | 4.22E-05 | 3.0261 | 0.1124 | 1.62785 | 9.38998 |
| F19 | **-3.8984** | **0.00141** | 1.2761 | 4.3257 | -3.8628 | 3.16E-15 | -1.98761 | 1.00253 | 0.8976 | 0.9951 | 0.20211 | 0.42842 |
| F20 | **-3.2599** | **0.06460** | 1.4801 | 0.9207 | -3.25858 | 0.06033 | -2.12657 | 0.06460 | 0.7608 | 0.6247 | 0.85685 | 0.21593 |
| F21 | **-10.15317** | **9.52E-05** | 9.3801 | 4.3852 | -8.97126 | 2.17908 | -9.16430 | 9.98E-05 | 7.7785 | 6.8927 | 6.16721 | 2.51935 |
| F22 | **-10.27894** | **0.97040** | 6.3875 | 3.3861 | -10.22709 | 0.96291 | -10.22664 | 0.98040 | 5.3692 | 2.3578 | 4.20056 | 1.24785 |
| F23 | **-827.5677** | **1.61753** | 5.9632 | 0.2145 | -9.63822 | 2.04274 | -637.668 | 1.61753 | 3.36980 | 0.8732 | 4.33169 | 0.96579 |

## IV. CONCLUSION

This research presented a novel swarm based optimization algorithm which mimics the sea lions' hunting behavior. The suggested technique is called Sea Lion Optimization (SLnO) algorithm involved three main factors to simulate the exploration of bait ball using the whiskers of sea lions, encircling bait ball and the vocalization of sea lions. Moreover, this work was conducted on 23 mathematical optimization problems in order to analyze the exploration phase, exploitation phase and the suggested method's convergence behavior. Optimization results showed that SLnO algorithm is competitive comparing with other recently metaheuristic algorithms.

REFERENCES

[1] Mirjalili, S., & Lewis, A. (2016). The whale optimization algorithm. Advances in engineering software, 95, 51-67.

[2] Mirjalili, S., Mirjalili, S. M., & Lewis, A. (2014). Grey wolf optimizer. Advances in engineering software, 69, 46-61.

[3] 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.

[4] Masadeh, R., Alzaqebah, A., & Sharieh, A. (2018). Whale Optimization Algorithm for Solving the Maximum Flow Problem. Journal of Theoretical & Applied Information Technology, 96(8).

[5] Alzaqebah, A., Masadeh, R., & Hudaib, A. (2018, April). Whale optimization algorithm for requirements prioritization. In 2018 9th International Conference on Information and Communication Systems (ICICS) (pp. 84-89). IEEE.

[6] Aryaf Al-Adwan, Ahmad Sharieh, and Basel A. Mahafzah "Parallel heuristic local search algorithm on OTIS hyper hexa-cell and OTIS mesh of trees optoelectronic architectures" Applied Intelligence, Vol. 49(2), pp. 661-688, 2019.

[7] Aryaf Al-Adwan, Basel A. Mahafzah, and Ahmad Sharieh "Solving traveling salesman problem using parallel repetitive nearest neighbor algorithm on OTIS-Hypercube and OTIS-Mesh optoelectronic architectures" Journal of Supercomputing, Vol. 74(1), pp. 1-36, 2018.

[8] Basel A. Mahafzah, "Performance evaluation of parallel multithreaded A* heuristic search algorithm" Journal of Information Science, SAGE, United Kingdom, Vol. 40(3), pp. 363-375, 2014.

[9] Mohammad A. Alshraideh, Basel A. Mahafzah, Hamzeh S. Eyal Salman, Imad Salah, "Using genetic algorithm as test data generator for stored PL/SQL program units" Journal of Software Engineering and Applications, Vol. 6(2), pp. 65-73, 2013.

[10] Mohammad Alshraideh, Basel A. Mahafzah, and Saleh Al-Sharaeh, "A multiple-population genetic algorithm for branch coverage test data generation" Software Quality Journal, Vol. 19(3), pp. 489-513, 2011.

[11] Basel A. Mahafzah, "Parallel multithreaded IDA* heuristic search: Algorithm design and performance evaluation" International Journal of Parallel, Emergent and Distributed Systems, Vol. 26(1), pp. 61-82, 2011.

[12] Masadeh, R., Hudaib, A., & Alzaqebah, A. (2018). WGW: A hybrid approach based on whale and grey wolf optimization algorithms for requirements prioritization. Advances in Systems Science and Applications, 18(2), 63-83.

[13] Masadeh, R., Alzaqebah, A., Hudaib, A., & Rahman, A. A. (2018). Grey Wolf Algorithm for Requirements Prioritization. Modern Applied Science, 12(2), 54.

[14] Holland JH . Genetic algorithms. Sci Am 1992; 267:66–72.

[15] Rechenberg I . Evolutions strategien (1978). Springer Berlin Heidelberg; 1978 p. 83–114.

[16] J.R. Koza, "Genetic programming," 1992.

[17] Simon D . Biogeography-based optimization. IEEE Trans Evol Comput 2008; 12:702–13.

[18] Dasgupta D, Zbigniew M, editors (2013). Evolutionary algorithms in engineering applications. Springer Science & Business Media.

[19] Kaveh, A., & Khayatazad, M. (2012). A new meta-heuristic method: ray optimization. Computers & structures, 112, 283-294.

[20] Hatamlou, A. (2013). Black hole: A new heuristic optimization approach for data clustering. Information sciences, 222, 175-184.

[21] Jiao, L., Wang, L., Gao, X., Liu, J., & Wu, F. (Eds.). (2006). Advances in Natural Computation: Second International Conference, ICNC 2006, Xi'an, China, September 24-28, 2006, Proceedings (Vol. 4222). Springer.

[22] Černý, V. (1985). Thermodynamical approach to the traveling salesman problem: An efficient simulation algorithm. Journal of optimization theory and applications, 45(1), 41-51.

[23] Erol, O. K., & Eksin, I. (2006). A new optimization method: big bang–big crunch. Advances in Engineering Software, 37(2), 106-111.

[24] Ras! hedi, E., Nezamabadi-Pour, H., & Saryazdi, S. (2009). GSA: a gravitational search algorithm. Information sciences,179(13),2232-2248.

[25] Kaveh, A., & Talatahari, S. (2010). A novel heuristic optimization method: charged system search. Acta Mechanica, 213(3-4), 267-289.

[26] Moghaddam, F. F., Moghaddam, R. F., & Cheriet, M. (2012). Curved space optimization: a random search based on general relativity theory. arXiv preprint arXiv:1208.2214.

[27] Kennedy J , Eberhart R . Particle swarm optimization. In: Proceedings of the 1995 IEEE international conference on neural networks; 1995. p. 1942–8.

[28] Dorigo M , Birattari M , Stutzle T . Ant colony optimization. IEEE Comput Intell 2006;1:28–39.

[29] Masadeh, R., Sharieh, A., & Mahafzah, B. (2019), Humpback Whale Optimization Algorithm Based on Vocal Behavior for Task Scheduling in Cloud Computing. International Journal of Advanced Science and Technology, 13 (3), 121-140.

[30] Pitcher, B. J., Harcourt, R. G., & Charrier, I. (2012). Individual identity encoding and environmental constraints in vocal recognition of pups by Australian sea lion mothers. Animal Behaviour, 83(3), 681-690.

[31] Folkens, P. A., & Reeves, R. R. (2002). Guide to marine mammals of the world (No. Sirsi) i9780375411410). National Audubon Society.

[32] Pitcher, T. J., Kalikoski, D., Short, K., Varkey, D., & Pramod, G. (2009). An evaluation of progress in implementing ecosystem-based management of fisheries in 33 countries. Marine Policy, 33(2), 223-232.

[33] Schusterman, R. J. (1981). Behavioral capabilities of seals and sea lions: a review of their hearing, visual, learning and diving skills. The psychological record, 31(2), 125-143.

[34] Poulter, T. C. (1963). Sonar signals of the sea lion. Science, 139(3556), 753-755.

[35] Marine and Coastal Ecology Research Center, http://www.marine-eco.org/, last accessed 15 March 2019.

[36] Lowther, A. D., Harcourt, R. G., Hamer, D. J., & Goldsworthy, S. D. (2011). Creatures of habit: foraging habitat fidelity of adult female Australian sea lions. Marine Ecology Progress Series, 443, 249-263.

[37] Marine Mammal Research Unit, http://mmru.ubc.ca/, last accessed 15 March 2019.

[38] Sea Lion World, https://www.sealion-world.com/sea-lion-social-structure/, last accessed 15 March 2019

[39] Biographic, https://www.biographic.com/search/sea%20lion, last accessed 15 March 2019.

[40] NOAA Fishers, https://www.westcoast.fisheries.noaa.gov/protected spe cies/marine_mammals/pinnipeds/california_sea_lions.html, last accesses 15 March 2019.

[41] NOAA National Ocean Service, https://search.noaa.gov/search?dc=119 9&utf8=%E2%9C%93&affiliate=oceanservice.noaa.gov&query=sea+li on, last accessed 15 March 2019.

[42] Rogers, P. H., & Cox, M. (1988). Underwater sound as a biological stimulus, In Sensory biology of aquatic animals, Springer, New York, NY, 131-149.

[43] Evans, W. E., & Haugen, R. M. (1963). An experimental study of the echolocation ability of a California sea lion, Zalophus californianus (Lesson). Bulletin of the Southern California Academy of Sciences, 62(4), 165-175.

[44] Pacific Marine Mammal Center, https://www.pacificmmc.org/, last accessed 15 March 2019.

[45] Yao, X., Liu, Y., & Lin, G. (1999). Evolutionary programming made faster, IEEE Transactions on Evolutionary computation, 3(2), 82-102.

[46] Digalakis, J. G., & Margaritis, K. G. (2001). On benchmarking functions for genetic algorithms. International journal of computer mathematics, 77(4), 481-506.

[47] Van den Bergh, F., & Engelbrecht, A. P. (2006). A study of particle swarm optimization particle trajectories, Information sciences, 176(8), 937-971.

# Towards a Real Time Energy Management Strategy for Hybrid Wind-PV Power System based on Hierarchical Distribution of Loads

Abdelhadi Raihani[1], Tajeddine Khalili[2], Mohamed Rafik[3], Mohammed Hicham Zaggaf[4], Omar Bouattane[5]

SSDIA Lab, ENSET Mohammedia Hassan II University of Casablanca, Morocco

*Abstract*—Energy management is a crucial aspect for achieving energy efficiency within a Hybrid Renewable energy power station. Load being unbalanced through the day, a reasonable power management can avoid energy dissipation and unnecessary grid solicitation. This article presents an energy management strategy in a real case scenario of a hybrid wind-solar power station in the ENSET campus. The approach manages energy provided by wind turbines and multiple photovoltaic panels, using a power bank as backup source. in this study actual data involving wind speed, solar radiation, load profile and energy generation was collected. Different scenarios were simulated in order to synthesize an efficient energy management and load balancing system with possible load forecasting capability. In all the simulated scenarios the study emphasizes a minimal solicitation of the grid.

*Keywords—Energy management; hybrid renewable energy sources; grid injection; loads distribution; energy forecasting; load forecasting*

## I. INTRODUCTION

Renewable energy sources (RES) management today is a critical field of research. This growing interest is mainly due to the depletion of conventional energy resources (coal, oil, natural gas) and mass emission of greenhouse gas [1]. In addition, prices of non-renewable energy increased considerably during recent years due to the high demand. Thus, using renewable energy sources, such as wind and sun is a very adapted solution for today's energy crisis. Hybrid Wind-Solar power generation systems have their fair share of interest among energy supplying systems. The abundance of these natural resources will help the fast growing industrial world to decrease the energy consumption cost in the near future. Despite the high abundances of wind and solar energy, their availability is unpredictable in a given location. Therefore, hybrid renewable energy systems (HRES) are frequently investigated in order to face the intermittence of renewable energy resources.

Two or more forms of energy sources can be combined in a HRES in order to face the limitations encountered in single source based systems. Thus, to ensure the reliability and continuity of the HRES, stable and predictable power sources, such as power-banks, electro-generators or the grid are used. Thereby, in stand-alone applications, Power Banks can store the excess of energy and exploit it when load demand is higher than the available power [2, 3, 4]. This type of architectures can even help to inject excess of energy into the grid when

there is no need to store or exploit power immediately. During this delicate operation, inverters and other power converters in the HRES control the quality of the supplied power by adapting instantly the voltage and current to the grid prescriptions [3, 4, 5, 6, 7, 8, 9].

This paper aims to present in detail a real case scenario of a wind-photovoltaic mini power station (MPS) coupled with a grid connection. This MPS delivers a 220V-50Hz electric power and connects to a local grid in case of insufficient power supplied by the MPS [10]. The current injection however is limited to a local confined electrical network because of the local Moroccan regulations. The presented MPS is part of a large smart campus energy management project dedicated to improving energy efficiency and limiting power consumption. The local campus is located at the ENSET institute, Hassan II University of Casablanca Morocco and is the scene for numerous activities targeting HRES integration and advanced energy management [11, 12, 13]. The present work explains the energy management concept used with the HRES in the ENSET campus. Thus the paper starts by presenting the HRES sources and the architectures applied to this case study. The sources models are then presented in order to establish the global energy management system. The study, afterwards, uses real local energy and load data to examine the efficiency and validity of different energy management approaches. However this paper describes a real case study, simulation was used to determine how enhanced energy management algorithms can upgrade the energy management in the campus.

## II. HYBRID PV-WIND SYSTEM ARCHITECTURE

According to numerous studies, different components can be used as basic sources for the HRES [14, 5, 15, 16, 17]. Their configuration adopts one of the main coupling architectures: DC or AC coupling architecture. The combination of different sources on the same bus poses major problems related to voltage regulation and power distribution. Thus the energy management requires a global optimization of the energy transfer in the HRES using a rigorous analysis of different control methods.

In the DC coupling configuration, all the HRES sources are connected to a common DC bus as shown in Fig. 1. The wind turbine in the one hand is coupled to the DC bus via an AC-DC converter; the photovoltaic panels on the other hand are connected to the same bus using multiple DC/DC converters. The power converters used ensure a proper voltage regulation

and an optimal power exploitation using the adapted MPPT systems. The power bank is connected to the DC bus through a bidirectional DC-DC converter to maintain a stable storage-supply balance at the rated capacity of the storage units and the optimal control of the DC bus. The charging regulators in the power bank control the charging level of the storage units (Lead-Acid batteries and super-capacitors) and keep an appropriate balance between load demand and power supply. These measures protect the power bank units against deep discharge and high frequency charging-discharging cycles.

In the DC coupling configuration, DC loads can be powered directly from the central DC bus and the AC loads are fed using appropriate DC-AC converters. Furthermore, this configuration, also referred to it as centralized DC-bus topology has proven to be a relatively simple and affordable solution in term of system control. However, numerous studies have reported some disadvantages related to the DC-coupling configuration, like a reduced system efficiency due to the notable energy losses caused by the power converters [18]. The fact also that HRES sources can't power AC loads directly, DC loads being unbalanced are major drag downs for this configuration [10,16].

In the AC-bus configuration, all renewable energy sources are connected to the main AC bus via DC-AC or AC-AC converters as shown in Fig. 2. Furthermore the injected AC voltage should be carefully regulated and synchronized taking at consideration the grid prescription. The bidirectional converter used with the power bank regulates the voltage delivered by the power bank and feed the storage units when energy excess is available [10, 16]. In this architecture, the AC-bus is the central component where all the energy is delivered. Thus, it is generally easier to face the energy deficit during peak periods by connecting or disconnecting the grid or an electro-generator. A basic performance comparison of the different bus-coupling architectures should primarily take at consideration the load demand variation during a typical day and how the HRES can cover it. Thus, the AC-coupling configuration is the most efficient one [16]. This configuration is generally adopted for its high reliability, modular and scalable structure, especially when the HRES is off-grid. In the present case, all end-users are AC loads. In addition, the loads are fairly dispatched among different locations, for these reasons the AC-coupling configuration was adopted. The MPS architecture for the present case study is presented in Fig. 3.

All energy sources feed a single AC bus trough appropriate power converters. A $50Hz$ voltage adaption powers the different loads that we categorized into two major sections as it will be detailed further along. To ensure a high reliability of the MPS, the system should possess a minimal probability of power failure. Thus, the grid can feed the AC bus with $230VAC - 50\,Hz$ voltage when there is an energy shortage in the HRES. Thus, on normal operating conditions, when there is enough energy, loads are powered using the HRES sources (standalone mode). However, if the load demand exceeds available energy supplied by the HRES sources, loads can draw the needed power out of the power bank if it is possible or the main grid if the storage units are on power outage. The energy excess produced by the HRES is redirected to the power bank to store energy without over charge. These technical needs require an efficient energy management, as proposed further along.



Fig. 1. DC-Coupling Configuration.



Fig. 2. AC-Coupling Configuration.



Fig. 3. MPS Architecture.

### III. HRES Model

Modeling the individual components of each source is first step for sizing a hybrid energy system [19]. This approach can help determine the ideal components characteristics and contribute to the concept of energy management in the MPS.

#### A. Wind Power Source

Wind energy is exploited in the ENSET campus using two wind turbines: Whisper 175 and E70 PRO. The characteristics of each turbine are given in Table I.

The energy delivered by the wind turbine can be calculated using the collected wind speed data. Thus, the wind speed output power can be generally calculated using the following equation:

$$P_w = \frac{1}{2}.C_p(\lambda,\beta).\rho.S_b.W_s^3 \qquad (1)$$

where $\rho$ is the air density (kg/m$^3$), $S_b$ is intercepting area of the rotor blades (m$^2$), $W_s$ is the average wind speed (m/s), $C_p$ is the power coefficient, a function of tip speed ratio ($\lambda$) and pitch angle ($\beta$). The theoretical maximum value of the power coefficient $C_p$ is 0.593, also known as Betz's coefficient. In our case, the pitch angle value is 0° [20]. The Tip Speed Ratio for a wind turbine is mathematically defined as:

$$\lambda = \frac{R.\omega}{W_s} \qquad (2)$$

where R is radius of turbine (m) and $\omega$ is angular speed (rad/s). It should be noted that the wind speed at any reference height can be extrapolated to find the wind speed at a different altitude using the following equation [21]:

$$\frac{W_s(z)}{W_s(z_0)} = \left(\frac{z}{z_0}\right)^\alpha \qquad (3)$$

Where $W_s(z_0)$ is the wind speed at the anemometer' height $z_0$, $W_s(z)$ is the wind speed at the height z and $\alpha$ is a power coefficient relative to the surface roughness.

Furthermore, the total wind turbine energy $E_{w,T}$, during a period $T$, can be calculated using the following equation:

$$E_{w,T} = \sum_{i=1}^{n} P_{w,i}.T_i \qquad (4)$$

TABLE I. WIND TURBINES CHARACTERISTICS

| Parameter | Value | |
|---|---|---|
| | *Whisper* | *E70 PRO* |
| Peak power (W) | 3200 | 5500 |
| Rate power (W) | 2400 | 4000 |
| Blades numbers | 2 | 3 |
| Blades diameters (m) | 4.26 | 4.3 |
| Cut-in wind speed (m/s) | 3.1 | 2 |
| Rated wind speed | 10 | 11 |
| Voltage (V DC) | 24/48 | 24/48 |
| Tower height (m) | 21 | 14 |

Where n indicates the number of sampling slot time $T_i$ during a large time period such as a month or a year.

#### B. Photovoltaic Power Source

The available photovoltaic (PV) power system consists of a multitude of solar panels of different technologies: (i) amorphous technology (PowerMax® 3.5), (ii) monocrystalline technology $(AC-270M/156-60S)$, and (iii) polycrystalline technology $(AC-270P/156-60S)$. Table II summarizes the main parameters of the used PV panels at Standard Test Conditions STC (Air Mass AM1.5, Irradiance $1000W/m^2$, Cell Temperature $25°C$):

Mathematic modeling of the physical phenomena behind the PV system behavior is a critical to establish a well-founded energy management concept. As widely discussed in a large number of papers [19, 22, 23], the mathematical equation describing the current-voltage characteristic of a PV cell is given by the implicit model using a single or double diode. In this case, the single diode model was adopted as represented in Fig. 4. We carefully adopted the same parameters as given by the PV constructor.



Fig. 4. Single Diode PV Model.

TABLE II. SOLAR PANELS CHARACTERITICS

| Parameter | Value | | |
|---|---|---|---|
| | *PowerMax® 3.5* | *AC-270M/156-60S* | *AC-270P/156-60S* |
| Maximum power Pm | 125 Wc | 270Wc | 270Wc |
| Current at maximum power point Impp | 2.981A | 8.80A | 8.71A |
| Voltage at maximum power point Vmpp | 42.4V | 30.94V | 31.12V |
| Short-circuit current Isc | 3.35A | 9.41A | 9.25A |
| Open circuit voltage Voc | 58V | 39.26V | 38.21V |
| Panel efficiency | 11.9% | 16.6% | 16.6% |
| Current temperature coefficient αsc | 0 mA/°C | 0.04%/K | 0.04%/K |
| Voltage temperature coefficient βoc | -170 mV/°C | -0.3%/K | -0.3%/K |
| Temp. Coefficient of Powerkp | -0.39 %/°C | -0.4%/K | -0.42%/K |

The mathematical equation combining the cell voltage $V_c$ and cell current $I_c$ is amply described in [19] and can be expressed as:

$$I_c = I_{ph} - I_{s1}.\left[exp\left(\frac{V_c + R_{sc}I}{a_1 V_{T1}}\right) - 1\right] - \frac{V_c + R_{sc}I}{R_{pc}} \qquad (5)$$

The PV cell generated current is noted $I_{ph}$, $I_{s1}$ and $I_{s2}$ are the reverse saturation current of the model's diode, $V_{T1}$ is the diode thermal voltage, $a_1$ represents the diode ideality constant. $R_{sc}$ and $R_{pc}$ are respectively the series loss resistance and shunt loss resistance of the cell. The output power of a PV cell is given by the following equation:

$$P_c = V_c.I_c \qquad (6)$$

To enhance the power, the PV array consists of $N_s$ and $N_p$ PV cells connected in series and parallel respectively. Therefore, the total power output for the global array will be:

$$P_{array} = N_s.N_p.P_c \qquad (7)$$

The output power of the PV module is primarily related to the solar radiation and the ambient temperature. Thus, the output power of the PV array (kW), denoted $P_{PV}$, is given by:

$$P_{PV} = P_R.f_{PV}.\frac{G_T}{G_{STC}}.\left[1 + k_p.(T_C - T_{STC})\right] \qquad (8)$$

Where $P_R$ is the rated power of the PV module under standard test conditions (kW), $f_{PV}$ is the PV derating factor, $G_T$ the solar radiation (kW/m$^2$), $G_{STC}$ incident radiation at standard test conditions (1kW/m$^2$), $k_p$ the temperature coefficient (0.004°C$^{-1}$, $T_C$ the PV cell operation temperature (°C) and $T_{STC}$ is PV cell temperature under standard test conditions (25°C) [24].

## IV. POWER BANK

In order to achieve a reliable and efficient HRES, a power bank system is used to ensure continuous delivery of power to the loads. The storage system is used as a backup option to store energy when there is a production excess. The power bank is specifically designed to supply power when wind speed and solar radiation are weak or during peak hours when there is a high demand of energy. This study uses a battery-super capacitors hybrid energy storage system (BSHS). Batteries are the primary energy buffer for long durations and normal exploitation, Super-capacitors provide energy during instantaneous massive power demand. This design is supposed to help HRES keep a stable output and meet sudden peak power consumption.

### A. Lead-Acid Batteries

Batteries modeling and for energy management purpose it relies on the determination state of charge (SOC) during exploitation. This criterion is a fundamental key in optimizing energy management and batteries lifetime cycles. Modeling the lead acid batteries used in the HRES [25, 26] should take at consideration the internal parameters of batteries such as the state of charge (SOC), storage capacity, the rate of charge/discharge, ambient temperature and the life cycle [26, 27, 28].

All battery circuit parameters are widely described in [23], where an equivalent circuit of the battery is represented. The SOC being the most important parameters of theme all, it is generally calculated between a starting time $t_0$ and the time of study t and is given by [24]:

$$SOC = SOC_0 + \int_{t_0}^{t}\left(\frac{I_{bat}}{C_{bat}}\right) dt \qquad (9)$$

where $SOC_0$ is the battery's state of charge at the starting time, $C_{bat}$ is the battery capacity (Ah) and $I_{bat}$ is the battery current (A). However, in order to preserve the battery's health, the SOC is generally subjected to the following constraints [24]:

$$SOC_{min} \leq SOC(t) \leq SOC_{max} \qquad (10)$$

Where $SOC_{min}$ and $SOC_{max}$ are the minimum and maximum allowable battery charge. Here the $SOC_{max}$ is the rated nominal capacity of battery bank ( $C_{bat}$ ) and $SOC_{min}$ is determined by depth of discharge (DOD), as described in the following equation:

$$SOC_{min} = (1 - DOD_{max}).C_{bat} \qquad (11)$$

$DOD_{max}$ is the maximum allowable depth of discharge which may decrease battery life [24, 26]. It should be noted that SOC can't be measured directly, but can be estimated as detailed in [29, 30].

### B. Super-Capacitors (SC)

Super-capacitors are components able to store energy directly as electric charges, contrary to batteries storage process which is based on chemical reactions. The super-capacitors have a higher power density compared to batteries and provide high power over a short time lapse. Furthermore, super-capacitors-based storage systems possess longer life cycles and are ideally suited for high power and short discharge applications [14, 31]. Several representative models of super-capacitors have been introduced in literature [23]. However, the classic equivalent circuit for the super-capacitor highlight the effect of the capacitance ($C$), the equivalent series resistance ($R_{ES}$), and the equivalent parallel resistance ($R_{EP}$) [23, 32, 33]. The total usable energy from a SC is given by:

$$E_{SC} = \frac{1}{2}.C_t.\left(V_i^2 - V_f^2\right) \qquad (12)$$

Where $V_i$ and $V_f$ are respectively the initial voltage before the discharging starts and the final voltage after the discharging ends; $C_t$ is the total SC system capacitance. $C_t$ can be the equivalent of series or parallel combination of capacitors. $C_t$ is given by:

$$C_t = \frac{n_p}{n_s}.C \qquad (13)$$

Where, $n_s$ is the number of SCs connected in series to meet the rated $DC$ voltage and $n_p$ is the number of SCs connected in parallel to achieve the required storage capacity.

### C. Hybrid Power Storage Approach

Using both batteries and super-capacitors in the power bank can enhance its performance. Thus, batteries are generally used to exploit the long term storage capability and super-capacitors

are combined during certain time laps where important loads are detected. The Batteries/Super-capacitors coupling can be configured in different ways [34, 35]. The architecture of coupling in the present work is represented in Fig. 5. This configuration uses two bidirectional $DC/DC$ converters and a $DC/AC$ converter for the $AC$ bus interfacing. Furthermore, this configuration guaranties a flexible control of the power bank during power failure.



Fig. 5.    Batteries-Supercapacitors Hybrid Coupling.

For the study in this paper, the batteries/super-capacitors hybrid system is composed of 6 batteries with different coupling possibilities ( $12V, 210AH, 1500A$ ) and 20 Supercapacitors ($1200\,F, 2.7V, 1.2Wh$).

## V.    ENERGY POTENTIAL ASSESSMENT AND LOAD VARIATION

The studied case in this paper was carried out, taking at consideration the wind and solar data provided by a meteorological station installed at the ENSET institute-Mohammedia, Morocco (location: 33°41'23" N, 7°23'23" W) [21]. The wind speed and solar radiation data for a typical day used to simulate HRES architecture performance and energy management options are respectively shown in Fig. 6 and 7.

The load variation is represented by a real case scenario of the ENSET institute. All information about heating, lightening, teaching and other electrical power consumption within the administrative department was carefully collected in order to simulate different scenarios.

Energy consumption in the institute is directly affected by the teaching activities, thus there is a noticeable decrease in power consumption during the July-September holiday period. The evaluated loads in this work are based on the campus buildings architecture. These blocks are named Block 1 to Block 6. Power Consumption is evaluated based on the time period when each block draws energy.



Fig. 6.    Wind Speed Data for a Typical Day.



Fig. 7.    Solar radiation Data for a Typical Day.

## VI.    ENERGY MANAGEMENT SYSTEM

### A.    Strategy Description

The power supplied by renewable energy sources can't be ensured continuously, the proposed HRES must include an integrated energy management system (EMS) to allow an optimal power supply as well in off-grid mode as in grid-connected mode. When the system is operating in off-grid mode, renewable energy sources and the power bank act as the primary source of energy. The grid, however, can always supply additional power when there is a shortage in the HRES. Thus, energy can either be absorbed by the HRES or injected into the grid, if the load power demand and conditional requirements are met. The primary objective is to ensure that the HRES feed the local loads using available renewable energy sources and minimizing the grid contribution as much as possible.

The main decision parameters for the power management strategies are the available instantaneous power from the renewable energy system (wind-sun), the instantaneous state of charge ($SOC_C^B(t)$) of the power bank and the load profile. Thus, the instantaneous available power ($P_{av}(t)$) provided by the hybrid system is defined as:

$$P_{av}(t) = P_w^G(t) + P_{pv}^M(t) + c.P_C^B(t) \tag{14}$$

$P_w^G(t)$ is the power provided by the two wind turbine systems:

$$P_w^G(t) = P_{w1}(t) + P_{w2}(t) \tag{15}$$

$P_{pv}^M(t)$ denotes the total power delivered by the PV modules:

$$P_{pv}^M(t) = \sum_{i=1}^{n} P_{pv_i}(t) \tag{16}$$

$P_C^B(t)$ is the supplied or absorbed power according to the BSHS discharging or charging process respectively according to the c value:

$$c \in [-1,1] \tag{17}$$

Expression (17) means:

$$
\begin{cases}
c = -1 \iff SOC_C^B(t) = SOC_{min}^{BC} & (17.1) \\
c = 1 \iff SOC_C^B(t) = SOC_{max}^{BC} & (17.2) \\
-1 < c < 1 \iff SOC_{min}^{BC} < SOC_C^B(t) < SOC_{max}^{BC} & (17.3)
\end{cases}
$$

The case (17.1) and case (17.2) represents respectively the charging and discharging process of the Power bank units. However, in the case (17.3), the hybrid storage units provide or absorb energy to maintain a minimal consumption from the grid. The instantaneous total loads power ($P_L^T(t)$) delivered by the sources can be expressed as:

$$P_L^T(t) = f.P_{av}(t) + (1-f).P_G(t) \qquad (18)$$

Where $P_G(t)$ is the grid power and f the portion of loads supplied by the HRES system ($0 \le f \le 1$). For example, when $f = 1$, the HRES sources supply the total requested power. Otherwise, if $f = 0$, the totality of the needed power is delivered by the grid.

In the proposed energy management system (EMS), it is possible to connect or disconnect the loads according to the available power of the HRES and Power Bank. Mainly, the energy management algorithm is made to achieve the followings objectives:

- Providing an instantaneous power supply adapting the generated energy to the loads.

- Minimizing the use of the grid.

- Improving the performance of the system by optimizing the use of the Power Bank and avoiding its critical solicitation.

It should be noted that the super capacitor has a faster charge/discharge speed. Thus, this component is less sensitive to high frequency charging cycles. The super capacitors are then the first components to be connected in order to minimize the involvement of batteries as much as possible. The different principles of this approach are described in the following flow chart in Fig. 8.

Furthermore different types of loads have to be managed according to the available HRES power. In order to achieve this goal, a Load Management Module (LMM) was carefully elaborated. The LMM approach is based on the basic principal of load prioritizing; this means that powering a load doesn't only depends on the maximum available energy of the HRES and power bank status, but also on the priority of the targeted load. Thus, the end-use loads are categorized into two major types: primary and secondary appliance. The primary loads are critical. Therefore, they have to be powered without any shifting in time. This type is given a higher priority and is allowed to operate at its scheduled time. The secondary appliance type however, can be shifted in time because the nature of the use doesn't require an immediate powering.

if each primary load $Fl_{i=1,2,\dots,n}$, of a building block $B_{b=1,2,\dots,6}$, has a rated power $P_{i,b}^{Fl}$ then the total power of the n fixed loads ($P_L^F(t)$), in a time t of the day, is expressed by:

$$P_L^F(t) = \sum_{b=1}^{6}\left(\sum_{i=1}^{n} P_{i,b}^{Fl}(t) \times s_{i,b}^{Fl}(t)\right) \qquad (19)$$

Where $s_{i,b}^{Fl}(t)$ is the state of $Fl_{i=1,2,\dots,n}$ for each block B in slot time $t$:

$$s_{i,b}^{Fl}(t) = \begin{cases} 1 \text{ if the load is powered} \\ 0 \text{ if the load is not powered} \end{cases} \qquad (20)$$

Likewise, the total power consumption of all secondary loads ($P_L^S(t)$) taking as consideration the state of each component, is calculated using the following equation:

$$P_L^S(t) = \sum_{b=1}^{6}\left(\sum_{j=1}^{j} P_{j,b}^{Sl}(t) \times s_{j,b}^{Sl}(t)\right) \qquad (21)$$

The instantaneous total power of loads($P_L^T(t)$), written on equation (18), can also be expressed by:

$$P_L^T(t) = P_L^F(t) + P_L^S(t) \qquad (22)$$



Fig. 8. Power Bank Management based on SoC Control.

Based on this approach, after supplying the primary loads with the required power, the LMM module, should supply as much energy as possible to the secondary loads. In order to achieve this function, the LMM module compares the excess of energy to the required power from secondary loads and disconnects their elementary units until the achievement of balance between the supplied power and total load demand. The following algorithm explains how exactly the LMM module works.

---

**Algorithm: Loads Management Module (LMM)**

*Inputs:*
- $\Delta P_F$ : The remaining available power
- $P_L^S$ : The total secondary loads demand
- $n$ : The number of the secondary loads
- $P_L^C$ : The critical loads power (no shiftable loads)

*Outputs:*
- $P_L^S$: The secondary loads demand that can be satisfied
- $f$ : The ratio of secondary loads which can be supplied with$\Delta P_F$

**BEGIN**
1. While $(P_L^S(t) > P_L^C)$ do
2. For (b=1 to 6) and (i=1 to n):
   - Update $P_L^S(t)$ according to (21)
   - Set switches $S_{i,b}^{SL}(t)$ to ensure $(\Delta P_F \geq P_L^S(t))$
   
   EndFor
3. $f = \dfrac{P_L^S(t)}{P_L^S}$
   
   EndWhile
4. Return$(P_L^S(t), f)$

**END**

---

Flowchart of the global energy management is represented in Fig. 9 in order to describe the general functioning of the system.

The global management algorithm basically calculates the total available renewable energy, if it is possible to supply both primary and secondary appliance loads, the power management system connects all loads. In the case of power shortage (when renewable energy can't cover the total energy demand) the LMM disconnects a certain number of secondary loads until achieving the appropriate balance between available energy and load demand $(\Delta P_F(t) \geq P_L^S(t))$. However, when available

renewable energy doesn't even cover the primary load, the energy deficit should be drawn from the grid according to (18). This approach ensures a continuous and uninterruptible supply to the primary load whose functioning is fundamental.

### B. System Modeling and Strategy Simulation

In order to study this strategy in details and determine which strategy scenario is the best in terms of efficiency and functionality, each component of the global energy system has been carefully modeled according to the real data collected from the ENSET campus. Fig. 10 illustrate the global modeled system

As it can be noticed in Fig. 10, six major blocks have to be modeled in order to state about the energy management strategies. The first block represents the power generated from the wind turbine. A physical approach has been adopted in order to simulate the available wind power. Thus, in the wind power energy block, wind turbines, $AC/DC$ , $DC/AC$ and MPPT converters have been used. However, the generated power is conditioned by the available wind potential. In the photovoltaic power source, PV panels, MPPT and $DC/AC$ converters have been used. The power delivered by the PV panels is also subject to the sun irradiation data collected. The power Bank was simulated using models of Lead-Acid batteries and super-capacitors, in addition the $DC/DC$ and $DC/AC$ converters adapt the delivered voltage. The power delivered by the power bank is conditioned by the $SoC$. Thus, when the $SoC$ is below 20%, power is drowning from the grid when the HRES are in shortage. Above 80% power is injected into the grid if there is an excess. Loads also have been modeled using the campus data collected from the concerned building. Furthermore, the grid has been modeled to supply or absorb power depending on the situation if there is an energy shortage or excess. Finally, the energy management system is responsible of finding the best strategy for managing power demand and energy generation. Fig. 11 presents the simulation diagram including the hybrid renewable energy sources, the power bank, the loads and the AC bus where different scenarios of energy management can be injected.



Fig. 9.   Basic Flowchart of the Global Energy Management.

Fig. 10. Global Model of the Studied System.



Fig. 11. Simulation Diagram of the Hybrid Renewable Power Station.

*a) Typical Week-End scenario:* Fig. 12 presents the evolution of the power supplied to primary loads during a typical week-end day using the load profile collected on the campus.

Fig. 13 depicts the evolution of the power delivered to secondary loads during the same typical week-end day.



Fig. 12. Primary Load Evolution Supply During a Typical Week-end Day.



Fig. 13. Secondary Load Evolution Supply During a Typical Week-end Day.

As it can be noticed on Fig. 12 and 13, primary loads are always powered whatever are the sources conditions, due to the criticality of the loads. Secondary loads however are disconnected when there is an energy shortage. Fig. 14 and 15 present respectively the quantity of power generated by Photovoltaic Panels and Wind Turbines.

Data used for simulating delivered power by the HRES was carefully collected on the campus during one year. The power presented in Fig. 14 and 15 is for a typical day using the load data in the buildings. Fig. 16 presents the global power generated by the HRES sources and the energy demand of the loads.

In Fig. 16, the black curve represents the total power generated by both sources PV panels and wind turbines. Green curve represents the power delivered to the secondary loads and the blue line represents the power delivered to the primary load. The red curve is the total power delivered to both loads, thus when the black curve is above the red one, there is enough energy within the HRES to supply both secondary and primary loads. However when the total energy demand is above the available power in the HRES, energy is withdrawn from the Grid. Fig. 17 presents the periods when energy is supplied by the grid.

In order to optimize the energy consumption within the campus, excess of energy generated by the HRES sources should be injected back to Grid when the Power Bank reaches the maximum storage capacity. Fig. 18 presents the energy exchange in the energy management system.

When the energy is below 0 (red line), power is withdrawn from the Grid. However when there is enough energy to cover the total load needs (power curve above 0), the energy excess is injected into the grid. For example, we can clearly see in Fig. 18 that on the middle of the day when both PV and wind potential are available, no power is withdrawn from the grid; however, it supplies most of energy during the first part of the night. The switching of the primary source is a very crucial aspect in the proposed energy management system. Therefore, the Power Bank helps diminishing the frequency of energy source switching, by providing an alternative during short periods of energy demand or temporary HRS disconnection. It is also the main destination where energy is routed before achieving energy excess.

Energy cost is also a key component in the energy management algorithm. Based on Fig. 16, it can be noticed that the secondary loads represented by the green curve are generally almost constant through the day. These loads represent routine work tasks that can be done at any time, automatically by machines or manually using electric tools. Thus these tasks can be shifted to any period of the day without any alteration. The management algorithm can help avoiding unnecessary power consumption during Peak hours. Thus it postpones the task when the not enough power available or when eminent power demand is about to start (power forecasting). Furthermore, the Power Bank is strictly reserved to the primary loads in order to avoid unnecessary energy dissipation. This goes against classic renewable energy systems where generally the Power Bank is not strictly controlled and is used in all energy routing.


Fig. 14. Power Delivered by the Photovoltaic Panels.


Fig. 15. Power Delivered the Wind Turbines.


Fig. 16. General Energy Flow in the Hybrid Renewable Energy Station.


Fig. 17. Energy Supplied by the Grid.


Fig. 18. Grid Energy Exchange During a Typical Week-end Day.

*b) Typical working day scenario:* On a typical working day scenario the global connected load is much more important. Fig. 19 represents the total connected load during typical working days. As it can be noticed in Fig. 19, power consumption on working days are much more important due to an excess of connected loads. Unlike the week-end scenario, during working days, secondary loads are not disconnected due to necessity of the constant use of the grid. However, a more economical approach can be achieved by disconnected as many secondary loads as possible on peak consumption period when power price is higher. Fig. 20 represents the global results obtained during a typical working day. The difference in the available energy during the day was taken into consideration, thus the wind and radiation data is from a different day.

It can be clearly noticed that the available sources in the campus are insufficient for covering the total load demand, especially in the middle of the day. In order to minimize the grid use, the available HRES should be used to compensate the energy drawn from the grid. Fig. 21 presents the energy exchange with grid.


Fig. 19. Total Loads Evolution During a Typical Working Day.


Fig. 20. Global Energy Flow During a Typical Working Day.


Fig. 21. Grid Energy Exchange During a Typical Working Day.

Because of the important load, on a typical working day, most of the energy is drawn from the grid. However the HRES helps reducing the grid's dependency. Nevertheless, the available wind and solar potential are fairly sufficient for covering the campus by upgrading the available HRES sources. The available wind potential can accommodate wind turbines ranging up to 12kW. This measure, along with a more important solar potential, can cover the global energy consumption of the campus, realizing consequently a standalone power status. In the meanwhile, energy consumption optimization can only targets loads categorization, energy exchange with the grid and power compensation in special cases.

## VII. CONCLUSION

In this paper, we analyzed the energy management in a hybrid renewable energy station. Models of each source including wind turbines, PV panels and the Power Bank have been presented in detail. Actual data has been collected in order to model the PV panels, wind turbines and loads in the different buildings. Then, the energy management algorithm has been explained in detail. The main role of the proposed system is to exploit the hybrid renewable energy sources as much as possible. Thus most of the energy is taken from the hybrid renewable sources and the Power Bank is used to accommodate short periods of energy decrease. However when there is an energy shortage prediction the loads are mainly powered using the grid. The actual collected data was implemented on a simulation environment to validate the energy management scenarios. Results have shown that using this technique can help retrieving the small energy quantities during extended period of loads exploitation. To resume, the basic energy has to be extracted from renewable energy sources, the algorithm helps detecting period of times when the Power Bank can help overcome the energy shortage. However, as last resort when it is not possible for the power bank to cover the whole period, the loads are connected to the grid according to their types and criticality. The present approach was mainly based on loads categorizing; two groups of loads (Primary and secondary) were introduced. We intend, in upcoming works to manage energy in the same campus building from a different standpoint, where each load can have a different prominence coefficient depending on the day, appliance and continuity of use.

REFERENCES

[1] Omar Ellabban, Haitham Abu-Rub, FredeBlaabjerg, "Renewable energy resources: Current status, future prospects and their enabling technology", Renewable and Sustainable Energy Reviews,39748–764, (2014). DOI: 10.1016/j.rser.2014.07.113

[2] Binayak Bhandari, Shiva Raj Poudel, Kyung-Tae Lee, Sung-HoonAhn, "Mathematical Modeling of Hybrid Renewable Energy System: A Review on Small Hydro-Solar-Wind Power Generation", International Journal of Precision Engineering and Manufacturing-Green Technology, Vol. 1, No. 2, pp. 157-173, (2014), DOI: 10.1007/s40684-014-0021-4

[3] R. Luna-Rubio, M. Trejo-Perea, D. Vargas-Vázquez, and G. J. Ríos-Moreno, "Optimal sizing of renewable hybrids energy systems: A review of methodologies" Solar Energy, 86, 1077-1088 (2012), DOI: 10.1016/j.solener.2011.10.016

[4] M. Dali, J. Belhadj, X. Roboam, "Hybrid Solare-wind system with battery storage operating in grid–connected and standalone mode: Control and energy management-Experimental investigation", Energy, 35 (2010) 2587-2595. DOI:10.1016/j.energy.2010.03.005

[5] P. G. Arul, V. K. Ramachandaramurthy, and R. K. Rajkumar, "Control strategies for a hybrid renewable energy system: A review",Renewable and Sustainable Energy Reviews, 42, 597-608, (2015), DOI: 10.1016/j.rser.2014.10.062

[6] B. Bhandari, K. T. Lee, G. Y. Lee, Y. M. Cho, and S. H. Ahn, "Optimization of hybrid renewable energy power systems: A review" International Journal of Precison Engineering and Manufacturing - Green Technology, 2, 99-112 (2015) DOI:10.1007/s40684-015-0013-z

[7] Rashid Al Badwawi, Mohammad Abusara, Tapas Mallick, "A Review of Hybrid Solar PV and Wind Energy System", Smart Science, Vol.3, No. 3, 127-138 (2015), DOI:10.1080/23080477.2015.11665647

[8] Tajeddine, K., Abdelhadi, R., Omar, B., & Hassan, O. (2017). A Cascaded H-Bridge Multilevel Inverter with SOC Battery Balancing. INTERNATIONAL JOURNAL OF ADVANCED COMPUTER SCIENCE AND APPLICATIONS, 8(12), 345-350.

[9] TajEddine, K, Abdelhadi, R., Omar, B., & Hassan, O. (2018, April). A new multilevel inverter with genetic algorithm optimization for hybrid power station application. In 2018 4th International Conference on Optimization and Applications (ICOA) (pp. 1-6). IEEE.

[10] T.S. Ustun, C. Ozansoy, A. Zayegh, "Recent developments in microgrids and example cases around the world—A review",Renewable and Sustainable Energy Reviews, 15 (2011) 4030–4041. DOI: 10.1016/j.rser.2011.07.033

[11] H.R. Ghosh, S.K. Nandi, "Techno-economical analysis of off-grid hybrid systems at Kutubdia Island, Bangladesh", Energy Policy,38(2):976–80 (2010).DOI: 10.1016/j.enpol.2009.10.049

[12] B.Y. Ekren, O. Ekren "Simulation based size optimization of a PV/wind hybrid energy conversion system with battery storage under various load and auxiliary energy conditions". Applied Energy 2009;86(9):1387–94.DOI: 10.1016/j.apenergy.2008.12.015

[13] Diaf, S., Belhamelb, M., Haddadic, M., and Louchea, A., "Technical and Economic Assessment of Hybrid Photovoltaic/Wind System with Battery Storage in Corsica Island," Energy Policy, Vol. 36, No.2, pp. 743-754, (2008).

[14] A. Chauhan, R.P. Saini, "A review on Integrated Renewable Energy System based power generation for stand-alone applications: Configurations, storage options, sizing methodologies and control", Renewable and Sustainable Energy Reviews, vol. 38 pp. 99–120, (2014). DOI:10.1016/j.rser.2014.05.079

[15] A. Maleki, A. Askarzadeh, "Comparative study of artificial intelligence techniques for sizing of a hydrogen-based stand-alone photovoltaic/wind hybrid system", International Journal of Hydrogen Energy, Vol.39, Issue 19 (2014), pp. 9973-9984. DOI: 10.1016/j.ijhydene.2014.04.147

[16] I. Patrao, E. Figueres, G. Garcerá, R. González-Medina, "Microgrid architectures for low voltage distributed generation", Renewable and Sustainable Energy Reviews, Vol.43, March 2015, pp. 415–424. DOI: 10.1016/j.rser.2014.11.054

[17] J.J. Justo, F. Mwasilu, J. Lee, J.W. Jung, "AC-microgrids versus DC-microgrids with distributed energy resources: A review", Renewable and Sustainable Energy Reviews, 24 (2013) 387–405,. DOI: 10.1016/j.rser.2013.03.067

[18] Nelson, D.B., Nehrir, M.H., Wang, C. "Unit sizing and cost analysis of stand-alone hybrid wind/PV/fuel cell power generation systems", Renewable Energy, Vol.31, Issue 10, August 2006, pp.1641–1656. DOI: 10.1016/j.renene.2005.08.031

[19] B. Bhandari, S.R. Poudel, K-T. Lee, S-H. Ahn, "Mathematical Modeling of Hybrid Renewable Energy System: A Review on Small Hydro-Solar-Wind Power Generation", International Journal of Precision Engineering and Manufacturing-Green Technology, Vol. 1, No. 2, pp. 157-173. DOI:10.1007/s40684-014-0021-4

[20] A. Raihani, A. Hamdoun, O. Bouattane, B. Cherradi, A. Mesbahi, "An Optimal Management System of a Wind Energy Supplier", Smart Grid and Renewable Energy, 2011, 2, 349-358 DOI:10.4236/sgre.2011.24040

[21] A. Raihani, A. Hamdoun, O. Bouattane, B. Cherradi, A. Mesbahi, "Toward an accurate assessment of wind energy platform of Mohammedia city, Morocco", IRACST – Engineering Science and Technology: An International Journal (ESTIJ), Vol.2, No. 5, October 2012.

[22] A. Ortiz-Conde, D. Lugo-Munoz, "An Explicit Multiexponential Model as an Alternative to Traditional Solar Cell Models With Series and Shunt Resistances", IEEE Journal of Photovoltaics, Vol. 2, No. 3, JULY 2012, pp. 261-268. DOI: 10.1109/JPHOTOV.2012.2190265

[23] P. Bajpai, V. Dash, "Hybrid renewable energy systems for power generation in stand-alone applications: A review", Renewable and Sustainable Energy Reviews 16 (2012) 2926– 2939. DOI:10.1016/j.rser.2012.02.009

[24] E. Dursun, O. Kilic, "Comparative evaluation of different power management strategies of a stand-alone PV/Wind/PEMFC hybrid power system", International Journal of Electrical Power & Energy Systems, Vol. 34, Issue 1, January 2012, pp. 81–89, DOI: 10.1016/j.ijepes.2011.08.025

[25] S. Bogdan, Z.B. Salameh, "Methodology for optimally sizing the combination of a battery bank and PV array in a wind/PV hybrid system". IEEE Transactions on Energy Conversion 11 (2), 367– 375. 1996

[26] J. F. Manwell et J. G. McGowan, "Lead Acid Battery Storage Model for Hybrid Energy Systems", Solar Energy, vol. 50, n°5, pp. 399-405, 1993.

[27] R. A. Jackey, "A Simple, Effective Lead-Acid Battery Modeling Process for Electrical System Component Selection", SAE SP, n° 2130, pp. 17-26, 2007.

[28] Cherif A., M Jraidi, A Dhouib, "A battery ageing model used in stand-alone PV systems", Journal of Power Sources, vol. 112, n 1, pp. 49-53, 2002.

[29] A.H. Anbuky, P.E. Pascoe, "VRLA battery state-of-charge estimation in telecommunication power systems", IEEE Trans. Ind. Electron., vol. 47, no.3, pp. 565-573, jun. 2000

[30] K. Kutluay, Y. Cadirci, Y.S. Ozkazanc, I. Cadirci, " A new online state-of-charge estimation and monitoring system for sealed lead-acid batteries in telecommunication power supplies", IEEE Trans. Ind. Electron., vol. 52, no.5, pp. 1315-1327, Oct. 2005

[31] A.Rebbani, O.Bouattane, L.Bahatti, M.Zazoui, "An Efficient Electric Charge Transfer Device for Intelligent Storage Units". Open Journal of Energy Efficiency, 3, 50-63, (2014). DOI: 10.4236/ojee.2014.33006

[32] O.C. Onar, M. Uzunoglu, M.S. Alam, "Dynamic modeling, design and simulation of a wind/fuel cell/ultra-capacitor-based hybrid power generation system", Journal of Power Sources, Vol. 161, Issue 1, (2006), pp. 707–722. DOI:10.1016/j.jpowsour.2006.03.055

[33] M. Uzunoglu, O.C. Onar, M.S. Alam, "Modelling, control and simulation of a PV/FC/UC based hybrid power generation system for stand-alone applications", Renewable Energy, Vol. 34, Issue 3, (2009), pp. 509–520. DOI:10.1016/j.renene.2008.06.009

[34] J. Lia, R. Xionga, Hao Mua, B. Cornélusseb, P. Vanderbemdenb, D. Ernstb, W. Yuanc, "Design and real-time test of a hybrid energy storage system in the microgrid with the benefit of improving the battery lifetime", Applied Energy, Vol. 218, pp.470–478, May 2018. DOI: 10.1016/j.apenergy.2018.01.096

[35] S.F. Tie, C.W. Tan, A review of energy sources and energy management system in electric vehicles, Renew. Sustain. Energy Rev. 20 (2013) 82– 102, DOI: 10.1016/j.rser.2012.11.077.

# Interaction between Learning Style and Gender in Mixed Learning with 40% Face-to-face Learning and 60% Online Learning

Anthony Anggrawan[1]
Computer Science Study Program
Bumigora University
Mataram, West Nusa Tenggara
Indonesia

Nurdin Ibrahim[2]
Suyitno Muslim[3]
Education Technology Study
Program, Jakarta State University
Jakarta, Indonesia

Christofer Satria[4]
Visual Communication Design Study
Program, Bumigora University
Mataram, West Nusa Tenggara,
Indonesia

*Abstract*—**Student learning styles are important factors that have a strong impact on student performance in learning outcomes. That is why each learning method will produce different learning outcomes for students who have different learning styles. According to the previous study concluded that mixed learning produces learning outcomes that are superior to online and face-to-face learning models, but the questions are how is the difference between learning outcomes between student learning styles in mixed learning, and whether there is an interaction between mixed learning models and student learning styles towards learning outcomes. This study provides a scientific answer solution, by conducting experimental research of mix learning with a mixture of 40% face-to-face material learning and 60% online material learning for the subject of Algorithms and Programming. Based on 2-way ANOVA, T, and SCHEFFE tests towards student learning outcomes in this study, it is found: there are differences in learning outcomes between students who have different learning styles, the learning outcomes of male students achieve better learning outcomes than female students, and there is an interaction between student gender and student learning styles towards learning outcomes, where with further tests, it was found that there is no difference in learning outcomes based on student learning styles of all students except students who have a visual learning style with male sex achieving superior learning outcomes than students who have auditory and kinesthetic learning styles.**

*Keywords*—*Online; face-to-face; mixed learning; algorithm and programming; learning outcome; interaction*

## I. INTRODUCTION

Apart from the learning model, whether face-to-face learning models, online learning or training learning, all lead to the same principle, namely to advance learning so as to produce better learning outcomes. Whereas to find out the learning model whether the learning outcomes are better in learning certain subjects, scientific research is needed.

The achievement of learning outcomes in the cognitive domain is not only limited to the full effect due to the influence of the learning model, but there are other factors that also contribute to the cognitive success of learning in each course. The learning styles of students, student gender, and courses studied are among the factors that contribute to the learning outcomes in addition to the learning model/method.

Learning styles, also known as cognitive styles or learning preferences show how students prefer ways of learning [1][2] and are characteristic behaviors that tend to be relatively stable over time [1]. Learning style is defined as a person's natural way, one's habits and something more suitable for someone in absorbing, processing and mastering new information and skills [3] and is an integral component of the learning environment [4].

There are various learning styles from students, namely visual, auditory and kinesthetic learning styles. Therefore, it is not surprising if the learning style among students is different each other. Students with a visual learning style absorb information very well in the form of visual information such as maps, images, diagrams, graphics and the like [5][6][3][7]. In face-to-face learning this visual type student appreciates the information written on the board and printed material in the textbook [3]. Auditory students prefer learning through the ear or hearing senses [5]. Auditory students are comfortable to study with lectures and discussions. Students with auditory learning style types remember well the reading or saying aloud. Meanwhile, students with kinesthetic learning styles are more pleased if the learning process takes place with an activity or is directly involved in the learning process, in the sense of not having to listen and read [5]. In the classroom, this type of student concentrates more with active teachers, and remembers well when given the freedom to participate in class activities.

Student learning styles are important factors that have a strong impact on student performance in learning outcomes [6]. Student learning styles are important cognitive characteristics that influence the way of learning. However, research conducted by Eudoxie confirms that the learning outcomes of students in the face-to-face learning model are influenced by student learning styles [4]. When learning styles match the learning methods used, better learning outcomes will be obtained, but on the contrary when learning styles and learning methods do not match, learning outcomes significantly deteriorate [8].

Learning methods that benefit a group of students do not mean that this applies to other groups of students [3]. Each learning method will produce different learning outcomes,

depending on the likes or dislikes of students with learning models and also whether or not the learning model matches the learning style of students. The fact of previous study does show that there is a positive influence on online and face-to-face learning models on student learning outcomes with a determination coefficient of 0.25, or in other words, the effect of online learning and face-to-face learning outcomes is 25% [9]. A research was found that online learning can replace face-to-face learning in the cognitive field [10]. Furthermore, It was revealed that the effect of mixed learning model on student learning outcomes is 34.81% [11].

The strength of face-to-face learning is the intense intensity of interaction with learners, facilitating the convenience of cooperative learning and also the clarity of learning material [12]. While online learning is not only interactive, it can also provide learning time according to interests of student (in asynchronous online learning) and lecturers (in synchronous online learning), centered on students and students who build a learning environment [13]. Besides that, online learning also has the ability to utilize various forms of multimedia: text, audio, silent and moving visuals, and other forms for learning purposes [14].

Moreover, according to Roblyer & Doering (2013), online learning has more effective interactions compared to face learning [15]. In face-to-face learning, at least there is interaction between students and instructors and with other students as well as interactions between students and learning materials that are taking place. Whereas in online learning interactions occur between students and learning material presented in learning modules and with other students in collaborative learning in asynchronous online learning. In synchronous online learning, student cognitive interaction with lecturers occurs as in face-to-face learning.

Related to the relationship between learning styles and student gender, a previous finding indicated, there are differences in learning styles of students who study online and students who study face-to-face, where gender is a factor that influences the relationship between learning styles and student involvement in learning [16].

Based on the above description, the questions that arise are whether there are differences in learning outcomes between male and female students, whether student learning styles provide different learning outcomes, and whether there is an interaction between learning styles and gender on student learning outcomes, if so how the results of interactions occur between student learning style and student gender to learning outcomes. So, in turn, it is clear that research into how the interaction between learning style and gender in student learning outcomes in face-to-face and online mixed learning becomes very relevant and important, this research provides solutions to the answers to these questions.

In mixed learning, it should be noted, the best portion of the mix of online learning in mixed learning is between 30% and 79% [17], but according to Agosto et. al. (2013), to get the best mixed structure of blended learning through "trial and error" the learning process [12]. Whereas Heather Kanuka & Liam Rourke (2013) emphasize that there is no standard provision about how much the mixed portion of online

learning in blended learning [18]. Mixing levels of mixed learning can be done on: learning activities, the mix portion (weight) of teaching materials and/or program modules [19].

The advantages of mixed learning are actually in harmony with the fact that students have a positive attitude and flexibility to adapt to mixed learning [20]. Substantially mixed learning provides better effectiveness than learning that only uses face-to-face learning methods [21]. Mixed learning provides two learning environments namely face-to-face learning environment and independent online learning, so it can be said that mixed learning is a representation of a combination of the advantages of online learning and the advantages of face-to-face learning. In other words, due to mixed learning combines face-to-face learning and online learning, so that definitely mixed learning activities take advantage of online and face-to-face learning patterns. According to Sleator (2010), the future learning patterns involve a combination of face-to-face experience and online learning [22]. That is why or the main reason, why this research was conducted on a mixture of learning face-to-face learning and online learning with a choice of 40% mix portion of face-to-face learning materials and 60% mix portion of online learning material.

## II. RESEARCH METHODOLOGY

### A. Participants and the Context of the Study

This research was conducted on mixed learning with a mixture of 40% face-to-face lesson and 60% online lesson in Algorithm and Programming subject matter. The main objective of this study is to know the interplay between learning style and student gender on student learning outcomes in mixed learning with a portion of a mixture of 40% face-to-face learning and 60% online learning.

The population of this study is a class of computer science study program students in the first semester of the 2017/18 academic year at Bumigora University in Mataram, Indonesia. The total number of experimental class students is 50 students randomly selected from the population. The online learning module with the Moodle platform has been designed according to the Semester Teaching Plan and has passed a formative test [23] prepared on computer server of Bumigora University. Every student can study asynchronous online learning module at anytime and anywhere (ubiquitous). The online learning module on the computer server contains the subject matter portion of 60% of the total subject matter. Whereas face-to-face learning is done by the teaching lecturer according to the Semester Teaching Plan which contains 40% of the subject matter section of the total subject matter.

### B. Learning Management System

In this study, the online learning module of the Algorithm and Programming courses used has been formatively evaluated and has been presented at the ICoCSIM international seminar in Mataram, Indonesia that will be published in the Scopus-indexed IOP proceeding [23]. The online learning module in this study uses the Moodle Learning Management System (LMS) application which is one of the best received by users in its segment to create efficient online learning sites. Specific Instructional Objectives, Time

Allocation, Learning Outcomes Indicators, Sub-teaching Materials, Subjects / Sub-topics, Sub-Subjects of Materials, Learning Methods, Learning Media, how to evaluate learning outcomes, and reference books that become learning references formulated in Learning Plans The semester is attached to the online learning module, and is also included with face-to-face lecturers, so as to realize the certainty of learning blended learning-1 mixture with a mix of face-to-face teaching materials of around 40% and online teaching materials of around 60%.

## C. Data Collection Procedure

Data collection is done by surveying each student's learning style, and measuring student learning outcomes with quizzes, midterm, and final semester examination, while distinguishing students' gender based on student names that match the student's electronic entry form. After the data is collected, a summative evaluation or statistical hypothesis testing of the data collected is carried out. Statistical hypothesis testing uses the SPSS statistical application program.

## D. Data Analysis

Because this research is inferential research where the research data is of type ratio, and is carried out on sample data, the research requirements testing is carried out, namely: test the normal data distribution and data homogeneity. Tests of validity and reliability and normality are carried out on learning outcome measurement instruments. While the instrument for measuring student learning styles uses standard VARK (Visual, Aural, Read/write, and Kinesthetic) instruments that have been tested for reliability and validity [4].

The 2-way ANOVA test was conducted to determine whether there was an interaction between student learning styles and gender on student learning outcomes, also whether there were differences in the influence of learning styles on student learning outcomes and whether there were differences in learning outcomes between male and female students.

The comparative analysis was carried out by T test of 2 independent samples in this study to find out the result of comparing learning outcomes between students who have male gender and students who have female gender in mixed learning. Because in the 2-way ANOVA test in this study there was an interaction between the learning style and gender of the students, then further tests were carried out using the SCHEFFE test to find out how the interactions between students' learning styles and student sex occurred.

Based on the previous discussion (in the Introduction), and also by paying attention to the lecturer style that is not the same as the student's learning style, it can make learning difficult for students who have different learning styles with lecturers in face-to-face learning [6], and by noting that online learning facilitates difference learning experiences and learning styles for diverse students [24]. Then the research hypothesis (H1) was decided as follows:

H1: There are differences in learning outcomes between students who have different learning styles towards learning outcomes in learning Algorithm and Programming subject.

With reference to the previous discussion (in Introduction) related to the influence of gender on learning outcomes and also with reference to the results of previous studies that there is a positive relationship between student learning styles and problem solving styles, and it is found that gender has an effect on the problem solving style by students [25], so the research hypothesis (H2) related to the influence of gender on learning outcomes is:

H2: There is a difference in learning outcomes between students who have male gender and students who have female gender in learning Algorithm and Programming subject,

Based on the previous discussion (in the Introduction) regarding interactions that occur in learning outcomes, and coupled with the facts that: (a). in mixed learning, classroom learning provides the social interactions needed for active learning, while online learning offers some flexibility, which is not commonly found in the classroom environment [20]; (b). online learning is a web-based learning environment in accessing learning materials, and realizing student and student interactions, with learning materials and with instructors at anytime and anywhere [26]; (c). scientists agree that the face-to-face classroom learning community offers real and meaningful interactions between students and teachers, where pure online learning cannot replace it [20], and research shows that the use of interactive computer technology in a collaborative approach to e-learning allows for specific educational purposes [27], hence the research hypothesis can be predicted (H3):

H3: There is interaction between gender and student learning styles towards learning outcomes in learning Algorithm and Programming subject.

Refer to the previous discussion in Introduction and taking into account that: visual experience is the main thing in learning to be able to understand and interact with the environment [28]; students with visual learning styles are not easily distracted (disturbed) with a noisy atmosphere; so that it can be decided the research hypothesis (H4) related to differences in learning outcomes between students diverse learning styles towards learning outcomes in mixed learning in this study is:

H4: There are differences in learning outcomes of mixed learning between students who have a visual learning style compared to students who have auditory and kinesthetic learning styles in learning Algorithm and Programming subject.

Moreover, it can be predited that students who have visual learning styles differ in learning outcomes compared to all students who have auditory and kinesthetic learning styles.

Facing a threat to internal validity, it is overcome by: (1) involving the appropriate face-to-face learning control group in mixed learning in this study, so that threats to the internal validity of history and maturation can be avoided; (2) students in online and face-to-face mixed learning have mostly equivalent initial cognitive abilities, where student samples are taken from high school graduate students who are equal (thus having equality in age and basic knowledge of

Algorithm and Programming lessons. So the internal threat of the validity of death or friction can be avoided.

Whereas, against threats to external validity, it is handled in the following ways: (1) avoiding "experimental or biased effects" i.e. deviations from experimental researchers, the mixed learning process is not carried out by researchers but by other lecturers. So that researchers become "blind" or "double blind" in influencing the results of studies. (2) Samples from mixed learning classes are random samples of representative populations, thus overcoming the threat of external validity from "selection-treatment interactions". (3) Mixed learning from this research is a new learning model conducted for Bumigora University students, and maintains that students do not know the purpose for research, so that effectively overcome the threat of reactive influence on external validity. (4) Students receive no more than one treatment, so there is no interaction between the previous treatment and after treatment. In other words the threat of external validity from various treatment disorders can be avoided. (5) The threat to external validity due to pretest treatment did not occur in this study, because this study did not carry out the pretest.

## III. RESULT AND DISCUSSION

Instruments for determining student learning styles using VARK questionnaires are distributed to students at the beginning of face-to-face learning in mixed learning. This instrument was chosen because it is quite widely used by previous researchers who examined the related learning styles published in scientific journals and their validity and reliability. Descriptive analysis conducted on the results of the VARK questionnaire is known that the number of students who have a visual learning style is 18 students, who have a auditory learning style is 22 students, and who have a kinesthetic learning style is 9 students, as presented in Table I.

Table II shows the gender frequency distribution of students who received mixed learning treatment in this study. The number of male students is 31 students and the number of female students is 18 students.

TABLE I. FREQUENCY DISTRIBUTION OF STUDENT LEARNING STYLE

| | | Frequency | Percent | Valid Percent | Cumulative Percent |
|---|---|---|---|---|---|
| Valid | Visual | 18 | 36.7 | 36.7 | 36.7 |
| | Auditory | 22 | 44.9 | 44.9 | 81.6 |
| | Kinesthetic | 9 | 18.4 | 18.4 | 100.0 |
| | Total | 49 | 100.0 | 100.0 | |

TABLE II. FREQUENCY DISTRIBUTION OF STUDENT GENDER

| | | Frequency | Percent | Valid Percent | Cumulative Percent |
|---|---|---|---|---|---|
| Valid | Male | 31 | 63.3 | 63.3 | 63.3 |
| | Female | 18 | 36.7 | 36.7 | 100.0 |
| | Total | 49 | 100.0 | 100.0 | |

Based on the results of the normality test using Shapiro-Wilk as shown in Table III, the significant value of the group of students with male gender is 0.131 and those with female gender are 0.514. Because the significant value is greater than the alpha value (significant value> 0.05), it can be concluded that the learning outcomes data of students learning mixed learning are normally distributed.

The significance of homogeneity of 0.935, 0.62 and 0.051 (> 0.05 of alpha value), shows the variable learning outcomes: quizzes, midterm test, and final exam in groups of students who get mixed learning are derived from homogeneous population data, with Levene Statistic values being 0.007, 3,659 and 4,021. So, the Levene test shows that mixed learning sample data is homogeneous, as shown in Table IV.

To ascertain the extent to which truth and trust in instruments measure student learning outcomes that have not been tested for validity and reliability that are used in this study, the validity test and calculation of reliability for the instruments used are used, as the results are shown in Table V and Table VI.

The validity test of the instrument of this research was carried out using Product Moment Correlation. Because testing validity for instruments that measure mixed learning outcomes as shown in Table V shows that Pearson correlation is 0.666, 0.960 and 0.977, it can be concluded that the research instrument for measuring mixed learning outcomes has high validity.

Cronbach's-Alpha was used in this study to measure the coefficient of internal consistency. The alpha coefficient for the three Cronbach's-Alpha items from the instrument reliability calculation of student learning outcomes in mixed learning is 0.745 (as shown in Table VI), indicating that items have good internal consistency (because after all, the reliability coefficient is 0.70 or higher of 0.70 is considered "acceptable" in most scientific studies).

Based on the 2-way ANOVA test as shown in Table VII, the significance value of student learning styles is 0.000 which is smaller than the alpha value (0.05). This means that learning styles affect student learning outcomes, or in other words there are differences in learning outcomes between students who have different learning styles towards learning outcomes. Thus the H1 hypothesis is accepted, or in other words the null hypothesis (H0) is rejected and the alternative hypothesis H1 is accepted. This means, there are differences in learning between students who have different learning styles towards learning outcomes in learning Algorithms and Programming subject.

TABLE III. OUTPUT RESULT OF NORMALITY TEST

| GenderBld1 | | Kolmogorov-Smirnov[a] | | | Shapiro-Wilk | | |
|---|---|---|---|---|---|---|---|
| | | Statistic | df | Sig. | Statistic | df | Sig. |
| Score Bld1 | Male | .154 | 31 | .058 | .947 | 31 | .131 |
| | Female | .101 | 18 | .200* | .955 | 18 | .514 |

TABLE IV.     OUTPUT RESULT OF HOMOGENEITY TEST

|  | Levene Statistic | df1 | df2 | Sig. |
|---|---|---|---|---|
| Quiz | .007 | 1 | 47 | .935 |
| Midterm | 3.659 | 1 | 47 | .062 |
| Examination | 4.021 | 1 | 47 | .051 |

The significance of the value of student gender (0.039) is smaller than the alpha value (0.05), meaning that the gender of the students significantly affected learning outcomes, or there was a difference in learning outcomes between students who have male gender and students who have female gender. Thus the H2 hypothesis is accepted, i.e. there are differences in learning outcomes between students who have male gender and students who have female gender in learning Algorithms and Programming subject.

While the significance value of the influence of gender and learning styles on student learning outcomes (0.017) is smaller than the alpha value (0.05). This indicates that there is interaction between gender and student learning styles towards learning outcomes, or gender and learning styles together have a significant effect on learning outcomes. So the H3 hypothesis is accepted, that is, there is an interaction between gender and student learning outcomes towards learning outcomes in the learning algorithm and programming subject.

Based on the results of comparison of learning outcomes between male students and female students using T test of 2 independent samples (as shown in Table VIII and Table IX), it is known that the significant value of the test results is 0.00 smaller than the alpha value of 0.05, this confirms that there is significant difference in learning outcomes between male and female students.

TABLE V.     VALIDITY TEST RESULT OF LEARNING OUTCOMES INSTRUMENT

|  |  | Quiz | Midterm | Examination | ScoreBld1 |
|---|---|---|---|---|---|
| Quiz | Pearson Correlation | 1 | .588** | .515** | .666** |
|  | Sig. (2-tailed) |  | .000 | .000 | .000 |
|  | N | 49 | 49 | 49 | 49 |
| Midterm | Pearson Correlation | .588** | 1 | .919** | .960** |
|  | Sig. (2-tailed) | .000 |  | .000 | .000 |
|  | N | 49 | 49 | 49 | 49 |
| Examination | Pearson Correlation | .515** | .919** | 1 | .977** |
|  | Sig. (2-tailed) | .000 | .000 |  | .000 |
|  | N | 49 | 49 | 49 | 49 |
| ScoreBld1 | Pearson Correlation | .666** | .960** | .977** | 1 |
|  | Sig. (2-tailed) | .000 | .000 | .000 |  |
|  | N | 49 | 49 | 49 | 49 |
| **. Correlation is significant at the 0.01 level (2-tailed). | | | | | |

TABLE VI.     REALIBITY TEST RESULT OF LEARNING OUTCOMES INSTRUMENT

| Cronbach's Alpha | N of Items |
|---|---|
| 745 | 3 |

TABLE VII.     THE RESULT OF THE 2-WAY ANOVA TEST ON THE VARIABLE OF LEARNING OUTCOME AND FIXED FACTORS OF LEARNING STYLES AND STUDENT GENDER

| Dependent Variable: ScoreBld1 | | | | | |
|---|---|---|---|---|---|
| Source | Type III Sum of Squares | df | Mean Square | F | Sig. |
| Corrected Model | 2014.296a | 5 | 402.859 | 13.839 | .000 |
| Intercept | 194554.820 | 1 | 194554.820 | 6683.576 | .000 |
| GenderBld1 | 132.524 | 1 | 132.524 | 4.553 | .039 |
| VAKBld1 | 893.167 | 2 | 446.584 | 15.342 | .000 |
| GenderBld1 * VAKBld1 | 262.640 | 2 | 131.320 | 4.511 | .017 |
| Error | 1251.704 | 43 | 29.109 |  |  |
| Total | 256275.000 | 49 |  |  |  |
| Corrected Total | 3266.000 | 48 |  |  |  |
| a. R Squared = .617 (Adjusted R Squared = .572) | | | | | |

Due to in the T test, the value of t is positive (4.213) and the average value of male students (7.63) is higher than the average value of female students (5.25), it can be concluded that students with male gender are superior to their learning outcomes compared to students with female gender in learning outcomes Algorithm and Programming course in mixed learning with 40% portion of face-to-face learning and 60% portion of online learning.

In the further test with the SCHEFFE test, as the results are shown in Table X, it was found that there were no differences in learning outcomes between male and female gender students with auditory and kinesthetic learning styles with all other students who had auditory and kinesthetic learning styles.

Likewise, there is no difference in learning outcomes between female gender students who have visual learning style with all students who have kinesthetic learning style, both male and female, and with male students who have visual learning style. All female students who have a visual learning style achieve different learning outcomes compared to all students who have auditory learning styles, both male and female. Further more, all male gender students who have a visual learning style differ in learning outcomes compared to all students who have auditory and kinesthetic learning styles, both male and female. This means that the H4 research hypothesis is accepted, namely there are differences in learning outcomes of mixed learning between students who have a visual learning style compared to students who have auditory and kinesthetic learning styles in learning Algorithms and Programming subject.

TABLE VIII.     THE AVERAGE VALUE OF LEARNING OUTCOMES OF THE RESULT OF THE T TEST TO COMPARE LEARNING OUTCOMES BETWEEN MALE STUDENTS AND FEMALE STUDENTS

| Group Statistics | | | | | |
|---|---|---|---|---|---|
| GenderBld1 | | N | Mean | Std. Deviation | Std. Error Mean |
| ScoreBld1 | Male | 31 | 75.5742 | 7.62670 | 1.36979 |
|  | Female | 18 | 67.0083 | 5.24565 | 1.23641 |

TABLE IX. THE RESULT OF THE T TEST OF COMPARISON OF LEARNING OUTCOMES BETWEEN MALE STUDENTS AND FEMALE STUDENTS

| Independent Samples Test | | Levene's Test for Equality of Variances | | t-test for Equality of Means | | | | | 95% Confidence Interval of the Difference | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | F | Sig. | t | Df | Sig. (2-tailed) | Mean Difference | Std. Error Difference | Lower | Upper |
| ScoreBld1 | Equal variances assumed | 3.060 | .087 | 4.213 | 47 | .000 | 8.56586 | 2.03330 | 4.47539 | 12.65633 |
| | Equal variances not assumed | | | 4.642 | 45.500 | .000 | 8.56586 | 1.84528 | 4.85041 | 12.28131 |

TABLE X. SCHEFFE TEST (POSTHOC ADVANCED TEST) OF VARIABLE LEARNING OUTCOMES WITH FIXED FACTORS LEARNING STYLES * STUDENT GENDER

| (I) InteractionBld1 | (J) InteractionBld1 | Mean Difference (I-J) | Std. Error | Sig. | 95% Confidence Interval | |
|---|---|---|---|---|---|---|
| | | | | | Lower Bound | Upper Bound |
| AFemale | AMale | -5.92 | 2.340 | .289 | -14.08 | 2.24 |
| | KFemale | -3.25 | 3.242 | .960 | -14.56 | 8.06 |
| | KMale | 1.00 | 3.009 | 1.000 | -9.49 | 11.49 |
| | VFemale | -7.20 | 3.009 | .352 | -17.69 | 3.29 |
| | VMale | -16.54* | 2.340 | .000 | -24.70 | -8.38 |
| AMale | AFemale | 5.92 | 2.340 | .289 | -2.24 | 14.08 |
| | KFemale | 2.67 | 3.085 | .979 | -8.08 | 13.43 |
| | KMale | 6.92 | 2.839 | .330 | -2.98 | 16.82 |
| | VFemale | -1.28 | 2.839 | .999 | -11.18 | 8.62 |
| | VMale | -10.62* | 2.116 | .001 | -18.00 | -3.24 |
| KFemale | AFemale | 3.25 | 3.242 | .960 | -8.06 | 14.56 |
| | AMale | -2.67 | 3.085 | .979 | -13.43 | 8.08 |
| | KMale | 4.25 | 3.619 | .924 | -8.37 | 16.87 |
| | VFemale | -3.95 | 3.619 | .943 | -16.57 | 8.67 |
| | VMale | -13.29* | 3.085 | .007 | -24.05 | -2.53 |
| KMale | AFemale | -1.00 | 3.009 | 1.000 | -11.49 | 9.49 |
| | AMale | -6.92 | 2.839 | .330 | -16.82 | 2.98 |
| | KFemale | -4.25 | 3.619 | .924 | -16.87 | 8.37 |
| | VFemale | -8.20 | 3.412 | .347 | -20.10 | 3.70 |
| | VMale | -17.54* | 2.839 | .000 | -27.44 | -7.64 |
| VFemale | AFemale | 7.20 | 3.009 | .352 | -3.29 | 17.69 |
| | AMale | 1.28 | 2.839 | .999 | -8.62 | 11.18 |
| | KFemale | 3.95 | 3.619 | .943 | -8.67 | 16.57 |
| | KMale | 8.20 | 3.412 | .347 | -3.70 | 20.10 |
| | VMale | -9.34 | 2.839 | .076 | -19.24 | .56 |
| VMale | AFemale | 16.54* | 2.340 | .000 | 8.38 | 24.70 |
| | AMale | 10.62* | 2.116 | .001 | 3.24 | 18.00 |
| | KFemale | 13.29* | 3.085 | .007 | 2.53 | 24.05 |
| | KMale | 17.54* | 2.839 | .000 | 7.64 | 27.44 |
| | VFemale | 9.34 | 2.839 | .076 | -.56 | 19.24 |

Based on observed means.
 The error term is Mean Square(Error) = 29.109.

*. The mean difference is significant at the 0.05 level.

## IV. Conclusion

Based on the results of statistical tests: ANOVA, T and SCHEFFE conducted on student learning outcomes in this study, it is known: that: (1) there are differences in learning outcomes between students who have different learning styles; (2) the learning outcomes of male students achieve better learning outcomes than female students; (3) there is an interaction between student gender and student learning styles towards learning outcomes (gender and learning style have a significant effect on learning outcomes); (4) there is no difference in learning outcomes based on student learning styles of all students except students who have a visual learning style with male sex achieving superior learning outcomes than students who have auditory and kinesthetic learning styles as well as students who have a visual learning style with female sex achieving superior learning outcomes than students who have auditory learning style.

Some constructive suggestions for the direction of future research are: (1) researching learning outcomes in blended learning with other mixed levels, so that it can be ascertained what level of mixed portions can produce better learning outcomes, including interactions that occur between learning styles and blended learning learning patterns in learning Algorithms and Programming and also in other lesson learning; (2) researching the differences in learning outcomes and interactions that occur from two or more blended learning patterns that have different mixed levels in other subjects besides the Algorithm and Programming subject.

## Acknowledgment

### References

[1] P. Fenrich, "Getting Practical with Learning Styles in ' Live ' and Computer-based Training Settings," Informing Science and Information Technology, vol. 3, 2006.

[2] J. S. Sitton, "Student Learning Style Preferences in College-Level Biology Courses: Implications for Teaching and Academic Performance," 2009.

[3] W. B. Rhouma, "Perceptual Learning Styles Preferences and Academic," vol. 09, no. 02, pp. 479–492, 2016.

[4] G. D. Eudoxie, "Learning Styles among Students in an Advanced Soil Management Class : Impact on Students ' Performance," pp. 137–144, 2011.

[5] D. Indriana, Variety of Teaching Media Tools: Knowing, Designing, and Practicing it. 2011.

[6] S. Psycharis, E. Botsari, and G. Chatzarakis, "Examining the Effects of Learning Styles, Epistemic Beliefs and the Computational Experiment Methodology On Learners' Performance Using the Easy Java," J. Educ. Comput. Res., vol. 51, no. 1, pp. 91–118, 2014.

[7] Ž. Pekić, "The Impact of Felder's Learning Styles Index on Motivation

and Adoption of Information Throught E-Learning," vol. 6, no. December, pp. 93–100, 2016.

[8] C. O. Leary and J. Stewart, "The Interaction of Learning Styles and Teaching Methodologies in Accounting Ethical Instruction," J. Bus Ethiccs, vol. 113, pp. 225–241, 2013.

[9] A. Anggrawan, "Correlation of Learning Models to Student Learning Outcomes on Java Programming Language Learning," J. Mantik Penusa, vol. 2, no. 2, 2018.

[10] A. Anggrawan and J. Qudsi, "Comparative Analysis of Online E-Learning and Face To Face Learning: An Experimental Study," Accept. Int. Conf. ICIC 2018 Palembang will be Publ. IEEE Explor. J., 2018.

[11] A. Anggrawan, N. Ibrahim, and C. Satria, "Effect of Blended Learning Patterns on Student Learning Outcomes in Learning Algorithms and Programming," Still Process being Publ. a reputable jou, 2019.

[12] D. E. Agosto, A. J. Copeland, and L. Zach, "Testing the Benefits of Blended Education: Using Social Technology to Foster Collaboration and Knowledge Sharing in Face-To-Face LIS Courses," J. Educ. Libr. Inf. Sci., vol. 54, no. 2, pp. 94–107, 2013.

[13] M. Simonson, S. Smaldino, M. Albright, and S. Zvacek, Teaching and Learning at a Distance: Foundation of Distance Education. 2012.

[14] M. K. Clark, Tom and Barbour, "Online, Blended, and Distance Education in Schools: Building Successful Program." Stylus Publishing, United States of America, 2015.

[15] A. H. Roblyer, M. D and Doering, Integrating Educational Technology in Teaching. Boston: Pearson, 2013.

[16] B. N. Martin and D. Garland, "Do Gender and Learning Style Play a Role in How Online Courses Should Be Designed ?," vol. 4, no. 2, pp. 67–81, 2005.

[17] I. E. Allen, J. Seaman, and R. Garrett, The Extent and Promise of Blended Education in the United States. 2007.

[18] H. Kanuka and L. Rourke, "Using blended learning strategies to address teaching development needs : How does Canada compare ?," Can. J. High. Educ., vol. 43, no. 3, pp. 19–35, 2013.

[19] P. Lieser and S. D. Taff, "Empowering Students in Blended Learning," J. Appl. Learn. Technolornal, vol. 3, no. 3, pp. 6–13, 2013.

[20] C. M. Tang and L. Y. Chaw, "Readyness for Blended Learning : Understanding Attitude of University Students," Int. J. Cyber Soc. Educ., vol. 6, no. 2, pp. 79–100, 2013.

[21] V. N. Johan, "Blended learning Educational programs Web-based instruction Face-to-face communication Computers in education," Comput. Educ., vol. 103, no. 2, p. p16, 2016.

[22] R. D. Sleator, "The evolution of eLearning Background , blends and blackboard ...," Sci. Prog., vol. 93, no. 3, pp. 319–334, 2010.

[23] A. Anggrawan, C. Satria, N. Ibrahim, M. Suyitno, and M. Zarlis, "Instructional Design for Online Learning of Algorithm and Programming," in The Third Internaitonal Conference on Computational Science and Information Management ICoCSIM-2019, 2019, pp. 1–5.

[24] M. A. Suparman, Modern Instructional Design: Teachers Guide and Innovators Education. Jakarta: Erlangga, 2014.

[25] S. Gholami and M. S. Bagheri, "Relationship between VAK Learning Styles and Problem Solving Styles regarding Gender and Students ' Fields of Study," J. Lang. Teach. Res., vol. 4, no. 4, pp. 700–706, 2013.

[26] N. Ibrahim and Hasbullah, Application of Blended Learning Learning at Kandanghaur-Indramayu Open Middle School. Jakarta Indonesia: State University of Jakarta, 2015.

[27] R. D. Quinn, "E-Learning in Art Education: Collaborative Meaning Making Through Digital Art Production," J. Art Educ., no. July, pp. 19–25, 2011.

[28] N. Ibrahim, The Remote Open Education Perspective, Theoretical and Application Studies. 2010.

# Frequency Reconfigurable Vivaldi Antenna with Switched Resonators for Wireless Applications

Rabiaa Herzi[1], Ali Gharsallah[4]

Unit of Research CSEHF Faculty of Sciences of Tunis
El Manar, Tunis 2092, Tunisia

Fethi Choubani[2], Mohamed-Ali Boujemaa[3]

INNOVCOM Laboratory, SUPCOM, University of
Carthage, Tunis, Tunisia

*Abstract*—In this paper, a frequency reconfigurable Vivaldi antenna with switched slot ring resonators is presented. The principle of the method to reconfigure the Vivaldi antenna is based on the perturbation of the surface currents distribution. Switched ring resonators that act as a bandpass filter are printed in specific positions on the antenna metallization. This structure has the ability to reconfigurate between wideband mode and four narrow-band modes which cover significant wireless applications. Combination of the bandpass filters and tapered slot antenna characteristics achieve an agile antenna capable to operate in UWB mode from 2 to 8 GHz and to generate multi-narrow bands at 3.5 GHz, 4GHz, 5.2 GHz, 5.5 GHz, 5.8 GHz and 6.5 GHz. The measurement and simulation results show good agreement. This antenna is an appropriate solution for wireless applications which require reconfigurable Wideband multi-narrow bands antenna.

*Keywords—Frequency reconfigurable; Vivaldi Antenna (VA); Ultra-Wideband (UWB); slot ring resonator; wireless applications*

## I. INTRODUCTION

Research in antenna development has attracted many researches to satisfy the requirements of modern wireless applications such as developing active, compact and miniaturized antennas that can group many services [1-3].

Developing of radio systems for multimode terminal applications which involve a combination of Wi-Fi, WLAN, Wimax, Bluetooth… is obligatory [4,5]. So, using systems with frequency reconfigurable operation and wideband spectrum sensing is a suitable solution due to its compactness, flexibility and its capability to operate over multiple bands [5-6].

The reconfigurable antenna is an antenna which has the capability to change dynamically its radiation characteristics such as its radiation patterns, frequency operations or its polarization [7-8]. Usually, the reconfigurability of antenna is achieved using PIN diodes, Varactors or MEMS to change its geometrical characteristics [1,4].

Because of their benefits of flexibility, the capability to reduce interferences and compactness, frequency reconfigurable antenna is a better alternative which can be used in cognitive radio and multi-mode applications [9]. More precisely, they are many modern systems which have great numbers of antennas that are operated at different frequencies. Therefore, frequency reconfigurable antenna, which can support many functions at various frequencies bands and can significantly decrease the hardware cost and size, is required.

There are many types of frequency reconfigurable antennas such as switching between different narrow bands, wideband to notch band reconfiguration, wideband to narrowband switching [10-11].

Achieving an antenna which has the capacity of wideband to multi-narrow bands reconfiguration is very important and essential for several applications such as a cognitive radio that uses wideband sensing and multi-bands communications [10, 12, 13].

Because of their better radiation performances as well as Ultra-wide bandwidth, elevated gain, and compact structure [14-15], Vivaldi antenna is the best selection to be used in wideband to multi-bands reconfiguration.

Developing of reconfigurable antenna between wideband and multiple narrow-bands has received considerable attention. In [16] and [17] two different switchable Vivaldi antennas have been investigated where several PIN diodes and capacitors are used which augment the complexity of the antenna design. In addition, the reconfigurability of both antennas has notably deteriorated the antenna gain which has decreased by about 2.5 dBi. A reconfigurable Vivaldi antenna was designed in [18], where the obtained agile narrow bands are close with inconstant radiation patterns over the operating frequency range. A wideband frequency agile patch antenna is proposed in [19]. The reconfigurability is achieved using Varactors which require high bias voltages. Very close frequency reconfigurable bands are attained from 1.47 to 1.84 GHz.

In this paper, a frequency reconfigurable Vivaldi antenna for wireless applications that has the capacity to switch from an UWB of 2-8 GHz to essential narrow bands is proposed. By inserting of switchable slot ring resonators using PIN diodes, several narrow bands are obtained. In Section 2, details of the UWB antenna and the proposed antenna design are described. The principle of the method used to achieve a wide to narrow-bands reconfigurable Vivaldi antenna is demonstrated in Section 3. In Section 4, the simulated and measured results obtained throughout the switch modes are explained and discussed. Finally, Section 5 presents the conclusion of this paper.

## II. UWA AND RECONFIGURABLE ANTENNAS DESIGN

The basic structure of the proposed antenna, which is shown in Fig. 1(a) is an UWB Vivaldi antenna that operates from 2 to 8 GHz. It consists of an exponentially tapered slot

printed on an FR-4 substrate that has 4.7 of permittivity and 1.575 mm of thickness. It is fed by a microstrip line printed on the other side of the substrate. The dimension of the antenna is L= 85 mm and W=70 mm. The exponential edge of the aperture antenna is determined using next equations:

Where f, $(x_1,y_1)$ and $(x_2,y_2)$ $Y = \pm(Ae^{fx} + B)$      (1)

$$A = \frac{y_2 - y_1}{e^{fx_2} - e^{fx_1}} \tag{2}$$

$$B = \frac{y_1 e^{fx_2} - y_2 e^{fx_1}}{e^{fx_2} - e^{fx_1}} \tag{3}$$

Are respectively the exponential factor, the peak point and the bottom point of the exponential edge.

Fig. 1(b) shows the wide-narrow bands agile antenna configuration. The Vivaldi antenna geometry is modified by inserting four slot resonators at specific positions in order to perturb the surface currents flow. Each resonator is formed by two symmetric ring slots connected by a rectangular gap. So, three pairs of ring slots are printed on the antenna metallization and coupled into the tapered slot through four pairs of gaps with 1 mm of width. The outer and the inner radius of the ring slots are respectively 5 mm and 3 mm. Moreover, ten PIN diodes are inserted in the opening of the perturbing slots to switch the function of the different resonators and control the current distributions.



(a)          (b)

Fig. 1. Simple Vivaldi Antenna (a) and Reconfigurable Vivaldi Antenna with PIN Diodes (b).

## III. PRINCIPLE OF RECONFIGURATION OF THE VIVALDI ANTENNA

The Ultra-Large Band that characterizes the Vivaldi antenna is the result of the tapered slot geometry that can be devised into two parts named the propagation and radiation parts. The propagation part guides the waves to the radiation area. Each level of the radiation part radiate at the corresponding frequency, where the width of the slot is about the half wavelength.

Generally, frequency reconfigurable Vivaldi antenna can be achieved by disturbing the surface currents flow which is attained by distorting the inner edges of the tapered slot. Moreover, inserting slot ring resonator can act as a filter which has the capability to pass a frequency band and block others. Single ring slot acts as a stop band filter, blocking the

frequencies that correspond to a quarter of the wavelength ($\lambda/4$), while double serial ring slots act as pass band filter which pass frequencies corresponding to a half wavelength ($\lambda/2$). Then, the resonance current path length (D) the stop band and pass band filters can be expressed respectively as following:

$$D_S = \frac{\lambda_{\text{eff}}}{4} \approx (P + s) \tag{4}$$

$$D_P = \frac{\lambda_{\text{eff}}}{2} \approx (P + 2s) \tag{5}$$

The resonance frequencies of the stop band and pass band filters can be determined as follows:

$$F_S = \frac{c}{4\sqrt{\varepsilon_{eff}}(P+s)} \tag{6}$$

$$F_P = \frac{c}{2\sqrt{\varepsilon_{eff}}(P+2s)} \tag{7}$$

Where c, $\varepsilon_{\text{eff}}$, P and S are, respectively, the velocity of light in the free space, the effective permittivity of the substrate, the perimeter of the ring slot which is equal to 31.4 mm and the width of the split.

Fig. 2 and Fig. 3 illustrate the structure and the response of a single and double rings slot resonator. It is clear that the single slot ring resonator products a stop band response around 2.14 GHz while the double rings resonator offers a pass-band filter at 3.65 GHz.

The equivalent circuit of the double slot rings resonators which consists of two associated stop band filters is shown in Fig. 4. Two LC circuits are connected in series, where each LC circuit shows a stop band filter.



(a)



(b)

Fig. 2. Configuration of Single Ring Slot Resonator (a) and the S-Parameters Response (b).

Fig. 3.    Configuration of Double Rings Slot Resonator (a) and the S-
Parameters Response (b).



Fig. 4.    Equivalent Circuit of the Double Rings Slot Resonator.

To investigate the effects of the slots ring resonator in the radiation characteristics of the Vivaldi antenna, a double rings slot resonator is inserted in the beginning of the radiation part. Fig. 5 presents the surface currents distribution of the Vivaldi antenna with and without printed resonator. It can be seen that the current before inserting perturbation follows the direction of the inner radiation edges of the antenna. On the other hand, the surface current is distorted by inserting of the ring slots; it is significantly reduced along the inner edges. So the major asset of this method resides in the capability to promote the disruption of the currents in the annular slots and stopped its repartition along the radiation edges situated after the disturbance. In addition, Fig. 6 compares the reflection coefficients of the simple and the modified antenna. It is clear that the impedance matching is distorted by inserting of the ring slots in the beginning of the radiation part, where the reflection coefficient is up to -10 dB over the frequency range. So, up to the perturbation, the surface currents distribution along the edges of the radiation area is blocked.

As can be seen, inserting a ring slot resonator has two roles: to generate a band-pass response at the frequency corresponding to the half wavelength and to block the surface currents distribution which makes the area situated up to the perturbation in the off state. So, reconfigurable Vivaldi antenna can be achieved by varying the level (Y) of the perturbation and the width of the resonant current path (D) that are the important parameters where the inserted resonator can operate as a filter which passes a range of frequencies and

blocks others. In this way, four modes are investigated. Fig. 7 shows four proposed structures with different positions of the printed resonator. Fig. 8 compares the reflection coefficients of the different structures where four narrow-band modes are achieved.



Fig. 5.    Surface Currents Distribution of (a) the Simple Vivaldi Antenna and
(b) Vivaldi Antenna with Annular Slots at 3.5 GHz.



Fig. 6.    Reflection Coefficients of the Antenna with and without
Perturbation.



Fig. 7.    Antenna Configuration with different Positions of the Printed
Resonator.

(a)



(b)



(c)



(d)

Fig. 8. Simulated Reflection Coefficients of four Configuration with different Positions of the Printed Resonator.

First, the level of the perturbation is about 16.5 mm and the current path length (D) is about 24 mm which practically corresponds to the half wavelength of 6.5 GHz resonance frequency. So, a narrow band at 6.5 GHz is obtained. Second, the elevation (Y) of the perturbation is increased, it is approximately 29 mm and the current path length (D) is 27 mm which means to the half wavelength of 5.2GHz. Fig. 8(b) shows the return loss of this configuration which presents a narrow bandwidth at 5.15 GHz. Third, in Fig. 8(c) two narrow bands are obtained at 4 GHz and 5.8 GHz where the perturbation level is about 47.5 mm and D is about 38 mm (approximately equal to the half wavelength of 4 GHz). Finally, the level (Y) of the perturbation is increased as shown in Fig. 7(d); it is about 53mm where the width (D) is 41 mm which is equal to the half wavelength of 3.5 GHz resonance frequency. This configuration can offer three bands; two narrow bands at 3.5 GHz and 5.2 GHz and wideband from 6.2 GHz to 7.78 GHz. It can be concluded that increasing the level of disturbance enlarge the radiating region of the Vivaldi antenna which results the radiation of the antenna at lower frequencies and its matching at wide bandwidths. As well, this method can be used to properly control the frequency operating of the antenna which facilitates its command and agility.

## IV. FOUR MODES RECONFIGURABLE VIVALDI ANTENNA

To demonstrate the functionality of the antenna to reconfigure between the four modes previously explained, the four perturbations will be printed on the antenna and switched using PIN diodes as shown in Fig. 9(a). The frequency reconfigurability of the antenna is achieved by setting in the OFF state of the diodes that command one perturbation and in the ON state the other diodes which command the opening of the other perturbations. In this way, one perturbation will be activated which can force the disruption of the surface current and stop its circulation up, this allows the function of the antenna in one mode. Fig. 9(b) compares the different reflection coefficients of the four switched modes where sex narrow and agile bands are obtained. The first mode (M1) takes place when the four diodes (D1, D2, D3 and D4), which control the first perturbation, are in the ON state. In this case, one narrow band is obtained at 6.25 GHz. In the second mode (M2) which is attained by switching only the diodes D5 and D6 in the ON state, one narrow band is obtained that resonates at 5.2 GHz. Once the diodes D7 and D8 are activated, the third mode (M3) occurs which can present two narrow bands at 3.97 GHz and 5.5 GHz. At last, the fourth mode (M4) is obtained when the superior perturbation is opened by the deactivation of diodes D9 and D10. This mode can offer two narrow bands at 3.5 GHz and 5.8 GHz. Because of the effect of PIN diodes, it can be observed that activation of the different modes using active switches has slightly affected the frequency response of the antenna. For example, activation of the fourth mode (M4) has result two narrow bands at 3.5 GHz and 5.8 GHz which is different to the result obtained in the precedent section with the passive configuration in mode M4.

(a)



(b)

Fig. 9.   Reconfigurable Vivaldi Antenna (a) and Reflection Coefficients or different Switching Modes (b).

has not greatly affected the radiation patterns of the antenna. Table I summarizes the simulated and measured radiation performances which are obtained through the different modes of switching.



(a)



(b)



(c)



(d)

Fig. 10.  Designed and Fabricated Prototype for different Modes: (A) Mode M1, (B) Mode M2, (C) Mode M3, (D) Mode M4.

To avoid the complexity of the bias circuit for PIN diodes and its effects in the antenna performances, four passive configurations of the defined modes are prototyped. Fig. 10 presents the designed and fabricated prototype for the different modes using ideal switches where the ON state of the PIN diode is modeled by a short-circuit while the OFF state is modeled by an open circuit. The surface current distributions excited at the resonance frequencies of the different modes using ideal switches are shown in Fig. 11. For the different cases, it is noticeably observed that the surface current distribution is more concentered at the active ring slot. Moreover, the active perturbation can stop the repartition of the surface currents on the upper part of the antenna which force the resonance at the desired frequency.

The simulated and measured reflection coefficients of the four switched modes are compared in Fig. 12. A good agreement between simulation and measured results is achieved, accepting slight difference is observed.

Fig. 13 plots an example of prototype antenna in the anechoic chamber, while Fig. 14 illustrates the simulated and measured radiation patterns in the E and H-planes at the different frequencies obtained throughout the four switching modes. Generally, the measured results are in conformity with the simulated ones, except for a little decrease of the measured radiation patterns as compared to the simulated results. For the different operating modes, almost constant radiation patterns are obtained with maximum gain between 7 dB and 8 dB for the different operating frequencies in the direction of 0° is obtained. It can be observed that the proposed configuration

Fig. 11. Surface Currents Distribution for: (a) Mode M1 Excited at 6.5 GHz, (b) Mode M2 Excited at 5.2 GHz, (c) Mode M3 Excited at 4 GHz, (d) Mode M4 Excited at 3.5 GHz.



Fig. 12. Measured and Simulated Reflection Coefficients of different Operating Modes: (a) Mode M1, (b) Mode M2, (c) Mode M3, (d) Mode M4.





Fig. 13. Example of Prototype Antenna in Anechoic Chamber.

TABLE I.  RADIATION PERFORMANCES OF THE VIVALDI ANTENNA FOR DIFFERENT SWITCHING MODES

| Mode | Resonance frequency (GHz) | Simulated Gain(dB) | Measured Gain (dB) |
|---|---|---|---|
| Mode M1 | 6.5 | 8.07 | 7.53 |
| Mode M2 | 5.2 | 9.54 | 8.41 |
| Mode M3 | 4<br>5.5 | 9.53<br>9.48 | 7.22<br>9.36 |
| Mode M4 | 3.5<br>7 | 7.53<br>9.26 | 7.05<br>8.15 |

## V. CONCLUSION

A frequency reconfigurable Vivaldi antenna with switched band pass resonator has been proposed. The Vivaldi antenna, which exhibits wideband of operation from 2 GHz to 8 GHz, can be reconfigured by perturbing the surface current disruption. So, four switched perturbations are printed on the antenna metallization which can successfully control the surface current flow. Four switching modes are obtained that can offer several narrow bands at 3.5 GHz, 4 GHz, 5.2GHz, .5.5 GHz, 6.5 GHz with satisfactory radiation patterns. In addition, measured and simulation results are in good agreement. Certainly, wideband multi-narrow bands reconfigurable antenna is extremely practical in more recent wireless applications that require dynamic agile frequency such as cognitive radio.

REFERENCES

[1] Rabiaa HERZI, Hsan ZAIRI, Ali GHARSALLAH: Reconfigurable Vivaldi Antenna with improved gain for UWB Applications. Microwave and Optical Technology Letters/ Vol. 58, No. 2, February 2016 DOI 10.1002/mop.

[2] M. Bouslama, A.Gharsallah, M. Traii , and T. A. Denidni: Beam-Switching Antenna with a New Reconfigurable Frequency Selective Surface, IEEE Antennas and Wireless Propagation Letters 15, p 1159-1162, 2016.

[3] A. K. Singh, R.K. Gangwar, B. K. Kanaujia: Wideband and compact slot loaded annular ring microstrip antenna using L-probe proximity-feed for wireless communications. International Journal of Microwave and Wireless Technologies. 2015 doi:10.1017/S1759078715000446.

[4] Ullah S, Ahmad S, Khan BA, Flint JA (2018). A multi-band switchable antenna for Wi-Fi, 3G Advanced, WiMAX, and WLAN wireless applications. International Journal of Microwave and Wireless Technologies 1–7. https://doi.org/10.1017/ S1759078718000776.

[5] S. sharma and C. charu tripathi: Wideband to concurrent tri-band frequency reconfigurable microstrip patch antenna for wireless communication. International Journal of Microwave and Wireless Technologies, 1-8. 2016 doi:10.1017/S1759078716000763.

[6] Cai, Y., Y. J. Guo, and T. S. Bird: A frequency reconfigurable printed Yagi-Uda dipole antenna for cognitive radio applications. IEEE Transactions on Antennas and Propagation, Vol. 60, No. 6, 2905-2912, Jun. 2012.

[7] Kim JY, Ha SJ, Kim D, Lee B, Jung CW: Reconfigurable beam steering antenna using U-slot fabric patch for wrist-wearable applications. J. Electromagn. Waves Appl. 2012; 26: 1545–1553.

[8] C. Yong Rhee, and all. Frequency-reconfigurable antenna for broadband airborne applications. IEEE Antennas and Wireless Propag. Lett. 13 (2014).

[9] YingsongLi, Wenxing Li, andQiubo Ye: ,A reconfigurable wide slot antenna integrated with sirs for UWB/multiband communication applications, Microwave and Optical Technology Letters Volume 55, Issue 1, pages 52–55, January 2013.

Fig. 14.  Radiation Patterns in the E-Plane (Left) and H-Plane (Right) at (a) 6.5 GHz, (b) 5.2 GHz, (c) 4 GHz, (d) 5.5 GHz, (e) 3.5 GHz, and (f) 7 GHz.

[10] Yingsong Li , Wenxing Li, andQiubo Ye: A compact circular slot UWB antenna with multimode reconfigurable band-notched characteristics using resonator and switch techniques, Microwave and Optical Technology Letters Volume 56, Issue 3, pages 570–574, March 2014.

[11] N. Ojaroudi, Y. Ojaroudi, S. Ojaroudi: Compact Ultra-Wideband Monopole Antenna with Enhanced Bandwidth and Dual Band-Stop Properties. International Journal of RF and Microwave Computer-Aided Engineering Vol. 25, No. 4, May 2015.

[12] Kalteh, A. A., G. R. DadashZadeh, M. Naser-Moghadasi, and B. S. Virdee: Ultra-wideband circular slot antenna with reconfigurable notch band function, IET Microwaves, Antennas & Propagation, Vol. 6, No. 1, 108-112, 2012.

[13] Y. Tawk and C. G. Christodoulou, Member, IEEE: A New Reconfigurable Antenna Design for Cognitive Radio, IEEE Antennas and Wireless Propagation Lettres, VOL. 8, 2009.

[14] MertKarahan, Demest S. Armagan Sahinkaya: A Reduced Size Antipodal Vivaldi Antenna Design for Wideband Applications, 2014 IEEE.

[15] Rabiaa HERZI, RamziGharbi, Hsan ZAIRI, Ali GHARSALLAH: A Tuneable Antipodal Vivaldi Antenna for UWB applications, 10th International Multi-Conference on Systems, Signals, and Devices (SSD-13), Hammamet, Tunisia, March 2013.

[16] M. R. Hamid, P. Gardner, P. S. Hall, and F. Ghanem: Switched-band Vivaldi antenna, IEEE Trans. Antennas Propagat., vol. 59, no. 5, pp.1472–1480, May 2011.

[17] M. R. Hamid, P. Gardner, P. S. Hall, and F. Ghanem: Vivaldi Antenna with Integrated Switchable Band Pass Resonator, IEEE Trans. On Antennas and Propagation, Vol. 59, No. 11, November 2011.

[18] T.L. Yim, S.K.A. Rahim and R. Dewan: Reconfigurable wideband and narrowband tapered slot Vivaldi antenna with ring slot pairs. Journal of Electromagnetic Waves and Applications, 2013 Vol. 27, No. 3, 276–287.

[19] F. Meng, S. K. Sharma, B. Babakhani: A Wideband Frequency Agile Fork-Shaped Microstrip Patch Antenna with Nearly Invariant Radiation Patterns. International Journal of RF and Microwave Computer-Aided Engineering/Vol. , No., 2016.

# Accuracy Performance Degradation in Image Classification Models due to Concept Drift

Manzoor Ahmed Hashmani[1], Syed Muslim Jameel[2], Hitham Alhussain[3], Mobashar Rehman[4], Arif Budiman[5]

Department of Computer and Information Sciences, Universiti Teknologi PETRONAS, Sri Iskandar, Malaysia[1, 2, 3]
Universiti Tunku Abdul Rahman, Kampar, Malaysia[4]
University of Indonesia, Jakarta, Indonesia[5]

*Abstract*—Big data is playing a significant role in the current computing revolution. Industries and organizations are utilizing their insights for Business Intelligence by using Deep Learning Networks (DLN). However, dynamic characteristics of BD introduce many critical issues for DLN; Concept Drift (CD) is one of them. CD issue appears frequently in Online Supervised Learning environments in which data trends change over time. The problem may even worsen in a BD environment due to the veracity and variability factors. The CD issue may render the DLN inapplicable by degrading the accuracy of classification results in DLN which is a very serious issue that needs to be addressed. Therefore, these DLN need to quickly adapt to changes for maintaining the accuracy level of the results. To overcome classification accuracy, we need some dynamical changes in the existing DLN. Therefore, in this paper, we examine some of the existing Shallow Learning and Deep Learning models and their behavior before and after the Concept Drift (in experiment 1) and validate the pre-trained Deep Learning network (ResNet-50). In future work, this experiment will examine the most recent pre-trained DLN (Alex Net, VGG16, VGG19) and identify their suitability to overcome Concept Drift using fine-tuning and transfer learning approaches.

*Keywords—Pre-trained networks; deep learning; concept drift; fine tuning; transfer learning*

## I. INTRODUCTION

Machine Learning (ML) is being actively investigated by researchers for the last few decades. However, the need for utilizing these ML models in production has introduced new issues and problems for the research community. Interestingly, the term Industrial Revolution 4.0 has not only changed the traditional business process, but also the traditional Machine Learning approach [1]. The essence of these new approaches is to harness the power of Big Data. Among other types of Big Data (text, audio, sensor data, images, etc.), images are significant. Images are being utilized for prediction and classification in health, education and industries (automobiles, agriculture, etc.) employing the Image Classification Models (SVM, ELM, CNN, OSELM, ACNNELM, etc.) [2]. Industrial Revolution 4.0 demands more robust and scalable Image Classification models. IR4.0 related applications need to classify input data stream in real time. Therefore, instead of batch processing (offline training), these models must learn online learning scenarios (learning with classification). However, in Online Learning scenarios, the statistical properties of images may vary at different time steps which

substantially decrease the performance in terms of accuracy in Image Classification Models. This phenomenon is also known as Concept Drift issue.

## II. RELATED WORK

In literature, several studies discussed and proposed solutions for the Concept Drift issue in online machine learning scenario that handled CD in the non-imaging dataset. However, handling of Concept Drift in Image Classification models is rarely reported. Shallow ML classification models (e.g., Extreme Learning Machine (ELM), Support Vector Machine (SMV), Multi-Layer Perception Neural Network (MLP NN), Hidden Markov Model, etc.) handle classification and regression problems efficiently in structured data [3] and are not feasible to handle the large Image datasets [4]. However, Deep Learning algorithms are a better suited to handle complex data streams and extract value with higher accuracy as compared to the conventional approaches. However, the issue of CD can be handled in online learning through new data adaption models/ modes/ concepts. Therefore, truly autonomous, self-maintaining and adaptive Image Classification models are needed. Some studies focused on CD and adaptation of systems during online sequential environment and urged researchers to further investigate in this direction [5].

A recent study [1] presented a comprehensive view of all existing Image Classification models to handle the Concept Drift issue. This study summarizes the available literature in the following categories:

- Non-Adaptive Shallow Learning Approaches
- Adaptive Shallow Learning Approaches
- Non-Adaptive Deep Learning Approaches
- Adaptive Deep Learning Approaches
- Adaptive Hybrid Deep Learning Approaches

Moreover, this study also proposes a Fully Adaptive Image Classification approach using Meta-cognitive principles by using an ensemble classifier approach. However, the selection of an optimized classifier is very critical. Therefore, in this paper, we have conducted several experiments on several Machine Learning models (including pre-trained Deep Learning models) to validate their better performance in the CD scenario.

## III. RESEARCH METHODOLOGY

This research study examines the multiple Concept Drift scenario (Real, Virtual and Hybrid) and its behavior with several Shallow Learning (SVM, ELM, OSELM) and Deep Learning models (CNN, ResNet-50) in online image stream scenarios. MNIST and CIFAR 10 are considered as benchmark datasets [6] for grey scale (2 channels) and RGB (3 channels) images [7]. In this study we carried out two different experiments to understand the behavior of Machine Learning models due to Concept Drift. They are defined below:

Experiment 01: To identify the effect of Concept Drift on Shallow and Deep Learning Models.

Experiment 02: To mitigate the adverse effect of Concept Drift through the pre-trained network.

Experiment 01 is dependent on two major parts. Firstly, the Machine Learning (ML) models (SVM, ELM, OSELM, and CNN) were trained and tested using MNIST dataset (to illustrate the normal behavior of those Machine Learning models for greyscale images). After that, we simulated the Real, Virtual and Hybrid CD and analyzed their effects on SVM, ELM, OSELM, and CNN. Secondly, we did the same experiment for CIFAR 10 dataset to analyze the CD impact on ML due to RGB images.

In Experiment 02, we used the pre-trained network ResNet-50. Firstly, the ResNet-50 was evaluated by using several natural images. The testing accuracy of these tests was recorded. However, we used CNN for feature extraction and SVM for the classification task. Three categories were used (Airplane, Ferry and Laptop among 101 categories) for simplicity (airplane =800 images, ferry=67 images, laptop=80). These classes were balanced to 67 each because they contained the class imbalance problem (check what if we don't balance these classes). We just retrained the last layer of ResNet-50 with additional three categories (airplane, ferry, and laptop). Note that these layers were trained using 1000 classes from ImageNet dataset. We divided the dataset into 30% training and 70% validation (through randomize technique). Caltex 10 image pixel value is 300x300x3 whereas ResNet-50 takes an input value of 224x224x3. Therefore, we used augmented image datastore (it adjusts image size as per input size and gray-scale into RGB). We can get the deep layer features by using activation (it is better to extract features through the before layer of classifier layer). A minibatch of size 32 was used for getting an optimized GPU utilization. It is better to arrange activation output as a column to obtain more/higher GPU optimization. We used Stochastic Gradient Descent (SGD) for cost optimization (a vector was extracted from CNN for this feature), the SVM classifier (fitccoc in Matlab) was trained by using CNN features. The mean accuracy was evaluated and displayed by using the Confusion Matrix. Each layer of CNN makes some contribution to the input image (by applying/updating weights or adjusting the activation function). We can see by visualizing the network filter`s weight.

### A. Dataset

In Machine Learning, MNIST is recognized as the benchmark dataset for greyscale images [6][7][8]. The MNIST dataset contained 70,000 handwritten images (28x28 pixel)

with 10 target classes [9]. However, for color images, CIFAR 10 [10] and CALTECH 101 [11] datasets are considered as benchmark for classification problem. Caltech 101, one of the most widely cited and used image datasets, was collected by Fei-Fei Li, Marco Andreetto, and Marc 'Aurelio Ranzato.

### B. Models

In Experiment 01, this study analyzed the performance accuracy for Support Vector Machine (SVM) [12], Extreme Learning Machine (ELM) [13], Online Sequential Extreme Learning Machine (OSELM) [14] and Convolutional Neural Network (CNN) [15]. However, Experiment 02 used network ResNet-50 [16] to investigate the behavior of pre-trained networks before and after CD scenario. In the ResNet-50 networks, initial layers are based on CNN structure, which is used to extract its own feature unlike SVM (which extracted hand-crafted features i.e. SURF, HOG, and SPARSE). The initial layers of CNN extract the basic image feature i.e. blob, edges, etc. ResNet-50 model is already trained from 1 million images of 1,000 classes using ImageNet [17] dataset (e.g., AlexNet, GoogleNet, VGG-16, and VGG-19). It uses CNN part for image feature extraction (It requires input image size 224x224x3). However, we used Support Vector Machine (SVM) as classification layer. We trained SVM classifier on three classes (discussed in methodology) from CALTECH 101 dataset.

### C. Software and Hardware

Both experiments were conducted on a single node with parallel processing graphic card (G-FORCE NVIDIA GPU TITAN XP) containing 3748 cores with 32 GB RAM. MATLAB Statistical Machine Learning Toolbox, Deep Learning Toolbox, and ResNet-50 API were utilized to perform these experiments.

## IV. RESULTS AND DISCUSSION

In this study, we conducted two different experiments to validate the performance degradation (in terms of accuracy) due to Concept Drift in SVM, ELM, OSELM, and CNN using Image Dataset.

Initially, ELM, SVM, CNN, and OSELM were trained by using MNIST (784 input pixels and 10 classes) dataset. Later, these trained models validated performance accuracy by testing dataset (a subset of MNIST). The results in Table 1 show that the accuracy in a normal condition is acceptable. The CNN, however, performs well among all other existing models. To simulate the Virtual Concept Drift, an additional feature Histogram of Gradient (HoG) was added to the testing dataset. After testing, the trained models (SVM, ELM, CNN and OSELM) degraded their classification accuracy by almost more than 20% (as shown in Tables 1 and 2).

To simulate the Real Concept drift scenarios, three different data streams (with a change in class boundary) were created. The existing trained models (SVM, ELM, CNN and OSELM) were tested using these three Data Streams (Data1, Data2 and Data3) (shown in Table 3). After testing from Data1, the performance accuracy did not change significantly (because Data1 is the same data stream on which these models were

trained). However, the classification accuracy decreased substantially after Data2 and Data3 (shown in Table 4).

Hybrid Concept Drift is the condition in which both Virtual and Real drifts take place simultaneously. Therefore, to simulate the Hybrid Drift, the properties both Virtual and Real drifts were added to the testing dataset (as shown in Table 5). The already trained models were tested before and after the Hybrid Drift conditions. The results indicated that, after Hybrid Drift, the decrease in the accuracy performance was more than the accuracy degradation in Virtual Drift and Real Drift (shown in Table 6).

Next experiment was conducted to analyze the behavior of SVM, ELM, CNN and OSELM models for 3 channels (RGB, color) images. In this experiment, only hybrid drift was simulated in this experiment because it contained the properties of both types of Concept Drift. Initially, the models were trained using Data1 (which is only RG channels). In order to maintain the input image size, we added 1 channel of Red and 2 channels of Green. The testing results were very nominal after training from less amount of dataset (CIFAR 10) (shown in Table 8). Thereafter, additional 1 channel (blue) was added to simulate the Hybrid Drift in testing data (shown in Table 7). After the occurrence of drift scenario, the accuracy performance of all models became worse to such an extent that even the accuracy of the CNN model was also degraded by almost 50%.

TABLE I. MNIST AND MNIST-HoG DATASETS TO SIMULATE VIRTUAL CD SCENARIO

| Dataset | Input | Output |
|---|---|---|
| MNIST | 784 | 10 classes (0-9) |
| MNIST-HoG | 784x81=865 | 10 classes (0-9) |

TABLE II. TESTING ACCURACY AFTER VIRTUAL DRIFT

| Models | Accuracy before CD (MNIST) | Accuracy after CD (MNIST-HoG) |
|---|---|---|
| | | |
| ELM | 92.7% | 70.2% |
| SVM | 89.0% | 68.8% |
| CNN | 99.7% | 78.4% |
| OSELM | 95.2% | 73.2% |

TABLE III. TMNIST DATA STREAMS TO SIMULATE REAL CD SCENARIO

| Dataset | Input | Output |
|---|---|---|
| MNIST (Data1) | 784 | 10 (0-9) |
| MNIST (Data2) | 784 | 20 (A-J, 0-9) |
| MNIST (Data3) | 784 | 4 (6-9) |

TABLE IV. TESTING ACCURACY BEFORE AND AFTER REAL DRIFT

| Models | Before CD (Data1) | After CD (Data2) | After CD (Data) |
|---|---|---|---|
| ELM | 92.12% | 65.76% | 69.56 |
| SVM | 88.67% | 70.68% | 68.78 |
| CNN | 99.76% | 76.76% | 73.45 |
| OSELM | 94.56% | 72.34% | 69.67 |

TABLE V. MNIST DATA STREAMS TO SIMULATE HYBRID CD SCENARIO

| Dataset | Input | Output |
|---|---|---|
| MNIST (Data1) | 784 | 10 (0-9) |
| MNIST (Data2) | 784+81=865 | 4 (6-9) |

TABLE VI. TESTING ACCURACY BEFORE AND AFTER HYBRID DRIFT

| Models | Before CD (Data1) | After CD (Data2) |
|---|---|---|
| ELM | 92 | 60.32 |
| SVM | 89 | 59.46 |
| CNN | 100 | 72.87 |
| OSELM | 95 | 69.59 |

TABLE VII. CIFAR 10 DATA STREAMS TO SIMULATE HYBRID CD SCENARIO

| Dataset | Input | Output |
|---|---|---|
| CIFAR-10 (Data1) | 32x32x3 (3072) {RG+G} | 10 |
| CIFAR-10 (Data2) | 32x32x (2+1) = (3072) {RG + B} | 5 |

TABLE VIII. TESTING ACCURACY BEFORE AND AFTER HYBRID DRIFT

| Models | Before CD (Data1) | After CD (Data2) |
|---|---|---|
| ELM | 70.49 | 25.26 |
| SVM | 77.68 | 34.75 |
| CNN | 97.8 | 49.50 |
| OSELM | 91.28 | 39.67 |

In experiment 2, a pre-trained network (ResNet-50) was validated in the Real Concept Drift condition. However, it can be noticed that the accuracy performance was more than all other existing models (shown in Table 9).

TABLE IX. MEAN ACCURACY OF RESNET-50 WITH SVM CLASSIFIER

| Models | Dataset | Mean Accuracy Before CD | Mean Accuracy After CD |
|---|---|---|---|
| ResNet-50 | Caltech 101 | 98.58 | 76.81 |

## V. CONCLUSION

It can be safely concluded from the results of this experimental study, that the Image Classification model degrades its performance (in terms of accuracy) due to Concept Drift. However, the Hybrid drift causes more accuracy degradation then Real or Virtual Drift. Moreover, the complexity of image dataset is also directly proportional to the accuracy degradation after Concept Drift (e.g., MNIST image's accuracy degradation is less than CIFAR 10). Interestingly, the CNN model showed higher accuracy than others in most of the experiments because it has extracted and not hand-crafted features (because it is extracted not hand-crafted features). However, the accuracy of pre-trained network (ResNet-50) is better than that of CNN. Fundamentally, ResNet-50 also uses CNN for feature extraction, whereas its feature extraction layers are already trained using ImageNet dataset (containing 1 million images of 1000 classes). Therefore, through these experiments, we can conclude that the pre-trained network

offers a better solution for handling Concept Drift at the classifier level. Nonetheless, there is a need of developing a dynamic adaptation approach (which will work along with classifier) to adapt new features of Image Data stream (in online learning scenario). In future work, this study will be extended towards the training and validation of pre-trained Deep Learning models (e.g., ResNet-50, AlexNet, VGG16, VGG19) in certain Concept Drift scenarios. For that reason, multiple CD scenarios will be added explicitly (by adding the new features and classes in input image dataset) to investigate the performance accuracy of other Deep Learning networks in CD environment. This experiment will help in understanding and analyzing the performance of pre-trained Deep Learning networks in a variety of CD scenarios.

### REFERENCES

[1] Jameel, Syed Muslim, et al. "A Fully Adaptive Image Classification Approach for Industrial Revolution 4.0." International Conference of Reliable Information and Communication Technology. Springer, Cham, 2018.

[2] Najafabadi, M M,. Villanustre, F,. Khoshgoftaar, T M,. Seliya, N,. Wald, R,. Muharemagic, E. "Deep Learning Applications and Challenges in Big Data Analytics", Journal of Big Data 2(1), 1 (2015).

[3] Kuncheva L.I. "Classifier Ensembles for Changing Environments". In: Roli F., Kittler J., Windeatt T. (eds) Multiple Classifier Systems. MCS 2004. LNCS, vol. 3077, pp. 1-15. Springer, Berlin, Heidelberg (2004).

[4] Kuncheva L.I. "Classifier Ensembles for Detecting Concept Change in Streaming Data". Overview and Perspectives. 2nd Workshop SUEMA, pp. 5-10 (2008).

[5] Zliobaite, I., Bifet, A,. Pechenizkiy, M,. Bouchachia, A.. "A Survey on Concept Drift Adaptation". ACM Computer Survey 46(4), pp. 1-37 (2014).

[6] Budiman, A., Fanany, M I., Basaruddin C.: Adaptive Online Sequential ELM for Concept Drift Tackling. Computational Intelligence and Neuroscience 2016 (20), 2016.

[7] Corrigan, Owen. An Investigation Into Machine Learning Solutions Involving Time Series Across Different Problem Domains. Diss. Dublin City University, 2018.

[8] Papernot, Nicolas. "Adversarial Examples in Machine Learning." (2017).

[9] Y. LeCun, C. Cortes. MNIST handwritten digit database [online] (2010).

[10] Snoek, Jasper, Hugo Larochelle, and Ryan P. Adams. "Practical bayesian optimization of machine learning algorithms." Advances in neural information processing systems. 2012.

[11] Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. "Imagenet classification with deep convolutional neural networks." Advances in neural information processing systems. 2012.

[12] Cortes, Corinna, and Vladimir Vapnik. "Support-vector networks." Machine learning 20.3 (1995): 273-297.

[13] Huang, Guang-Bin; Zhu, Qin-Yu; Siew, Chee-Kheong (2006). "Extreme learning machine: theory and applications". Neurocomputing. 70 (1): 489–501.

[14] G.-B. Huang, N.-Y. Liang, H.-J. Rong, P. Saratchandran, N. Sundararajan, On-line sequential extreme learning machine; M.H. Hamza (Ed.), IASTED International Conference on Computational Intelligence, IASTED/ACTA Press (2005), pp. 232-237

[15] LeCun, Yann. "LeNet-5, convolutional neural networks". Retrieved 16 November 2013.

[16] He, Kaiming, et al. "Deep residual learning for image recognition." Proceedings of the IEEE conference on computer vision and pattern recognition. 2016.

[17] Russakovsky, O., Deng, J., Su, H. et al. Int J Comput Vis (2015) 115: 211. https://doi.org/10.1007/s11263-015-0816-y.

# Study and Design of a Magnetic Levitator System

Brian Meneses-Claudio[1], Avid Roman-Gonzalez[3]
Image Processing Research Laboratory (INTI-Lab)
Universidad de Ciencias y
Humanidades
Lima, Perú

Zeila Torres Santos[2]
Interdisciplinary Research Center Science and Society
(CIICS), Universidad de Ciencias y
Humanidades
Lima, Perú

*Abstract*—**Magnetic levitation is one of the mechanisms that is at the forefront of technology. It is used in its most basic form in educational teaching, where the principles of physics converge that have as their principle electromagnetism and the fields created by existing poles that repel according to a quantity of initial current, giving instructive ideas of how the theoretical formulas work, giving life to a practical visual system. The current use on a large scale are the Maglev trains of Japan or superconductivity, being the realization of the quantum effects visualized at the moment of cooling the sample. The electronic circuit tends to be stable because, when using a high-power current, a Triac is needed to compensate the electrical flow provided by the operational amplifier and, therefore, stabilize with the photodiode when activated with the Led diode. Our purpose is to create a circuit that identifies the values of the electronic components that allow reaching equilibrium, with input and output variables that indicate the position and height of the object to be levitated.**

*Keywords*—*Magnetic levitator; electromagnet; electronic circuit; differential potential*

## I. INTRODUCTION

The magnetic and gravitational force, will allow us to understand the Levitation phenomenon of great importance and visually enrich.

The Magnetic Levitator [1] teaches us not only how magnetic fields are generated [2] and the repulsion force that the system must have in order to avoid the force of gravity that all bodies are affected. This is the starting point to understand how a body could levitate. For this purpose, will be necessary to create a functional circuit with electronic devices and coils, when the electric current passes, allows us to see the object suspended. An example of utility is the super-fast public transport trains in Japan.

At the beginning of the seventeenth century, researchers such as Volta, W. Gilbert and others-maintained ideas that Electricity and Magnetic were understood separately, until at least the nineteenth century that Oersted showed a connection between the magnetic field and the electric field using a compass and a conductive thread. Then, in 1861 with Maxwell, he joined the concepts with specific formulas that were fundamental to enrich that area, concepts such as magnetic induction, magnetic permeability, magnetic field, automatic inductance and more. The implementation of the formulation to the construction of the levitator came together with the theory of the Electromagnetism and depending on the type of levitator that you want to work come together

mathematical equations that help to discover with accuracy the physical quantities of voltages, sensitivities, electronic units and other useful components in the process [3].

Levitate an object is possible, and there are variants within the electromagnetic theory to perform this process. Already by itself only an electromagnet and its ferromagnetic core forms the basis for levitating an object, with an array of polarities of the metallic object. A current flow through the wound coils of the electromagnet generating a magnetic field, being that the ferromagnetic core provides a path of reluctance in which the magnetic field is concentrated which induces an attractive force on the object to be levitated. Now, if we understand the theory of Superconductivity [4], we know that cooling an object at sufficiently low temperatures causes a distance to be separated from its initial base by positioning itself at a height that we can call a gap, until the temperature again decays. There is levitation by repulsion and suspension, in both cases a study of field forces must be performed, because positioning the object at the midpoint that will be the equilibrium-stable point to levitate, will depend a lot on the circuit, the coils and the distance between them, how to calibrate until the objective is achieved. When we speak of coils through them, an electric field flows generating a force field, requiring an adjustment in the flowing current.

The power electronics has a very large study area, the most important being the control of the current, because of this TRIACS are used in addition to resistors for the opposition to the current [5].

Photodiodes are used for the activation and opening of systems by means of the excitation of light being directly proportional, meaning that when the photodiode is excited with the intensity of light, its impedance will decrease increasing the circulation of current [6].

The following research work is structured as follows: In Section II, the methodology for the design of a magnetic levitator circuit will be presented, for which, equations are proposed on the magnetic force for the estimation of electronic components required for magnetic levitator also shows the design of the simulated electronic circuit in the Proteus 8 software. In Section III, the faults that were obtained at the time of the simulation are briefly explained indicating that no current losses were obtained. In Section IV, the discussion of the research work is presented, as well as the problems that were obtained for the design of the magnetic levitator circuit. Finally, in Section V, the conclusions of the research work are presented.

## II. METHODOLOGY

### A. Magnetic Levitator

Since magnetism is a physical property of response to the applied magnetic field, it is measured through the Lorentz Force, given by the following equation:

$$\vec{E} = \frac{\mathbb{F}}{q}$$

What produces the magnetic field are the electric currents or electric field $\vec{E}$ that varies with time, being the charge $q$ in Coulomb that moves with a given electric force $\mathbb{F}$ in Newton.

Therefore, if we join these two fields being one dependent on the other, we obtain the electromagnetic force by the following equation.

$$\vec{F} = q\,\vec{E} + q\,\vec{v} \times \vec{B}$$

Where $\vec{B}$ it is the magnetic field vector.

As the equations begin to arise, the system could be linear, non-linear, MISO (multiplexer), SIMO (demultiplexer), SISO (single input and single output), variant, invariant or depending on the inputs and outputs of the system, forming a set of state equations that describe the dynamics of the system, which in a general way describe the movement and functionalities in conjunction with the variables that make up the system [7]:

$$\dot{x}(t) = f(x, u, t) \tag{1}$$

$$\ddot{x}(t) = h(x, u, t)$$

These systems equation is depending on the position ($x$), current-voltage ($u$) and time ($t$). A suspension levitator requires sensors, controllers, coils and the passage of current through a circuit structured specifically for an object to be suspended at a certain distance and weight. In Fig. 1, the design that gives us an idea of the forces that occur is shown.

The magnetic force equation that is derived from [8]:

$$\vec{f}_{ml}(\vec{\imath}, x) = \frac{i^2}{2}\frac{dL}{dx} \tag{2}$$

Is assumed the dependence equation: $L = L_1 + \frac{L_0 x_0}{2}$ put this into the Equation (2) we obtain:

$$\vec{f}_{ml}(\vec{\imath}, x) = \frac{i^2}{2}\frac{d\left(L_1 + \frac{L_0 x_0}{2}\right)}{dx} \tag{3}$$

$$\vec{f}_{ml}(\vec{\imath}, x) = \frac{L_0 x_0}{2}\left(\frac{i}{x}\right)^2$$



Fig. 1. Forces that are Applied in a Magnetic Levitation.

From the last equation we have the magnetic force constant $k = \frac{L_0 x_0}{2}$

Being $i$ the current in the levitation coils, which control the force $\vec{f}(\vec{\imath}, x)$ and $x$ the gap-height or the distance of separation. Adding forces, following Newton's second law and considering the suspended object along $x$ axis, we have:

$$\sum \vec{F} = m\vec{a} = m\ddot{y} = \vec{f}_{ml}(\vec{\imath}, x) - \vec{F_g} = 0$$

$$m\ddot{y} = k\left(\frac{i}{x}\right)^2 - mg \tag{4}$$

Where $i$ the current in the coil of the electromagnet is, $g$ is the gravitational constant, $k$ is the magnetic force constant and $m$ is the mass of the object.

The levitator consists of a magnet that create a magnetic field and electromagnets that control that magnetic field. It is necessary to vary the value of the electromagnetic force by adjusting the current that passes through the electromagnet because the electromagnet is responsible for generating the electromagnetic force that allows the levitation of the object, as it shown in Fig. 2.

Each magnet has two poles: the north and the south. Experiments show that opposites are attracted and the same poles repel each other. Four cylindrical magnets are placed in a square and have the same polarity, around forming a magnetic field up to push any magnet, which has the same pole and in the middle of them. Together with four levitation coils, placed equidistantly and symmetrical magnets, it is possible to create an opposite magnetic field.

Sets of dynamic equations (Newton's and Kirchhoff's equations) that provide theory to the system are given by the following set of equations using equation (4) we have the equation of mechanical system and the differential equation of the circuit.

$$m\frac{d^2 x}{dt^2} = mg - k\left(\frac{i}{x}\right)^2 \tag{5}$$

$$L\frac{di}{dt} = e(t) - R\,i \tag{6}$$

Where $L$ is the inductance of the electromagnet that de, depends on the position of the object and a control variable, because is to assume that $L$ varies inversely with respect to the object´s position $x$.



Fig. 2. Magnetic Suspension Levitator.

The infrared photo-emitter and photo-receiver pair are the sensors used, which indicate that the amount of light received increases, which occurs when the metallic object moves away from the electromagnet, the output of the photo-receiver increases and vice versa. The equation that describes the sensor output $v_s$ is:

$$v_s = k_s h + \bar{v}_s \tag{7}$$

where $k_s > 0$ indicates the gain of the sensor and $\bar{v}_s$ is constant and $h$ is the gap (Fig. 1).

The voltage applied to the electromagnet having the same linear form:

$$i = k_v v + \bar{\imath} \tag{8}$$

where $\bar{\imath}$ is the current for the sensor output to be $\bar{v}_s$.

The controller (photodiode) uses the sensor signal to adjust the voltage applied to the actuator and thus maintain levitation. For MISO systems or for each dynamic systems, where the knowledge of the enter states variables determine the behavior of the system in equation (1) see the set of state equations, linearizing we have:

$$\dot{x}_1(t) = A(t)x(t) + B(t)e(t) \tag{9}$$

$$\dot{x}_2(t) = C(t)x(t) + D(t)e(t) \tag{10}$$

Constructed the state equation we use the follow system:

$$x_1 = h, \quad x_2 = \frac{dh}{dx}, \quad x_3 = i \tag{11}$$

We re write equation (5-6):

$$\dot{x} = x_1, \quad \dot{x}_1 = g - \frac{k}{m}\left(\frac{x_2}{x_1}\right)^2, \quad \dot{x}_2 = \frac{e}{L} - \frac{R}{L}x_2 \tag{12}$$

We find the linearized state equations at the equilibrium point for temporary invariant systems, in the state space, where the variables $\tilde{x}$ y $\tilde{e}$ are the variables subtracting the equilibrium positions.

$$\tilde{x} = \mathbb{A}\tilde{x} + \mathbb{B}\tilde{e} \quad , \mathbb{A} = \begin{bmatrix} 0 & 1 \\ \frac{2g}{h_0} & 0 \end{bmatrix} , \quad \mathbb{B}\begin{bmatrix} 0 \\ \frac{-2g.k_v}{i_0} \end{bmatrix} \tag{13}$$

Where $\mathbb{A}$ and $\mathbb{B}$ are called matrix state.

## B. Electronic Circuit Design of the Magnetic Levitator

For the design of the circuit, the following components shown in Table I was needed.

The Operational Amplifier is in the position of an inverter subtractor because in the circuit a high impedance is required to pass the current to stabilize it and maintain its constant flow, they are used more in devices with very weak signals, but where it requires a constant current flow. In addition, it requires 2 voltage sources, one in positive and one in negative for its operation, both with the same current.

$$V_{OUT} = V_2\left(\frac{(R3+R1)R4}{(R4+R2)R1}\right) - V_1\left(\frac{R3}{R1}\right) \tag{14}$$

Where:

- R1 = 20 kΩ
- R2 = 20 kΩ

- R3 = 4.7 kΩ
- R4 = Photodiode (Resistance converter based on the received light intensity), it is considered null.
- V1 =12 v & V2 = -12V

$$V_{OUT} = 2.079\ V$$

The output voltage of the Operational Amplifier is 2.079 V, only and when the photodiode is activated, it will be supplied to the TRIAC to send it to the coil.

The TRIAC or Triode [9], is a switch capable of switching alternating and continuous current, it is used when the control of the current flow in a circuit is required as shown in Fig. 3, and it also requires an activation voltage, being 2,079V obtained from the previous analysis. In this case, it will be used as an actuator controlled by the photodiode and the LED, will identify the presence of an object and activate the current flow for the creation of the magnetic field.

The coil is where current will flow, it is also where a The main function of the coil is to keep the current stable until it discharges on its own, meaning that, if there is a sudden change in the current, the coil will dissipate the magnetic field for periods of time until finally extinguished. For this reason, a diode is connected in antiparallel with the coil, in this way and thanks to the coil inverting the direction of the current, it will circulate through the diode. Thanks to this flow, the inductor will be discharged in a controlled manner and will not appear over voltages. The circuit is shown in Fig. 4 using all the electronic devices.

Finally, to find the power in direct current the following formula will be used:

$$P(t) = R \cdot I^2 = \frac{V^2}{R} \tag{15}$$

Which R is the equivalent resistance of the circuit and I is the controlled input current of the circuit.

TABLE I. ELECTRONIC DEVICES

| | |
|---|---|
| Resistors | - 1 x 270 Ω<br>- 2 x 20 kΩ<br>- 1 x 4.7 kΩ<br>- 1 x 5.6 kΩ<br>- 1 x 1 kΩ |
| Operational Amplifier | - LM358 |
| Led Diode | - 2 x Diode Led |
| TRIAC o Triode | - IRF540 |
| Electric Capacitor | - 1000 μf<br>- 4.7 μf |
| Diodes | - PHOTODIODE<br>- 1N5401 |
| Coil | - Coil 3 mH |
| Voltage Regulator | - 12 volts y 2 amps<br>- -12 volts y 2 amps |

Fig. 3. Operation of the TRIAC Against the Current.



Fig. 4. Design of the Magnetic Levitator Circuit.

## III. RESULTS

The electronic circuit of the magnetic levitator does not have loss of current due to the filters applied in the circuit to control the alternating current. In addition to establish the magnetic body within the range of the photodiode, providing stability and necessary use of current, the system is activated.

The simulation of the circuit was done through the software Proteus 8, in which it does not provide connection or current distribution errors, besides controlling the high level of current with power circuits.

## IV. DISCUSSION

The research work confirms the use of electronic components without the need for a programmable microcontroller, the main purpose is to use electronic components based on their main function and identify the voltage and current values that run the circuit to know the final result for the power in the coil.

In electronic magnetic levitator circuit, it was considered the current level to get the optimum electromagnetic field in the coil, therefore, an OPAM (Operational Amplifier) was used to control the current and maintain the electric flow constant. In addition, the TRIAC is required front a very high current to preserve the level of current and take care of the other important electronic devices such as the coil. In case, it does not have these electronic devices implemented in the circuit, it cannot maintain the flow and stability of the current.

## V. CONCLUSIONS

As a future work, we want to implement the electronic circuit and also add a magnetic levitator by repulsion to identify the sum of the forces exerted on an object when subjected by 2 electromagnetic fields.

The use of the TRIAC is based on the level of operating current of the circuit, compared to a transistor because as a maximum it can work at 1 ampere, thus limiting the desired power in the coil to generate the electromagnetic field.

REFERENCES

[1] Nayak A., Controller design for magnetic levitation system, Thesis for Master, Department of Electrical Engineering NIT, India, April 2015.

[2] Giancoli, Douglas C.; Physics Principles with Applications, Person Education Limited, England, Chapter 16-21, 2016.

[3] M. Ahmed, M. F. Hossen, M. E. Hoque, O. Farrok, and M. Mynuddin; "Design and Construction of a Magnetic Levitation System Using Programmable Logic Controller," Am. J. Mech. Eng. Vol. 4, 2016, Pages 99-107, vol. 4, no. 3, pp. 99–107, May 2016.

[4] F. Grilli, A. Morandi, F. De Silvestri, and R. Brambilla; "Dynamic modeling of levitation of a superconducting bulk by coupled H-magnetic field and Arbitrary Lagrangian-Eulerian formulations," Italy, 2018.

[5] J. A. Guijarro Solórzano and J. T. Vivar Martínez; "Diseño e Implementación de un Levitador Electro-Magnético basado en un Control PID utilizando LABVIEW," Universidad Politécnica Salesiana, Ecuador, 2015.

[6] S. Wojtczuk, X. Zhang, and W. MacNeish; "Enhanced visible near-infrared photodiode and non-invasive physiological sensor," 2016.

[7] Castro-Beltrán J. S., Vergara Ramírez C. F, Herrera Guayazan J S; Comparison of a linear and nonlinear control on magnetic levitator; Ingenierías USBMed, vol. 9 (1), pag. 112-118, 2018.

[8] Silviu Folea, Cristina Muresan, Robin De Keyser, Clara Ionescu; Theoretical Analysis and Experimental Validation of a Simplified Fractional Order Controller for a Magnetic Levitation System, IEEE Transactions on control systems technology, 2015.

[9] S. Bárez-López, M. J. Obregón, R. Martínez de Mena, J. Bernal, A. Guadaño-Ferraz, and B. Morte; "Effect of TRIAC in the treatment of Mct8," MCT8 Symp., 2016.

# Detection of Suspicions of Varicose Veins in the Legs using Thermal Imaging

Brian Meneses-Claudio[1], Witman Alvarado-Diaz[2]
Image Processing Research Laboratory (INTI-Lab
Universidad de Ciencias y
Humanidades Lima, Perú

Avid Roman-Gonzalez[3]
Senior Member, IEEE Image Processing Research
Laboratory (INTI-Lab Universidad de Ciencias y
Humanidades Lima, Perú

*Abstract*—Varicose veins also known as venous insufficiency, are dilated veins due to an accumulation of blood that occurs in different parts of the body, the most common are in the legs, in addition to having a higher index in women for clothing style that they use. Varicose veins are classified by grades ranging from I to IV and can cause pain, itching, cramps and even ulcers if they are treated in time. Not all varicose veins can be visible superficially, many of them begin inside of the skin. According to the WHO (World Health Organization) 10% of the world population has varicose veins. That is why the detection of suspicions of varicose veins in the legs was raised in this research work, first a thermal image will be obtained using the FLIR ONE Pro thermal camera following a necessary protocol of distance and temperature range. The thermal image is processed in MATLAB to identify the segments of the histogram of the thermal image, to obtain the area of the highest temperature indicating the presence of varicose vein in the subject's leg. The segmentation of the areas with the highest temperature was obtained as a result to be overlaid on the real image, showing the real image with the varicose vein segment found in the thermal image processing.

*Keywords*—*Thermal image; varicose veins; detection; image processing; image segmentation*

## I. INTRODUCTION

Varicose veins are inflamed veins that can be seen under the surface of the skin, majority appears in the legs because it is where it exerts more force to charge the weight of the torso, but are not the only places where varices can manifest. In addition, they can cause mild pain, blood clots, skin sores and itching [1].

According to the WHO (World Health Organization) , varicose veins are superficial, cylindrical or vascular veins and can be caused by many factors such as: sedentary lifestyle, pregnancy, exposure to heat, overweight, wear tight clothes and shoes, etc. [2]. Its main function is to prevent the return of blood to the heart continuously, so the veins of the body tend to degenerate and more if you live sedentary, the most common areas are the legs[1].

According to the WHO, varicose veins are a very common problem that almost 10% of the world population suffers, the rate is higher in women, in addition, the risk of developing varicose veins increases with age, with 35% of active people and increases between 50 to 60% when it comes to a sedentary lifestyle [3].

Varicose veins are classified by grades from I to IV [4], usually begin as an aesthetic problem showing a thin turquoise blue lines, giving the sensation of itching, heaviness and fatigue, then when going up grade, varicose veins can be appreciated in the surface of the skin with small swellings [5] and finally if they are not treated in time or the damaging factors of varicose veins continue to be applied, they can produce ulcers, internal circulation failures and inflammations of large areas in the leg [2].

The detection of suspected varicose veins in the legs early can help prevent further progress, although it is a medical condition that progresses slowly, is aggravated when they are shown superficially, in addition to feeling itching, cramping, etc.[6]; which are some symptoms of early varices. If in the event an early detection and healing process is not followed, the symptoms are ulcers, inflammation and swelling, bleeding from the veins near the skin and finally sensibility in the legs [7].

The main objective of the research work is the detection of suspicions of varicose veins while they are still not shown superficially to indicate that a healing process is required. It is detected through the segmentation of thermal images captured in the legs.

Thermography is a science of the study of temperature variation means that it can be applied where the problem can be revealed through a thermal difference. Currently, there are thermal cameras that connect to Smartphone and are used for areas such as aviation, medicine, construction, electronics, etc [8].

The thermal images show the thermal composition of a body or object depending on the temperature range to which it has been programmed, these images can be processed with image processing software because they are compatible [9]. In addition, the use of these images is useful for the detection of internal pathologies being not superficially visible, that is why it has been used for the detection of suspicions of varicose veins in the legs.

The following research work is structured as follows: In Section II, the image processing methodology for thermal image segmentation will be presented. In Section III, the

---

[1] https://medlineplus.gov/spanish/ency/article/000203.htm [3]

[2] https://www.aeev.net/varices.php [4]

results will be shown by a thermal image and the segmentation of the varicose vein overlaid on the real image. In Section IV, it will present the discussion of the research work. Finally, in the Section V, it will present the conclusion and a future plan of the research work.

## II. METHODOLOGY

In this part, each part of the segmentation of the thermal image is developed for the detection of suspicions of varicose veins, which consist of the acquisition of the image, image processing, segmentation of the areas with higher temperature and finally the segmented image overlaid on the real image [10].

The study of the images of the thermal camera was also used in the vascular disorder [11], making a segmentation of the grayscale histogram in the MATLAB program, then the bodies of the histogram will identify the hottest areas, then convert the image to the HSV scale to know the intensity of the image and finally show it separately, the blue being the hottest zones.

The stages of the system are shown in Fig. 1 where the processes by which the thermal image will be subjected:

### A. Image Acquisition

For the acquisition of an image, we use the FLIR ONE Pro, it is a camera capable of capturing thermal images, being compatible with IOS and Android mobile devices. It has a thermal sensor that will help to sectorize the temperature levels of the object being pointed. This device has two lenses as shown in Fig. 2, the upper lens captures the images and the lower lens captures the thermal images. In addition, the C to Micro USB connector converter is shown because an Android device was used to capture the thermal images [12].

Table I shows the operating characteristics of the FLIR ONE PRO [13].

The FLIR ONE Pro is a device capable of analyzing and visualizing at a distance the temperature distribution of complete surfaces, it was born as an idea for industrial areas due to the complexity of the detection of equipment failures, currently used in the field of the detection of problems not only of equipment but also of human bodies.

### B. Image Processing

In this stage, the thermal image is in 3 dimensions Fig. 3(b), so, using the MATLAB program, we convert it to gray scale.

The software makes a averaged sum by taking the values of each pixel of the image of each dimension and multiplying with a normed value as shown below:

$$0.2989 * R + 0.5870 * G + 0.114 * B \tag{1}$$

Then, there is the histogram of the image Fig. 3(c), identifying the segments of the thermal image, to observe the histogram, we use the following code:

$$imhist(thermal\_image\_gray)$$



Fig. 1. Process Flow Diagram for the Detection of Suspected Varicose Veins in the Legs.



Fig. 2. FLIR ONE Pro Thermal Camera and Type C to Micro USB Connector Adapter.

TABLE I.        CHARACTERISTICS OF THE FLIR ONE PRO

| FLIR ONE Pro | |
|---|---|
| Temperature Range | -20 °C – 400 °C |
| Compatibility | IOS and Android Devices |
| Distance | Functional Distance l to 1.8 meters |
| Weight | 36.5 g. |
| Dimensions | 68 x 34 x 14 millimeters |
| Thermal Resolution | 160 x 120 |
| Operating Time | 1 hora |

Each segment obtained from the histogram is to identify the peaks that can be considered as identifiable objects in the image. After identifying the segments, they will be overlaid to obtain the hottest zone of the image Fig. 3(d).

Next, finding the edges of the hottest zone, for that case, the Roberts Method was used, this method was chosen because it can easily identify the diagonal edges, besides filling that space Fig. 3(e). Roberts' method uses a filter that focuses on each pixel through the following formula:

$$\frac{df}{dx} = f(x + 1, y) - f(x, y) \qquad (2)$$

$$\frac{df}{dy} = f(x, y + 1) - f(x, y)$$

Where to locate the pixel (x, y) that is in gray scale within the range of 0 to 255, if the areas have a constant intensity, will turn them into 0, thus giving the edges where the varicose vein is drawn.

The next step is to improve this segmentation because when going through the Roberts Method, segmentation will always be obtained with noise. To improve the edges, we will use the Morphological Structure Method, its formula is shown below:

$$\delta_B(X) = X \oplus B = \{x | \ X \cap B_x \neq \emptyset\} \qquad (3)$$

Where it is indicated that X will travel through the whole image, when it passes through B, it will provide information about the data of the neighbors of that pixel, converting it to the maximum value of the environment of that neighborhood defined by the element of the structure [14]. The values that are around each pixel are called neighborhood. Then when the pixel has a neighborhood of values different from it, it will take the maximum value of that neighborhood; thus, improving the edges of the varicose vein.

Finally, fill in the border of a transparent color to overlay the segmented image of the varicose veins and the real image. Fig. 3(f), the whole process is shown in Fig. 3.



Fig. 3.   Varicose Veins Segmentation. (a) Real Image of the Tibial. (b) Thermal Image of the Real Image captured by the FLIR ONE PRO. (c) Histogram of the Gray Scale Thermal Imaging. (d) Union of the Segments of the Image Histogram. (e) Identification of Edges with the Roberts Method. (f) Superpose the Segmented Varicose Vein to the Real Image.

## III. RESULTS

The thermal images were acquired based on a protocol because all the images are required to have the same characteristics and therefore the same processing. In Table II, the characteristics of the thermal images are presented:

Fig. 5. Characteristics of Thermal Mages

| Thermal Images | |
|---|---|
| Distance | 30 cm |
| Temperature Range | 15 °C – 40 °C |
| Connected Device | Moto G Play |
| Place where the Image was captured | -Tibialis Anterior.<br>-Tibialis.<br>-Soleus. |
| Rest Time | 15 to 20 minutes of rest. |



Fig. 4. Thermal Imaging and Segmentation of the Varicose Vein Overlaid on the Real Image.

Not all varicose veins are shown superficially [15], because with the thermal image a better visualization of the involved areas of varicose was obtained. In addition, the difference in temperature between zones without any problem and varicose veins is confirmed. Fig. 4 shows the thermal image and the segmentation of the varicose overlaid on the real image of the user's leg.

In the research work [16], using thermal imaging to find varicose validating the use of thermography in the diagnosis of a venous thrombosis, indicating the area where is seen more red because there is a higher temperature variation in those points. They also indicated that the varicose vein detection process took a period of 3 weeks to find out if varicose veins were seen in that area, thus ruling out recurrent muscle warmth in the legs. Also, they used an ultrasound to reconstruct an image based on temperature, which is why they show thermal images.

## IV. DISCUSSION

The research work confirms the use of thermal images for the detection of suspicion of varicose veins due to the variation of temperature in the legs. In addition, not all varicose veins are visible on the surface of the skin, the majority is always internal.

The device to which the FLIR ONE Pro will connect is not important because it will only use as an image buffer since the FLIR ONE Pro has no internal storage. Therefore, it is only required that it is a Smartphone and has a C or Micro USB input; in addition to having the APP of the FLIR ONE Pro to capture the thermal images.

For the image processing, the real image and the thermal image were required because it was necessary to superimpose the segmented area of the varicose veins to the real image at the end of the process.

The resting time was very variable due to the fact that the legs were on the floor and did not stay static because there was always movement, being indicative of abnormal warming that may arise in the leg that is why the segmentation of the histogram was followed of the image looking for the hottest areas.

## V. CONCLUSIONS

It is concluded that varicose veins can be detected through thermal imaging in an efficient and fast way because only the

images need to be captured and then the software will automatically segment it, obtaining the area where the varicose vein is located in addition to the size of the varicose vein.

It is concluded that a protocol was established because it was adjusted to the size of the user's leg in addition to having the same characteristics for each image and calibrate the FLIR ONE Pro in the required temperature range.

As a future work, we want to measure the depth of the varicose veins because if the varicose vein is very deep, it is an indication that it does not need to be superficial to be serious. In addition, not only to analyze varicose veins, also another type of pathology that can be identified through the temperature difference.

### REFERENCES

[1] F. Deng, Q. Tang, G. Zeng, H. Wu, N. Zhang, and N. Zhong, "Effectiveness of digital infrared thermal imaging in detecting lower extremity deep venous thrombosis," Med. Phys., vol. 42, no. 5, pp. 2242–2248, Apr. 2015.

[2] M. Figueroa Pérez and C. Vergaray Falcón, "Conocimiento Sobre Várices Y Medidas Preventivas En Miembros Inferiores En El Profesional Enfermero De Centro Quirúrgico De Una Clínica Privada, Septiembre," Universidad Privada Cayetano Heredia, 2017.

[3] URAC American Accreditation HealthCare Commission, "Insuficiencia venosa," 2010. [Online]. Available: https://medlineplus.gov/spanish/ency/article/000203.htm. [Accessed: 16-Apr-2019].

[4] "Varices - Area Pacientes - AEEVH," 2012. [Online]. Available: https://www.aeev.net/varices.php. [Accessed: 22-Apr-2019].

[5] D. Paola Ortiz Renata Carvalho and D. J. Leandro Pérez Segura Julián Javier, "Enfermedad Venosa Superficial Crónica de Miembros Inferiores: Epidemiología, Anatomía Y Fisiopatología Enfocada a Latinoamérica," 2015.

[6] S. S. Sami, D. Harman, K. Ragunath, D. Böhning, J. Parkes, and I. N. Guha, "Non-invasive tests for the detection of oesophageal varices in compensated cirrhosis: systematic review and meta-analysis Key summary.

[7] Radiologial of Society North America, "Insuficiencia venosa (venas varicosas o várices)," 2017. [Online]. Available: https://www.radiologyinfo.org/sp/info.cfm?pg=varicose-veins. [Accessed: 26-Apr-2019].

[8] C. P. Naranjo Eraso and P. A. Vásquez Suárez, "Diagnóstico termográfico preventivo para lesiones músculo esqueléticas más comunes en futbolistas.," Universidad Politécnica Salesiana, Ecuador, 2019.

[9] S. Verdaguer D and J. C. Gana A, "Enfrentamiento de pacientes pediátricos con várices esofágicas," Rev. Med. Chil., vol. 144, no. 7, pp. 879–885, Jul. 2016.

[10] E. Yaka, S. Yılmaz, N. Özgür Doğan, and M. Pekdemir, "Comparison of the Glasgow-Blatchford and AIMS65 Scoring Systems for Risk Stratification in Upper Gastrointestinal Bleeding in the Emergency Department," Acad. Emerg. Med., vol. 22, no. 1, pp. 22–30, Jan. 2015.

[11] S. Bagavathiappan et al., "Investigation of peripheral vascular disorders using thermal imaging," Br. J. Diabetes Vasc. Dis., vol. 8, no. 2, pp. 102–104, Mar. 2008.

[12] A. Yilmaz, K. Shafique, N. Lobo, X. Li, ¡ Teresa Olson, and M. A. Shah, "Target-Tracking in Flir Imagery using Mean-Shift and Global Motion Compensation.

[13] "FLIR ONE Pro Thermal Imaging Camera for Smartphones | FLIR Systems," 2017. [Online]. Available: https://www.flir.com/products/flir-one-pro/. [Accessed: 22-Apr-2019].

[14] Carlos Platero, "Procesamiento morfológico.

[15] Annette Madison R, Francisca Honold G, Fernanda Castro V, Juan Escobar B, Violeta Rivas P, and Francisco Barrera M, "Prevención de resangrado de várices esofágicas en pacientes con cirrosis a quienes se les aplicó stents de diámetro pequeño vs terapia médica basada en titulación hemodinámica," Gastroenterol. latinoam 2018, vol. Vol 29, pp. 81–86, 2018.

[16] M. Shaydakov and J. Diaz, "Effectiveness of infrared thermography in the diagnosis of deep vein thrombosis: an evidence-based review," J. Vasc. Diagnostics Interv., vol. Volume 5, pp. 7–14, Feb. 2017.

# Skyline Path Queries for Location-based Services

Nishu Chowdhury[1], Mohammad Shamsul Arefin[2,*]

Computer Science and Engineering Department, Chittagong University of Engineering and Technology
Chattogram, Bangladesh

*Abstract*—A skyline query finds objects that are not dominated by another object from a given set of objects. Skyline queries help us to filter unnecessary information efficiently and provide us clues for various decision making tasks. In this paper, we consider skyline queries for location-based services and proposed a framework that can efficiently compute all non-dominated paths in road networks. A path *p* is said to dominate another path *q* if *p* is not worse than *q* in any of the *k* dimensions and *p* is better than *q* in at least one of the *k* dimensions. Our proposed skyline framework considers several features related to road networks and return all non-dominated paths from the road networks. In our work, we compute skylines considering two different perspectives: business perspective and individual user's perspective. We have conducted several experiments to show the effectiveness of our method. From the experimental results, we can say that our system can perform efficient computation of skyline paths from road networks.

*Keywords—Skyline queries; trip planning; location-based services*

## I. INTRODUCTION

Given a k-dimensional database DB, a skyline query retrieves a set of skyline objects, each of which is not dominated by another object. An object p is said to dominate another object p` if p is not worse than p` in any of the k dimensions and p is better than p` in at least one of the k dimensions. Fig. 1 shows a typical example of skyline. The table in Fig. 1 is a list of five routes, each of which contains two numerical attributes–"Cost" and "Distance". In the list, R2 and R5 are dominated by R3, while others are not dominated by any other routes. Therefore, the skyline of the list is {R1, R3, and R4}. Such skyline results are important for users to take effective decisions over complex data having many conflicting criteria. A number of efficient algorithms for computing skylines from the database have been demonstrated in the literature [1, 2, 3, 4, 5, 6].

Location-based services (LBSs) use positioning technology and traditional map information to furnish mobile users with new sorts of on-line services.

Location-based services in road network are becoming more popular. With rapid growth of technology, skyline queries on road networks [14, 15, 16, 17] have attracted much attention now a days.

Traffic jam refers to a long line of vehicles stuck in a jam. It is a common problem in the big cities and towns like Dhaka city of Bangladesh. Many factors such as less number of roads, lack of modern proper traffic management systems, narrowness of the roads, and increase of vehicles are the main causes of traffic jams in cities like Dhaka. These traffic jams

are creating many problems such as not reaching in time at offices, ambulance carrying patients cannot reach at the hospitals in time etc.

In such a scenario, a well-developed location-based service that focuses on the road conditions such as traffic jam, number of passengers and cost to the destination can give some comforts to the people by choosing skyline routes from which people can select their desired paths based on their preferences.

Each road in a road network has multiple-path criteria such as the distance of the road, the travel time through that road, the number of travellers and the number of traffic. The last two factors vary according to time. Before starting a journey, a traveller may want to know about the conditions of the road taken on their destination at a specific point of time. He/she may also want to know trip cost and other conditions of the roads.

In this paper, we apply skyline queries to support location-based services for road networks. In our approach, at first, a user needs to choose his pick-up point and the destination point. Based on the choice of the source and the destination by a user, our system then finds all alternate routes from source to destination. Next, each route is represented with several features such as traffic conditions, travelling time, travelling costs, number of passengers available through that each routes etc. After representing each route with a number of features, we apply skyline queries to filter dominated routes and to return only useful routes for the users. From the return results, a user can select his desired path such as less cost path or less traffic path.

The remainder of this paper is organized as follows: Section II provides a brief review of related work. We provide motivating examples at Section III. Section IV describes different concept related to the paper. We provide detail description of our proposed approach at Section V. In Section VI we present the experimental results. Finally, we conclude our paper at Section VII.



| ID | Cost | Distance |
|----|------|----------|
| R1 | 3 | 8 |
| R2 | 5 | 4 |
| R3 | 4 | 3 |
| R4 | 9 | 2 |
| R5 | 7 | 3 |

(a) Roads

(b) Skyline

Fig. 1. Skyline Example.

*Corresponding Author

## II. Related Work

Since the introduction of skyline queries in 2001, there are many works related to skyline queries considering different settings.

The Block-Nested Loops Algorithm (BNL) [4], which is the easiest skyline query method. Its objective is to build a candidate skyline set. This calculation investigates every data point with each other data point in the dataset. The BNL calculation requires each data point in the database be checked and tried for predominance; consequently, the time required for calculation increments with the volume of information.

The DAC calculation [4] separates information into groups and at that point leads skyline query in each group. The results are consolidated to acquire a definitive result.

The SaLS aalgorithm [2] utilizes an element acquired from the raw information as a threshold value with which to filter and dispose data points.

The BBS algorithm [18] is at present the most well-known skyline query algorithm. The BNL and DAC algorithms require that a large portion of the data points be processed all together to complete the skyline query comparison. Conversely, BBS utilizes an index structure for the identification of skyline points, which diminishes the number of points that must be tested all together to process a query.

Previous studies in which skyline queries were utilized to check road networks can be classified into those attempting to recognize skyline landmarks and those trying to distinguish skyline paths.

Deng et al. [8] presented the idea of searching for skyline landmarks in street network. The skyline landmark query recognizes landmarks that coordinate user criteria when user is going on a road network. For instance, when a user travels on a road network, skyline landmark query encourages him/her those points that are adjacent. The algorithm in this work characterizes landmark attributes as static or dynamic. Static properties have fixed values. Dynamic properties have variable attributes. The algorithm initially distinguishes static skyline landmarks on their static attribute values. At that point, when users perform to check, the algorithm distinguishes all unique skyline landmarks dependent on their dynamic attribute values, and consolidates query points with the static skyline landmarks. At last, the algorithm can recover skyline landmarks that fit all characteristics.

Huang and Jensen [13] proposed an alternate skyline landmark search concept from that of [8]. They contended that users' movement in road networks ought to be founded on a recently settled way. The algorithm in this work was like that proposed by Deng et al. [8], which utilized the ideas of static and dynamic attributes to identify skyline landmarks. The main contrast between the algorithms is the attribute calculation method. Deng et al. [8] considered the separation between the landmark and the inquiry area of the user, while the researchers of this work consider the separation between the landmark and the path preset by the user.

Tian et al. [21] presented the idea of skyline paths. Their proposed algorithm would utilize the edge attributes of a road network to discover all skyline paths between the user-specified starting vertex and goal. The algorithm would first decide a single skyline path between a starting vertex and goal whose summation of all attributes values is the most reduced among all ways. At that point, the algorithm would recognize other skyline paths by (1) a greedy algorithm to locate a relay vertex between starting vertex *s* and goal *t*. If skyline path domination was available after adding the two values, at that point *a* can't be a piece of a skyline path. In this instance, the algorithm again employs the greedy algorithm to identify other possible relay vertices or identify the next relay vertex following *a*.

Kriegel et al. [15] utilized the greedy algorithm to distinguish a possible relay vertex among *s* and *t*. Kriegel et al. [15] utilized a reference vertex to help estimations. By utilizing such a technique, they proclaimed the strategy proposed in this work was quicker than that proposed in crafted by [21].

Many researchers have looked to broaden the works in [15] and [21]. Aljubayrin et al. [1] examined the issue of skyline trips on different POI classes. Hsu et al. [10] connected the possibility of a skyline path to the arranging of treks to beat the conventional problem of acquiring multicriteria answers. Yang et al. [23] consolidated GPS history information in their inquiries to enable the user to design their skyline route under time-varying vulnerability. Unfortunately, these works don't consider aggregate attributes in road networks, which make them inapplicable to the issues tended to in this examination.

A new concept M-tree structure is described in [7]. A. Guttman al. [9] describes dynamic index structure called an R-tree and W. Son al. [20] describes spatial skyline queries for dynamic environment.

In [11], they focus on processing the continuous skyline query in road networks. They design a grid index to effectively manage the information of data objects. They proposed several algorithms combined with the grid index to answer the skyline queries.

In [12], they overcome the specific assumptions that each object is static in road networks. They focus on processing the CKNSQ over moving objects with uncertain dimensional values in Euclidean space and the velocity of each object (including the query object) varies within a known range.

Sheng et al. [19] present external memory algorithms for solving the skyline problem its variants in a worst-case efficient manner. They proved that the running time can be improved if some dimensions have small domains.

In [22], they bring out novel information by analyzing bulky databases to consolidate users experience to find place of interest. They use Apriori algorithm for identifying hidden association among item sets from large databases of user checking in data and to construct the route analogous to the key terms provided by user.

## III. Motivating Examples

Consider the graph of Fig. 2 that represents a road network where *L1* is considered as a source location and *L2* is a

destination location. Each vertex in Fig. 2 represents a junction i.e. dropping and/or pickup point and each edge represents a connection between two vertices. There are four values associated with each edge those are travel time, cost, distance and passengers, respectively. For example, edge (*L1, l1*) has values (3.18, 0.35, 0.7, 4), which indicates that the required time to reach from L1 to l1 is 3.18, cost of is 0.35, the distance between *L1* and *l1* is 0.7 and number of available passengers in is 4. In Fig. 2 if we use the route <*L1, l5, l6, l2, l8,l9,l4, L2*> to reach from *L1* to *L2,* we need total time 3.18 + 5 + 2.27 + 3.18 + 4.54 +4.09 + 1.81 = 24.07 and cost is 0.35 + 0.55 + 0.25 + 0.35 + 0.5 + 0.45 + 0.2 = 2.65. Here, total distance is 0.7 + 1.1 + 0.5 + 0.7 + 1 + 0.9+0.4 = 5.3 and the number of available passengers in this route is 4 + 5+ 7 + 6 + 8 + 10 + 10 = 50.

In this paper, we have considered two different scenarios. One is for business purpose and another is for individual user's perspective. These scenarios are explained below.

### A. Business Perspective

In applications such as Pathao and Uber, one trip can only be allotted to one passenger request at a specific time. In contrast, microbuses and cars have the capacity to carry five to eight passengers, respectively. Let us assume that a service provider has a microbus, which is a 10-seater vehicle. Suppose this person plans a trip from location *L1* to *L2*. Before starting the journey, by utilizing our method this person can find shorter route as well as a faster route. Our method also suggests a route having a large number of passengers, whereas for a fast route this method provides a route with a shorter distance and lesser traffic. Here, the passengers are those whose destination location is the same as that of the service provider and the start location belongs to the list of suggested routes. Thereby, the service provider can choose its preferable route and accept the passenger request for the same.

Table I shows the distance information on all routes and Fig. 3 presents the traffic and passenger conditions of two alternate routes. In Fig. 3, the graph lines are represented by three colours: green represents light traffic, orange represents medium traffic, and red indicates high traffic.

Hence, compared to multiple routes, it is necessary to find a desired route that is not dominated by any other route. In detail, a route is preferable to visitors if it is not dominated by any other route. The information on routes is given below. This information is collected from Google Map API.



Fig. 2. Example of a Graph Representing a Road Network.

TABLE I. INFORMATION ON ROUTES

| Route | Starting | Ending | Distance | Locations |
|-------|----------|--------|----------|-----------|
| Route 1 | L1 | L2 | 5.3 Km | L1,l5,l6,l2,l8,l9,l4,L2 |
| Route 2 | L1 | L2 | 4.8 Km | L1,l1,l2,l3,l4,L2 |
| Route 3 | L1 | L2 | 5.6 Km | L1,l5,l6,l2,l3,l4,L2 |
| Route 4 | L1 | L2 | 4.5 Km | L1,l1,l2,l8,l9,l4,L2 |

The traffic and passenger conditions of two routes are graphically represented in Fig. 3. The route line colour changes with time. From this figure, we can say that there is a light traffic for Road 1 from 8.00 a.m. to 12 p.m., and Road 2 will be free after 2 p.m. It is also observed that there is heavy traffic for Road 1 from 3 p.m. and for Road 2 from 11 a.m. to 2 p.m. This graph is also helpful in tracing a medium traffic condition. Road 1 has medium traffic from 12 p.m. to 3 p.m. and Road 2 from 8 a.m. to 11 a.m. We can also calculate our travel cost from route distance.

Table II represents the traffic and passenger conditions of every location for Route 1 and 2. Each row in the table represents traffic and passenger conditions. For measuring passenger we use normalize value. These are helpful in identifying the most interesting and preferable route.

Another graphical representation is given below. Fig. 4, represents the passenger condition with traffic. In this figure, blue colour Route 1 and green Road 2.

From Table II we find skyline points H3, H4.



Fig. 3. Traffic Condition of Two Routes with Respect to Time.

TABLE II. ROAD CONDITIONS FOR DIFFERENT TIME INTERVALS

| ID | Route | Travel Time | Passenger | Traffic | Distance (Km) |
|----|-------|-------------|-----------|---------|---------------|
| H1 | Route 1 | 8.00 am | 50 | Low | 5.3 |
| H2 | Route 1 | 2.00 pm | 40 | High | 5.3 |
| H3 | Route 2 | 8.00 am | 41 | High | 4.8 |
| H4 | Route 2 | 9.00 pm | 35 | Medium | 4.8 |

Fig. 4.    Traffic and Passenger Conditions.

### B.  Individual Perspective

Suppose a visitor wishes to travel from Location *L1* to *L2*. Before starting their journey, he/she wishes to know about the route and traffic conditions as well as the cost of travel. Then, this method will provide him/her with traffic information using Google Map Traffic API and calculate the cost by calculating the fuel cost per litre, the mileage of their vehicle and the distance between their start and end locations. From the resulted dataset a user can easily filter routes according to choice.

Table III represents road condition for specific interval. Each row in the table represents traffic and cost, which are helpful to identify the most interesting and preferable route.

From Table III, we find that cost is changing with user. Moreover, cost depends on user vehicle's fuel cost and mileage.

In this paper, we compute a method that can help service providers to choose their desired route from our resulted skyline routes. Our location-based computation method can significantly find the appropriate route, based on the dataset. By this way, our method is useful for individual trip planning and transport service business planning.

TABLE III.    TRAVEL COST

| ID | User | Start point | End point | Route | Distance (Km) | Traffic | Cost |
|----|------|-------------|-----------|-------|---------------|---------|------|
| H1 | user 1 | *L1* | *L2* | Route 1 | 5.3 | High | 2.65 |
| H2 | user 1 | *L1* | *L2* | Route 2 | 4.8 | Low | 2.4 |
| H3 | user 1 | *L1* | *L2* | Route 3 | 5.6 | High | 2.5 |
| H4 | user 1 | *L1* | *L2* | Route 4 | 4.5 | High | 2.25 |
| H5 | user 2 | *L1* | *L2* | Route 1 | 5.3 | High | 3.65 |
| H6 | user 2 | *L1* | *L2* | Route 2 | 4.8 | High | 3.4 |
| H7 | user 2 | *L1* | *L2* | Route 3 | 5.6 | Low | 3.5 |
| H8 | user 2 | *L1* | *L2* | Route 4 | 4.5 | Low | 3.35 |

## IV.  PRELIMINARIES

Consider a database *DB* with *N* attributes and *k* objects. Let *a1, a2,...,aN* be the *N* attributes of *DB*. We consider that smaller values in each attribute are better and that each attribute has positive values.

### A.  Skyline Queries

Skyline query is a decision-supporting mechanism that highlights the best options among vast data.

An example is given below:

In Fig. 5, we have some points in a two-dimensional space, as shown above, then we define a point *p* that will dominate point *q* provided its coordinates are larger than that of *q*. In this example, there is a point *p* that dominates several other points. So what is the skyline point? Skyline points are points that are not dominated by any other points present in the dataset. They are also called maximal points. If you connect these with horizontal and vertical lines, then you will get skyline points.

Let *L* denote a set of all locations. Each location has an ID and a spatial coordinate *l = (xy)*. Let us suppose *A* is a category attribute. In our research work, passenger and traffic are category attributes. So, we denote the coordinates of location *L* by *l. L*; *l.a* represents the value of attribute *A*.

Definition 1 (Dominance Relationship): Given two objects *a* and *a'* exist, then object *a* is said to dominate *a'* if *a < a'* for all the attributes.

Definition 2 (Skyline Query): Skyline query is the set of objects that cannot be dominated by any other object. Given point *p, r ∈ D*. If *p < r*, then *p* belongs to the skyline set.

### B.  Multi-Attribute Network Graph (MAG)

Graph *G (V, E, W)* is a multi-attribute network graph, where *V* denotes a set of vertices, *E* a set of edges, and *W* weight vector. In Fig. 6, nodes define profiles of activity, roles and actors etc. Edges define the relationship among those nodes or entities and weight defines the behaviour of the edges.



Fig. 5.    Skyline Points.



Fig. 6.    Multi-Attribute Network Graph (MAG).

## V. METHODOLOGY

Our method comprises two modules: the first module delves into the business perspective and the second into the individual perspective. Fig. 7 describes the proposed framework. In each module, the user can provide the source and destination addresses while prioritising a specific destination based on his/her choice.

The business perspective and individual perspective operates in three processes or modules: processing module, query execution module and output module. Our functional algorithm, which parses the dataset and the filters, is known as the processing module. In terms of both perspectives, it works in five steps: first, it measures the geolocations of the start and end locations. Upon completion, an iteration process continues to measure all alternate routes from the source location to the destination location. Thereafter, it calculates traffic, trip cost or passengers based on the dataset. Thereafter, the process migrates into the query execution module, where a resulted dataset is generated imposing skyline queries. Through these processes, we get the dominant paths that are filtered later on the system output, which shows the result of these potential paths.

The most naïve approach to locating skyline paths in a road network is to identify all of the paths between the origin and destination in the network, calculate attributes of the paths, and perform a dominance check of all the attributes. The process of estimating traffic, cost and passenger are given in below.

### A. Traffic Estimation

Fig. 8 describes the traffic condition at a specific time. Suppose we wish to assess the traffic conditions in all alternate routes from *L1* to *L2* at 2 pm. In this framework, car has been used as a transport mode. We get four routes from the given graph: Route 1 comprises *L1, l5, l6, l2, l8, l9, l4, L2* and Route 2 comprises *L1, l1, l3, l4, L2*. Similarly Route 3 contains *L1, l5, l6, l2, l3, l4, L2* and Route 4 contains *L1, l1, l2, l8, l9,l4,L2*. For assessing the traffic conditions at a specific time, Google Map Traffic API is used. For example, if someone wants to assess the traffic condition from *L1* to *L2* at 2 pm, then the system counts all alternate routes that he/she can take to reach the destination. Thereafter, it uses the latitude and longitude of a distance at every 0.5 km interval and checks the location at each iteration. Whenever a new location returns, we measure the traffic condition at those points by employing Google Map Traffic API. It provides the standard time and the time required to reach one's destination, and then it stores all the data on the latter for every 0.5 km interval. In this way, we can obtain all the data on the time taken for all alternate routes. In this figure, the blue-coloured text represents the standard time (in minutes) taken to travel from one location to another. Another colour represents the time required to travel from one location to another. In this figure, three different colours are used: orange is used to represent medium traffic, green to indicate low traffic and red for heavy traffic. When the standard time is equal to the required time, the given time interval contains medium traffic. If the required time is low, it indicates the presence low traffic. Otherwise, the presence of heavy traffic is indicated.

By this way, we gather data for our dataset. Tables IV, V, VI and VII represent the required time for Route 1, Route 2, Route 3 and Route 4 respectively.

Now, we calculate the total required time for each route. From above dataset we find Route 1, Route 2, Route 3 and Route 4 require 24.07, 21.79, 25.43 and 20.43 minutes respectively.

### B. Cost Estimation

Table III represents user wise cost for each route. Here, we represent how cost is changing with distance, mileage and fuel consumption. In our method, cost measures by using the following formula.

$$\text{Cost} = (\text{mileage/per ltr fuel cost}) * \text{Distance} \qquad (1)$$



Fig. 7. System Architecture.

TABLE IV.    REQUIRED TIME AND DISTANCE FOR ROUTE 1

| Source | Destination | Required Time (min) | Distance (Km) |
|--------|-------------|---------------------|---------------|
| L1 | l5 | 3.18 | 0.7 |
| l5 | l6 | 5 | 1.1 |
| l6 | l2 | 2.27 | 0.5 |
| l2 | l8 | 3.18 | 0.7 |
| l8 | l9 | 4.54 | 1 |
| l9 | l4 | 4.09 | 0.9 |
| l4 | L2 | 1.81 | 0.4 |

TABLE V.    REQUIRED TIME AND DISTANCE FOR ROUTE 2

| Source | Destination | Required Time (min) | Distance (Km) |
|--------|-------------|---------------------|---------------|
| L1 | l1 | 4.54 | 1 |
| l1 | l2 | 2.27 | 0.5 |
| l2 | l3 | 6.81 | 1.5 |
| l3 | l4 | 6.36 | 1.4 |
| l4 | L2 | 1.81 | 0.4 |

TABLE VI.    REQUIRED TIME AND DISTANCE FOR ROUTE 3

| Source | Destination | Required Time (min) | Distance (Km) |
|--------|-------------|---------------------|---------------|
| L1 | l5 | 3.18 | 0.7 |
| l5 | l6 | 5 | 1.1 |
| l6 | l2 | 2.27 | 0.5 |
| l2 | l3 | 6.81 | 1.5 |
| l3 | l4 | 6.36 | 1.4 |
| l4 | L2 | 1.81 | 0.4 |

TABLE VII.    REQUIRED TIME AND DISTANCE FOR ROUTE 4

| Source | Destination | Required Time (min) | Distance (Km) |
|--------|-------------|---------------------|---------------|
| L1 | l1 | 4.54 | 1 |
| l1 | l2 | 2.27 | 0.5 |
| l2 | l8 | 3.18 | 0.7 |
| l8 | l9 | 4.54 | 1 |
| l9 | l4 | 4.09 | 0.9 |
| l4 | l2 | 1.81 | 0.4 |

## C. Passenger Estimation

Fig. 9 describes the condition of passenger at specific time; suppose, we need to assess the condition of passenger of all alternate routes from location *L1* to *L2*. The passenger condition for specific time for each location can be assessed through passenger request. Fig. 8 shows the passenger condition. From this Fig. 8, we have found four alternate routes: Route 1 comprises L1, l5, l6, l2, l8, l9, l4, L2 and Route 2 comprises L1, l1, l3, l4, L2. Similarly  Route 3 contains L1, l5, l6, l2, l3, l4, L2 and Route 4 contains L1, l1,

l2, l8, l9,l4,L2. For example, we want to measure number of passengers for Route 1.At first, we count passengers of all location of Route 1, whose destination location is *L2*. Now get the maximum value from these locations. Suppose the value is *P*. We use the following formula to normalize location wise passengers.

$$P_i = P + 1 - P_i \qquad (2)$$

Tables VIII, IX, X and XI represent the condition of passengers for Route 1, Route 2, Route 3 and Route 4, respectively.

Now, calculate the condition of passenger for each route. Route 1, Route 2, Route 3 and Route 4 has 50, 41, 40 and 51 passengers, respectively.



Fig. 8.    Traffic Condition and the Distance of all Alternate Routes from *L1* to *L2*.

TABLE VIII.    PASSENGER CONDITION FOR ROUTE 1

| Source | Destination | Number of Passengers |
|--------|-------------|----------------------|
| L1 | l5 | 4 |
| l5 | l6 | 5 |
| l6 | l2 | 7 |
| l2 | l8 | 6 |
| l8 | l9 | 8 |
| l9 | l4 | 10 |
| l4 | L2 | 10 |

TABLE IX.    PASSENGER CONDITION FOR ROUTE 2

| Source | Destination | Number of Passengers |
|--------|-------------|----------------------|
| L1 | l1 | 9 |
| l1 | l2 | 8 |
| l2 | l3 | 9 |
| l3 | l4 | 5 |
| l4 | L2 | 10 |



Fig. 9.    Passenger Condition from L1 to L2.

TABLE X.    PASSENGER CONDITION FOR ROUTE 1

| Source | Destination | Number of Passengers |
|--------|-------------|----------------------|
| L1 | l5 | 4 |
| l5 | l6 | 5 |
| l6 | l2 | 7 |
| l2 | l3 | 9 |
| l3 | l4 | 5 |
| l4 | L2 | 10 |

TABLE XI.    PASSENGER CONDITION FOR ROUTE 2

| Source | Destination | Number of Passengers |
|--------|-------------|----------------------|
| L1 | l1 | 9 |
| l1 | l2 | 8 |
| l2 | l8 | 6 |
| l8 | l9 | 8 |
| l9 | l4 | 10 |
| l4 | L2 | 10 |

### D. Computing Skyline

Here, we generate a dataset that measure attributes such as traffic, passenger, cost and distance.

Let's consider a scenario. Suppose source is *S* and destination is *D*. There are ten alternate routes from *S* to *D*. We denote traffic condition as low, medium and high and define them as 1, 2, and 3 respectively. Table XII represents the route condition for a specific time.

From this dataset we need desire routes. By using BBS [18] algorithm we get our skyline routes. Fig 10 describes the BBS algorithm.

Using BBS algorithm, we get our skyline routes as *R1*, *R2*, *R3* and *R4*. From this method, a user can easily find his/her desire route in proficient and appropriate way.  If one wants a large passenger, low traffic and low cost route, then he/she can get the desired routes from the resulted routes.

TABLE XII.    ROAD CONDITION FOR A SPECIFIC TIME

| Route | Traffic | Cost | Passenger | Distance (Km) |
|-------|---------|------|-----------|---------------|
| R1 | 1 | 55 | 20 | 20 |
| R2 | 3 | 60 | 9 | 10 |
| R3 | 2 | 50 | 10 | 15 |
| R4 | 1 | 45 | 8 | 30 |
| R5 | 3 | 100 | 20 | 50 |
| R6 | 2 | 120 | 30 | 60 |
| R7 | 2 | 110 | 50 | 70 |
| R8 | 2 | 130 | 40 | 50 |
| R10 | 2 | 150 | 50 | 30 |
| R11 | 3 | 140 | 60 | 40 |

**Algorithm 1:** BBS

**Input:** A dataset $D$ (r-tree).

**Output:** The Set of skyline points of dataset $D$.

1. $S=\emptyset$ // list of skyline points
2. insert all entries of the root $R$ in the heap
3. **while** heap not empty
4. remove top entry $e$
5. **if** $e$ is dominated by some point in $S$ discard $e$
6. **else** // $e$ is not dominated
7. **if** $e$ is an intermediate entry
8. **for** each child $e_i$ of $e$
9. **if** $e_i$ is not dominated by some point in $S$
10. insert $e_i$ into heap
11. **else** // $e$ is a data point
12. insert $e_i$ into $S$
13. **end while**
14. **end**

Fig. 10. BBS Algorithm for Skyline Computation (Adapted from [18]).

## VI. EXPERIMENTS

We have implemented our proposed system in .Net Framework. We have performed the experiment in a simulation environment of a PC running on windows OS having an Intel(R) Core i7, 1.73 GHz CPU and 4 GB main memory. Due to the lack of real data, we evaluate our proposed algorithm using synthetic datasets only.

Fig. 11 shows the results when we consider Route 1, Route 2, Route 3 and Route 4. We observe that with the increases of distance, number of passengers varies.

Fig. 12 shows that with the increase of distance, number of routes also increases.

In Fig. 13, when we consider two (2D), three (3D), four (4D), and five (5D) features. We observe that with the increases of routes, there is very slight increase in computation time. This is because during the computation process, time increases with the increase of number of routes. We can also observe that computation time gradually increases if the number of features increases.



Fig. 11. Passenger Varies with Distance.



Fig. 12. Number of Routes Varies with the Distance.



Fig. 13. Time Varies with Number of Routes.



Fig. 14. Skyline Points Varies with Number of Routes.

Simultaneously, it is also observed that skyline points increase with the number of routes and number of features. Fig. 14 represents how skyline points increase with the number of routes.

## VII. CONCLUSION

With the rapid growth of civilization, traffic is seen to increase day by day. Therefore, collecting traffic information, passenger condition and cost calculation has become a popular method. Our experimental results demonstrate that the proposed algorithm is scalable enough to compute the skyline path for a specific time. The proposed approach can easily expand for recommendation. In this work, we performed different analyses on synthetic data. In future, we aim to expand large passenger route methodology in more efficient way and find desire route based on user preference. So that we can get skyline points in more proper ways. We also want to trace the vehicle movement and position in a more efficient and effective way.

### REFERENCES

[1] S. Aljubayrin, Z. He, and R. Zhang, ''Skyline trips of multiple POIs categories'', in Proc. Int. Conf. Database Syst. Adv. Appl. (DASFAA), 2015, pp. 189–206.

[2] I. Bartolini, P. Ciaccia, and M. Patella, ''SaLSa: Computing the skyline without scanning the whole sky'', in Proc. Int. Conf. Inf. Knowl. Manage. (CIKM), 2006, pp. 405–414.

[3] N. Beckmann, H.-P. Kriegel, R. Schneider, and B. Seeger, ''The R*-tree: An efficient and robust access method for points and rectangles'', in Proc. ACM Special Interest Group Manage. Data (SIGMOD), 1990, pp. 322–331.

[4] S. Borzsony, D. Kossmann, and K. Stocker, ''The skyline operator'', in Proc. IEEE 17th Int. Conf. Data Eng. (ICDE), Apr. 2001, pp. 421–430.

[5] Z. Chen, H. T. Shen, X. Zhou, and J. X. Yu, ''Monitoring path nearest neighbor in road networks'', in Proc. ACM Special Interest Group Manage. Data (SIGMOD), 2009, pp. 591–602.

[6] J. Chomicki, P. Godfrey, J. Gryz, and D. Liang, ''Skyline with presorting'', in Proc. IEEE Int. Conf. Data Eng. (ICDE), Mar. 2003, pp. 717–719.

[7] P. Ciaccia, M. Patella, and P. Zezula, ''M-tree: An efficient access method for similarity search in metric spaces'', in Proc. Int. Conf. Vary Large Data Bases (VLDB), 1997, pp. 426–435.

[8] K. Deng, Y. Zhou, and H. Tao, ''Multi-source skyline query processing in road networks'', in Proc. IEEE Int. Conf. Data Eng. (ICDE), Apr. 2007, pp. 796–805.

[9] A. Guttman, ''R-trees: A dynamic index structure for spatial searching'', SIGMOD Rec., vol. 14, no. 2, pp. 47–57, 1984.

[10] W. T. Hsu, Y. T. Wen, L. Y. Wei, and W. C. Peng, ''Skyline travel routes: Exploring skyline for trip planning'', inProc. IEEE Int. Conf. Mobile Data Manage. (MDM), Jul. 2014, pp. 31–36.

[11] Y.-K. Huang, C.-H. Chang, and C. Lee, ''Continuous distance-based skyline queries in road networks'', Inf. Syst., vol. 37, no. 7, pp. 611–633, 2012.

[12] Y.-K. Huang and Z.-H. He, ''Processing continuous K-nearest skyline query with uncertainty in spatio-temporal databases'', J. Intell. Inf. Syst., vol. 45, no. 2, pp. 165–186, 2015.

[13] X. Huang and C. S. Jensen, ''In-route skyline querying for location-based services'', in Proc. Web Wireless Geograph. Inf. Syst. (W2GIS), 2004, pp. 120–135.

[14] S. Jang and J. Yoo, ''Processing continuous skyline queries in road networks'', in Proc. Int. Symp. Comput. Sci. Appl., Oct. 2008, pp. 353–356.

[15] H.-P. Kriegel, M. Renz, and M. Schubert, ''Route skyline queries: A multipreference path planning approach'', in Proc. IEEE Int. Conf. Data Eng. (ICDE), Mar. 2010, pp. 261–272.

[16] F. Li, D. Cheng, M. Hadjieleftheriou, G. Kollios, and S.-H. Teng, ''On trip planning queries in spatial databases'', in Proc. Int. Symp. Spatial Temporal Databases (SSTD), 2005, pp. 273–290.

[17] S. Pan, Y. Dong, J. Cao, and K. Chen, ''Continuous probabilistic skyline queries for uncertain moving objects in road network'', Int. J. Distrib. Sensor Netw., vol. 2014, Mar. 2014, Art. no. 365064.

[18] Papadias, Y. Tao, G. Fu, and B. Seeger, ''An optimal and progressive algorithm for skyline queries'', in Proc. ACM Special Interest Group Manage. Data (SIGMOD), 2003, pp. 467–478.

[19] C. Sheng and Y. Tao, ''Worst-case I/O-efficient skyline algorithms'', ACM Trans. Database Syst., vol. 37, no. 4, 2012, Art. no. 26.

[20] W. Son, S.-W. Hwang, and H.-K. Ahn, ''MSSQ: Manhattan spatial skyline queries'', Inf. Syst., vol. 40, pp. 67–83, Mar. 2014.

[21] Y. Tian, K. C. K. Lee, and W.-C. Lee, ''Finding skyline paths in road networks'', in Proc. ACM Int. Conf. Adv. Geograph. Inf. Syst. (SIGSPATIAL), 2009, pp. 444–447.

[22] Y.-T. Wen, K.-J. Cho, W.-C. Peng, J. Yeo, and S.-W. Hwang, ''KSTR: Keyword-aware skyline travel route recommendation'', in Proc. IEEE Int. Conf. Data Mining (ICDM), Nov. 2015, pp. 449–458.

[23] B. Yang, C. Guo, C. S. Jensen, M. Kaul, and S. Shang, ''Stochastic skyline route planning under time-varying uncertainty'', in Proc. IEEE Int. Conf. DataEng. (ICDE), Mar./Apr. 2014, pp. 136–147.

# IoT-Enabled Door Lock System

Trio Adiono[1], Syifaul Fuada[2]

School of Electrical Engineering and Informatics, Institut Teknologi Bandung, Jln. Ganesha No.10, Gd. Ahmad Bakrie (LABTEK VIII) Lt. IV, ZIP 40116, Indonesia[1]
Program Studi Sistem Telekomunikasi, Universitas Pendidikan Indonesia[2]

Sinantya Feranti Anindya[3], Irfan Gani Purwanda[4]
Maulana Yusuf Fathany[5]

University Center of Excellence on Microelectronics, Institut Teknologi Bandung, Jln. Tamansari No.126, IC Design Laboratory, Gd. PAU Lt. IV
ITB Campus, 40132

*Abstract*—**This paper covers the design of a prototype for IoT and GPS enabled door lock system. The aim of this research is to design a door lock system that does not need manual input from user for convenience purpose while also remaining secure. The system primarily consists STM32L100 microcontroller as its core, TIP102 transistor that controls 12 $V_{DC}$ solenoid, and Xbee module to communicate with the smart home's host and receive status regarding user's GPS position. The system is tested by measuring the user's distance from the predetermined location using GPS coordinate captured by an Android application, which serves to test whether the system is able to operate as intended and measure the device's power usage. The test result shows that the device is able to work based on GPS coordinate data received, using 42.3 mA and 587 mA current in idle and active modes, respectively.**

*Keywords—Internet of Things; smart home; smart lock*

## I. Introduction

The Internet of Things (IoT) is a concept revolving around global information network that consists 'things' such as smart devices, sensors and actuators, and even smaller networks with their own identities and ability to self-configure and make their own decisions within certain extent, whether individually or collectively [1-2]. The advent of IoT foreshadows the future where everything and everyone is connected to the Internet through all devices they possess, whether from computer, smart phone, and other consumer devices. Objects within IoT network can also communicate using various communication technologies such as WiFi, Bluetooth, near-field communication, and many more.

One of the most common applications of the IoT technology is for supporting a smart home system. Recent smart home systems make good model on how an IoT architecture behaves, as objects in a house are connected to a gateway in wireless manner to communicate whether with each other and the house's inhabitants. As a smart home system is intended to enhance the inhabitants' living quality in terms of comfort of convenience, the IoT is utilized to at least enable easier control and monitoring over the home devices.

A smart home system can be categorized into at least one type based on its functionality: health care, entertainment, energy, and/or security [3]. Among these functionalities, security becomes one of the most crucial factors behind installation of the system. Enhancing home security using smart home can be done in a number of ways, including but not limited to installation of smart, customized door lock. Several examples of smart door lock implementation have already existed, such as camera-based door security systems [4-6], passwords [6-7], smart card [8], and proximity or location detection [9]. Each of the aforementioned methods has its own strengths and weaknesses, such as interoperability value of devices.

In this research, a GPS-based smart door lock is designed. The ultimate aim of this research is to design a door lock system that does not need manual input from user for convenience purpose while also remaining secure. However, for this publication, the scope is limited for examining the feasibility of utilizing GPS for location-based lock to achieve the aforementioned goal. Furthermore, the scope of this paper will be mostly limited to the lock hardware side, while the server and the Android application code handling the GPS tracking will only be briefly explained. This device is part of MINDS smart home system we designed, which also include RGB lamp [10], curtain controller [11], humidity and temperature sensor [12], fan controller [13], and infrared remote [14]. The smart home is controlled using Android-based application detailed in [15]. This paper details the design of the smart lock design and implementation, as well as the testing result.

This paper is divided into following sections: Section I presents the research background. Section II discusses the methodology that cover: highlighted research limitation, specification of developed system, hardware, system workflow, design of protocol, android and server. Section III reports the research result while the research conclusion and future work are given in Section IV and Section V, respectively.

## II. System Design

### A. Research Limitation

As mentioned above, this work is more focuses on the hardware part including defined diagram block, electronic circuit and its assembly, whereas the software part including system flowchart, android apps, and server are presented briefly, and it will be elaborated in detail for other publications. The test limit on the functional test (locked and unlocked condition) and power consumption measurement using digital multimeter. Further performance metrics will be explored in future.

## B. System Specification

As the designed system is part of MINDS smart home system, the overall architecture is similar to other devices designed for the aforementioned smart home. The smart door lock is one of the nodes supported by MINDS system, which is connected to a Raspberry Pi-based host that serves as a 'bridge' for the nodes to receive commands from user device in the form of Android-based smart phone. As for the smart phone, it is connected to the host through either Bluetooth or the Internet (via cloud server). The architecture of the MINDS system is depicted in Fig. 1.

## C. Hardware Design

The door lock works like an on/off switch, which is controlled based on user's proximity from the door. The lock system utilizes a battery-powered STM32L100 microcontroller as its core, which is used for controlling 12 $V_{DC}$/630 mA solenoid and an Xbee module to receive signals from the central host. The solenoid is controlled using TIP 102 transistor by utilizing its cut-off and saturation modes to switch the solenoid's state. Furthermore, the system is equipped with components such as DC power jack, microUSB port for development purpose, switch, and reset and mode buttons. The block diagram of the door lock's hardware is depicted in Fig. 2, while the circuit diagram for the solenoid control is depicted in Fig. 3.

As the transistor is required to work in saturation and cut-off region, the value of components such as resistor must be determined in such way so they can support these operation regions. According to TIP 102's datasheet, the base current ($I_B$) for saturation region is 1.9 mA. As such, to achieve the current, the equation (1) is used for calculating the base resistor value.

$$I_B = \frac{Vout - VBE(on)}{RBasis} \tag{1}$$

The value of the $V_{out}$ is based on STM32L100's output port voltage, which is 3.3 $V_{DC}$, while the voltage between base and emitter ports of the transistor is based on the transistor's Darlington configuration, hence 1.4 $V_{DC}$ (0.7 $V_{DC}$+0.7 $V_{DC}$). Based on the calculation, it can be inferred that 1k$\Omega$ resistor is required to make the base current into 1.9 mA. Aside from the base current, the maximum power the transistor can handle also needs to be considered, especially when the solenoid is active. According to the datasheet, the transistor is able to handle 80 watt of power dissipation. Assuming the $V_{CE}$ during saturation is 2 $V_{DC}$, then the solenoid power is:

$$P(transistor) = V_{CE}(sat) * I_C = 2V * 630mA = 1.26Watt$$

## D. System Workflow

This section describes two separate programs required to operate the lock: the firmware inside the microcontroller and the mobile phone (Android) software to transmit GPS coordinate of the user to the host. The flowchart for the firmware is depicted in Fig. 4, while the flowchart for the Android application is depicted in Fig. 5.



Fig. 1. Architecture of MINDS Smart Home System.



Fig. 2. Block Diagram of the Door Lock System.



Fig. 3. Circuit Diagram for the Solenoid Control.



Fig. 4. Flowchart for Firmware for the Door Lock Control.

Fig. 5. Flowchart for Android Application for Door Lock Control, Reproduced from [15] under Permission.



Fig. 6. Packet Data Structure for Door Lock Control, Reproduced from [16-17] under Permission.

### E. Protocol Design

The message protocol used to transmit GPS data to the host is based on the protocol elaborated in [16]. The specific structure of the message for door lock control is depicted in Fig. 6 (consisting the packet header, address, packet init, Data payload, and checksum).

### F. Android Application and Server Design

The door lock works by having the Android application that serves as user interface of the entirety of MINDS system to periodically detect GPS coordinate of the user, then send it to the server (and host) to be compared with the house's coordinate. If the distance is 10 meters or less, the application will send command to both the server to unlock the house, which will be relayed to the host then to the lock hardware. Likewise, if the user's distance with the house is more than 10 meters, then the system will be locked. The system will send a warning if an unauthorized attempts to access the device within the house is made while the house is locked.

### III. RESULTS AND DISCUSSION

### A. System Implementations

Based on the block diagram depicted in Fig. 2 earlier, the printed-circuit board for the controller is designed with the result depicted in Fig. 7. This controller design is similar to the controller from previous project [13]. The reasoning behind the similar design is to add interoperability value to the system while also enhancing its performance.

To control the solenoid lock, the $0^{th}$ pin of STM32L1000's GPIO-B port is used as the output port. The software implementation for the port configuration is elaborated in Table I, while the control process is elaborated in Table II.

As previously explained, this system works together with the host and MINDS Android application. Fig. 8 displays the comparison between the MINDS application during normal (unlocked) condition and locked condition. The application

captures the user's location every 5 minutes using GPS, then send the result to the server (and host), after which the host will send command to the lock device. If the house is supposed to be locked, then a warning will be sent if there is any unauthorized attempt to control the device. The code implementation for the GPS tracking is depicted in Table III, while the code implementation for the control and warning by host is depicted in Tables IV, V and VI.



Fig. 7. The Appearance of PCB DC12V Controller for MINDS Door Lock Device at Actual Size (40.21 mm x 39.84 mm).

TABLE I. GPIO-B PORT CONFIGURATION IN STM32L100 MICROCONTROLLER

```
// GPIO Init
GPIO_InitTypeDef GPIO_InitStructure;
RCC_AHBPeriphClockCmd( RCC_AHBPeriph_GPIOB, ENABLE); // Enable
GPIOB clock
// GPIO Configuration
GPIO_InitStructure.GPIO_Pin = GPIO_Pin_0;
GPIO_InitStructure.GPIO_Mode = GPIO_Mode_OUT;
GPIO_InitStructure.GPIO_OType = GPIO_OType_PP;
GPIO_InitStructure.GPIO_PuPd = GPIO_PuPd_UP;
GPIO_Init(GPIOB, &GPIO_InitStructure);
```

TABLE II. SOFTWARE IMPLEMENTATION FOR THE DOOR LOCK ON/OFF CONTROL

```
if(rawlock==0x00) GPIO_LOW(GPIOB,GPIO_Pin_0);
else if(rawlock==0x01) GPIO_HIGH(GPIOB,GPIO_Pin_0);
```



(a)                          (b)

Fig. 8. The Appearance of MINDS Application During: (a) Normal, Retrieved from [18] under Permission and (b) Locked Conditions.

TABLE III.    CODE IMPLEMENTATION FOR GPS TRACKING IN ANDROID APPLICATION

```
if (head.equals("LC")) {
    if (SensorService.latcur != 0.0 && SensorService.lngcur != 0.0) {
        try {                                           // send position
            AMQP.BasicProperties props = new
            AMQP.BasicProperties
            .Builder()
            .correlationId(SensorService.phone_id)
            .build();
            JSONObject obj = new JSONObject();
            obj.put("head", "LC");
            obj.put("homeid", SensorService.home_id);
            obj.put("lat", Double.toString(SensorService
                                            .latcur));
            obj.put("lng", Double.toString(SensorService
                                            .lngcur));
            StringWriter toSent = new StringWriter();
            obj.writeJSONString(toSent);
            String tosendstring = toSent.toString();
            Toast.makeText(getApplicationContext(),Double
            .toString(SensorSer
            vice.latcur) + " " +
            Double.toString(SensorService.
            lngcur),
            Toast.LENGTH_SHORT).show();
            SensorService.channel.basicPublish("",
                SensorService.home_id+"AES",
                props, SensorService
                .AESencrypt(tosendstring)
                .getBytes("UTF-8"));
            SensorService.lngcur = 0.0;
            SensorService.latcur = 0.0;
        } catch (Exception e) {
Toast.makeText(getApplicationContext(), "Failed
            to Update Location",
        Toast.LENGTH_SHORT).show();
        }
    }
}
```

TABLE IV.    SOFTWARE IMPLEMENTATION FOR THE RASPBERRY PI HOST LOCATION COMPARISON

```
elif (sent['homeid'] == homeid and sent['head'] == 'LC') :        print (sent['lat'],
sent['lng'])
        cnx = mysql.connector.connect(user='root', password='28031995',
            host='localhost', database='home'+homeid)
        cursor = cnx.cursor()
        cursor.execute("SELECT latitude, longitude FROM info")
        data = cursor.fetchone()
                # Get the user's coordinate in float form
        latf = float(sent['lat'])
        lngf = float(sent['lng'])
                # If the coordinate is (0, 0), assume the user's phone is deactivated
        if (latf == 0.0 and lngf == 0.0):
            cursor.execute("UPDATE users SET location = 'in' WHERE
                phoneid = %s", (props.correlation_id,))
        # Check if user's radius is above 100 meters
        elif (sqrt(math.pow(latf-data[0]),2) + math.pow(lngf-data[1]),2)) >
            0.001) :
            cursor.execute("UPDATE users SET location = 'out'
                    WHERE phoneid = %s", (props.correlation_id,))
            # Check if user is within radius
        cursor.execute("SELECT EXIST (SELECT location FROM users
WHERE location = 'in')")
        # If not, assume the user is away
        if (not(cursor.fetchone()[0])):
            tosend['homeid'] = homeid
            tosend['head'] = 'SC'
            tosend['name'] = 'lock_scen'
        ch.basic_publish(exchange='', routing_key=sent['homeid']+'AES',
            properties=pika.BasicProperties(correlation_id =
            props.correlation_id),
            body=AESencrypt(json.dumps(tosend),key))
    except:
            print(sys.exc_info()[0])
    channel.basic_consume(callback, queue=queue_name,
    no_ack=True)
    channel.start_consuming()
```

TABLE V.    SOFTWARE IMPLEMENTATION FOR RASPBERRY PI HOST DEVICE LOCKING

```
elif (sent['homeid'] == homeid and sent['head'] == 'SC') :
    try :
        print (sent['name'])
        cnx = mysql.connector.connect(user='root',
                password='28031995', host='localhost',
                database='home'+homeid)
        cursor = cnx.cursor()

        if (sent['name'] == 'lock_scen') :
                cursor.execute("UPDATE info SET lockstatus =
                'lock'")
                cursor.execute("SELECT homekey FROM info")
                ch.basic_publish(exchange='amq.topic',
                        routing_key=homeid,
                        body=AESencrypt(json.dumps(sent),
                        cursor.fetchone()[0]))
                print('set lock')

        if (sent['name'] == 'unlock_scen') :
                cursor.execute("UPDATE info SET lockstatus =
                'unlock'")
                cursor.execute("SELECT homekey FROM info")
                ch.basic_publish(exchange='amq.topic',
                        routing_key=homeid,
                        body=AESencrypt(json.dumps(sent),
                        cursor.fetchone()[0]))
```

TABLE VI.    SOFTWARE IMPLEMENTATION FOR WARNING FROM RASPBERRY PI HOST

```
# COMMAND TO DEVICES        elif (sent['homeid'] == homeid and
sent['head']              ==              'CO')                :
    try:
        #CHECK                  LOCK                    STATUS
        cnx = mysql.connector.connect(user='root',
                password='28031995',
host='localhost',
    database='home'+homeid)
        cursor              =            cnx.cursor()
        cursor.execute("SELECT   lockstatus   FROM   info")

        #   SEND    WARNING    IF    HOUSE    IS   LOCKED
        if   (cursor.fetchone()[0]      ==     'lock')     :
            message             =                    {}
            message['head']            =               'WA'
            cursor.execute("SELECT  homekey  FROM  info")
            ch.basic_publish(exchange='amq.topic',
                    routing_key=homeid,
            body=AESencrypt(json.dumps(message),
                        cursor.fetchone()[0]))
        else :
            run_command(sent['address'],  sent['type'],
                    sent['command'])
            cursor.execute("SELECT  homekey  FROM  info")
            ch.basic_publish(exchange='amq.topic',
                    routing_key=homeid,

            BODY=AESENCRYPT(JSON.DUMPS(SENT),
                    CURSOR.FETCHONE()[0]))
                CURSOR.CLOSE()
                CNX.CLOSE()
                GC.COLLECT()
```

## B. Functional Test

The testing of the system is conducted on a door miniature (Fig. 9). To test the door, user is required to stand within determined distances with the door. In this case, the threshold is set to 10 meters from the door. According to the result as in Table VII, the door lock system that has been developed can work properly as expected. It can lock and unlock wirelessly.

Fig. 9.    A Setup for Functional Test of Door Lock.

TABLE VII.    POWER MEASUREMENT DOOR LOCK DEVICE

| Commands on GUI | Device Condition | Results |
|---|---|---|
| Lock | The door is locked | √ |
| Unlock | The door is unlocked | √ |

### C. Power Measurement

To measure the power consumption, simply we used digital multi-meter to know the current flow during two conditions (idle and process). The result of the test is elaborated in Table VIII. The power consumption in idle condition is 507.6 mW that obtained from 12 $V_{DC}$ * 42.3 mA, while in process condition is 7044 mW.

TABLE VIII.    POWER MEASUREMENT OF DOOR LOCK DEVICE

| Input voltage | Current | |
|---|---|---|
| | *Idle condition* | *Process condition* |
| 12 $V_{DC}$ | 42.3 mA | 587.5 mA |

### IV.    CONCLUSION

In this paper, a prototype of location-based smart door lock system is designed. The system utilizes the user's GPS coordinate that is captured from a mobile application, which is then sent to a smart home system's central host to enable or disable the door lock based on the user's proximity to the door's designated GPS coordinates. Based on the testing conducted, it can be concluded GPS coordinates can be used for controlling door lock. However, further study is required to improve the quality of the system, whether in terms of power efficiency, area tracking and indoor accuracy, and further increase the security.

### V.    FUTURE DIRECTION

For the next phase of the work, Android's geo-fencing implementation will be studied to examine its effectiveness on improving the lock's accuracy.

REFERENCES

[1]  B. Lorenzo, et al., "A Robust Dynamic Edge Network Architecture for the Internet-of-Things", arXiv preprint arXiv:1710.04861, 2017.

[2]  P.P. Ray, "A survey on Internet of Things architectures", J. of King Saud University-Computer and Information Sciences, Vol. 30(3), pp. 291-319, July 2018.

[3]  T.D.P. Mendes, et al., "Smart Home Communication Technologies and Applications: Wireless Protocol Assessment for Home Area Network Resources", Energies, Vol. 8(7), pp. 7279-7311, 2015.

[4]  R. Manjunatha and R. Nagaraja, "Home Security System and Door Access Control Based on Face Recognition", Int. Research J. of Engineering and Technology (IRJET), Vol. 4(3), pp. 437-442, 2017.

[5]  S. Kavde, et al., "Smart Digital Door Lock System using Bluetooth Technology", Int. Conf. on Information, Communication & Embedded Systems (ICICES) 2017.

[6]  N.A. Hussein and I. Al Mansoori, "Smart Door System for Home Security using Raspberry Pi 3", Proc. of the Int. Conf. on Computer and Applications, pp. 395-399, 2017.

[7]  C. Vongchumyen, et al., "Door Lock System via Web Application", Proc. of the 5th Int. Electrical Engineering Congress, pp. 1-4, March 2017.

[8]  Y.C. Yu, "A Practical Digital Door Lock for Smart Home", Proc. of the IEEE Int. Conf. on Consumer Electronics (ICCE), pp. 1-2, 2018.

[9]  S. Jensen, "Proximity Door Locking", Master Thesis, Technical University of Denmark.

[10]  T. Adiono, M.Y. Fathany, S.F. Anindya, S. Fuada, and I.G. Purwanda, "Wirelessly Control for RGB Lamp End-Device: Design and Implementation," IEEE Region 10 Conf. (TENCON), pp. 2066-2070, October 2018.

[11]  T. Adiono, S.F. Anindya, S. Fuada, and M.Y. Fathany, "Curtain Control Systems Development on Mesh Wireless Network of the Smart Home," Bulletin of Electrical Engineering and Informatics (BEEI), Vol. 7, No. 4, pp. 615-625, December 2018.

[12]  T. Adiono, M.Y. Fathany, S. Fuada, I.G. Purwanda, and S.F. Anindya, "A Portable Node of Humidity and Temperature Sensor for Indoor Environment Monitoring," Proc. of the 3rd Int. Conf. on Intelligent Green Building and Smart Grid (IGBSG), pp. 1-5.

[13]  T. Adiono, M.Y. Fathany, S.F. Anindya, S. Fuada, and I.G. Purwanda, "Development of wireless fan speed control using smartphone for smart home prototype," Unpublished.

[14]  T. Adiono, S.F. Anindya, S. Fuada, and M.Y. Fathany, "Developing of General IrDa Remote to Wirelessly Control IR-based Home Appliances," Proc. of the IEEE 7th Global Conf. on Consumer Electronics (GCCE), PP. 461-463, 2018.

[15]  T. Adiono, S.F. Anindya, S. Fuada, K. Afifah, and I.G. Purwanda, "Efficient Android Software Development using MIT App Inventor 2 for Bluetooth-based Smart Home," Wireless Pers Commun, Vol. 105(1), pp. 233-256, March 2019.

[16]  T. Adiono, B. Tandiawan, and S. Fuada, "Device Protocol Design for Security on Internet of Things based Smart Home," Int. J. of Online Engineering (i-JOE), Vol. 14(7), pp. 161-170, 2018. M.Y. Fathany and T. Adiono, "wireless protocol design for smart home on mesh wireless sensor network," Proc. of int. symp. On intelligent signal processing and communication systems, pp. 42-477, November 2015.

[17]  T. Adiono, et al., "Design of smart home mobile application with high security and automatic features," Proc. of the 2018 3rd Int. Conf. on Intelligent Green Building and Smart Grid (IGBSG), pp. 1–4, 2018.

[18]  T. Adiono, "Intelligent and secured software application for IoT based smart home," Proc. of the 2017 IEEE 6th Global Conf. on Consumer Electronics (GCCE), pp. 1-2, 2018.

# Low-Cost and Portable Ground Station for the Reception of NOAA Satellite Images

Antony E. Quiroz-Olivares[1], Ntalia I. Vargas-Cuentas[2], Guillermo W. Zarate Segura[3], Avid Roman-Gonzalez[4]

Image Processing Research Laboratory (INTI-Lab), Universidad de Ciencias y Humanidades, Lima, Peru

*Abstract*—**Currently, in Peru, the study of satellite images is increasing because it has the Earth observation satellite PeruSat-1. However, the cost of implementing a ground station is very high; for this reason, it is baffling that each university has its station. In the present work, the design and implementation of a low-cost portable earth station for the reception of meteorological satellite images is proposed in an automatic way, using accessible electronic devices such as Raspberry Pi 3b +, Software Defined by Radio (SDR) and an antenna double cross four dipoles, in this way encourage the study of satellite images in schools and universities. The results obtained show the viability of this project.**

*Keywords—Software defined by radio; Raspberry Pi; meteorological images; antenna; dipoles*

## I. INTRODUCTION

In Peru, the study of satellite images is in advance with the National Centre for Operations of Satellite Images (CNOIS for its Spanish acronym) that provides the images from the PeruSat-1 Earth Observation Satellite to different public and private entities. Every image distribution takes in to account the security and national interests. CNOIS belong to the Peruvian Space Agency (CONIDA for its Spanish acronym for National Commission of the Aerospace Research and Development National Commission). CONIDA promotes the development, research, and dissemination of space science and technology, contributing to the socio-economic development of the nation. The cost of the ground station for the PeruSat-1 was more or less $9.7 million. For the use of these satellite images, the interested entities must assign two employees as representatives of the institution, who have to manage the request and reception of CNOIS products. The reception of these images is carried out by a direct link from the ground station to the PeruSat-1 [1].

The PeruSat-1 is an Earth observation satellite acquired by the Peruvian government. It is used for multiple productive sectors as the mining sector, agricultural, civil defense, environmental, among others [2].

Regarding Earth observation satellite, there are also meteorological satellites that are in charge of supervising the atmospheric time and the displacement of the clouds around the world [3].

Meteorological satellites are separated into two types, those of fixed observation of territory that is the geostationary orbit and those that observe throughout the world in a day, they have a polar orbit and are free use [4].

Since 1978 the National Oceanic and Atmospheric Administration (NOAA) together with the National Aeronautics and Space Administration (NASA) of the United States have built a series of Environmental Satellites in Polar Orbit (POES), these satellites transmit images of the terrestrial or marine and atmospheric surface captured by their sensors in real time. Sending these measurements for climate modeling around the world used to anticipate natural disasters and weather conditions 3-7 days before they occur [6] [7].

The POES satellites transmit two types of images in real time, one by a direct link in the S-band, which is the High-Resolution Picture Transmission (HRPT) and the other by the VHF band which is the Automatic Picture Transmission (APT). Images from APT are the third part of the resolution of the HRPT; one currently has three satellites that transmit in APT, these are NOAA15, NOAA18 and NOAA19. The transmission of these images is given by the frequency of 137 MHz [4].

Currently, there are raspberry pi based systems for the reception of NOAA APT images, one of these systems is used manually as users who use it have to be present when passing NOAA satellites [5].

Currently, the implementation of a ground station for the reception of satellite images comes very expensive and in turn needs to have qualified personnel for its management, for this reason, universities do not have a ground station for the reception of satellite images.

In the present work, one proposes to implement a portable and low-cost system (approximately 120 dollars of investment), for obtaining meteorological satellite images - of APT type - automatically. The idea is to facilitate to users in the handling of a portable ground station, motivating to the study of the climatic patterns of its locality. For it, one will use a raspberry pi b+, a Software Defined Radio (SDR) and an arrangement of a double cross antenna with which we will receive the signal APT to later decode it and to store it in the computer for its later study.

In Section II, one describes the steps for the ground station implementation. Section III shows the obtained results. Section IV describes the discussion and the future perspectives of our ground station.

## II. METHODOLOGY

The operation of the proposed system can be seen in the following block diagram in Fig. 1.

Fig. 1.    Block Diagram of the Proposed System.

*A.  Space Segment*

The space segment is composed by the NOAA weather satellites. These satellites capture real-time images of the place where they pass and transmit it in HRPT and/or APT format [4].

*B.  Communication Module*

The system receives APT satellite images thanks to the double cross antenna of four dipoles Fig. 2, the design of the antenna is ideal for the reception of this type of signal; in addition, its components for its construction can be acquired in any hardware store [8].

Software Defined Radio (SDR) is responsible for tuning the frequency of the NOAA satellite signal; in Fig. 3 one see the SDR that controls all its physical functions by software [9], this device is used for educational purposes, as it is a low-cost system that can help understand the functioning of real communications.



Fig. 2.    Double Cross Antenna.



Fig. 3.    Software-Defined Radio.

*C.  Control Module*

The raspberry 3b+ Fig. 4, is responsible for controlling the SDR for the reception of satellite images, this is affordable minicomputer capable of running the applications needed for the portable and automatic ground station. Compared to previous versions, the raspberry 3b+ has increased in processing speed from 1.2 GHz to 1.4 GHz, also improved wireless connectivity, now works in dual band 2.4 GHz and 5 GHz [10].

*D.  Display Module*

For the visualization of the received images, it is necessary the installation of the following applications:

*1) Raspbian:* The Raspbian is an operating system based on Linux oriented to computers with few hardware resources such as the raspberry [11], in Fig. 5 one see the graphical environment of the Raspbian operating system.

*2) RTL-SDR:* It is an application that is responsible for the control of SDR, such as the frequency at which it must be tuned and the bandwidth needed to receive the information [12]. Table I shows the frequencies at which NOAA satellites transmit information.

*3) SOX:* It is a set of audio tools that allow us to manage the flow of received audio. This tool will enable us to record and store the audio captured by the SDR [13].

*4) Predict:* It is an open source program that provides satellite positioning in real time[1], the use of this application is only by console, unlike another similar program, the prediction does not need so much hardware resource, so it is perfect for this type of projects. In Fig. 6 one visualizes the initial configuration.

---

[1]https://www.qsl.net/kd2bd/predict.html
[2]https://wxtoimgrestored.xyz.

Fig. 4.   Raspberry pi 3b+.



Fig. 5.   Graphical Environment of the Raspbian Operating System.



Fig. 6.   Initial Configuration of the Predict Application.

TABLE I.        FREQUENCIES OF NOAA SATELLITES

| Satellite | Receiving frequency (MHz) | Wavelength (m) |
|-----------|---------------------------|----------------|
| NOAA 15 | 137.620 | 2.18 |
| NOAA 18 | 137.9125 | 2.175 |
| NOAA 19 | 137.100 | 2.188 |

*5) WXtoimg Restored:* It is a decoder of meteorological satellites, that allows us the edition and visualization of the received images in the form of audio, between its main features one can mention that it is possible to work with command lines to execute it automatically in the second plane[2]. In Fig. 7 one visualizes the graphical environment of the WXtoImg, this way of using the application is decided by the user.

*E.  Execution Process*

One installs the SDR USB drivers, then install the RTL-SDR application along with the SOX audio tools. To locate the satellites, one installs the Prediction program, where we have to place the coordinates where the portable ground station will be located, and finally one install WXtoimg restored, this program will be in charge of decoding the captured audio and converting it into a satellite image. With the applications already installed one proceed to create the scripts that will be responsible for doing the whole process automatically. The automatic system starts at midnight downloading a list of satellites that will pass during the day; the Predict application provides this information.

With the list of satellites NOAA15, NOAA18 and NOAA19 already programmed, the script calculates the time in which each satellite will pass in the day.

The script programs of the RTL-SDR as well as of the SOX will be executed in a precise instant that the satellite passes through our location.

The system will start tuning to the 137MHz frequency, beginning the recording of an audio file until the satellite is over.

After passing the satellite, the WXtoimg restored application will run, decoding the recorded audio file and converting it into an image, saving it in a raspberry pi b+ folder as shown in Fig. 8.



Fig. 7.   WXtoImg Application.

Fig. 8. Creation of Audio Files.

As soon as the NOAA satellite passes through our location, it sends the APT signal to Earth.

The four-dipole double-cross antenna receives the electromagnetic wave from the satellite and converts it into an electrical signal.

Software Defined Radio (SDR) demodulates the electrical signal from the antenna and converts it into a digital signal that is sent to the USB port of the raspberry pi b+.

The raspberry pi b+ decodes the information uploaded by the SDR and converts it into a satellite image.

### III. RESULTS

The implementation of the portable ground station has a total cost of 120 dollars, proving that it is economical compared to other ground stations, also because of its small size facilitates installation as seen in Fig. 9.

The obtained results with this system give us an average of 7 satellite images per day automatically. The system has demonstrated that it can emulate a ground station to get satellite images without presenting difficulties for the user in the operation of the equipment.

The obtained images have coverture of 1793km x 3200km; this allows us to observe most of our territory both land and sea.

The image quality depends on how clear the sky is when the satellite passes through our location. When there are clouds, one receives distorted images like Fig. 10; however, when the sky is completely clear, the received image is sharp, as seen in Fig. 11. Also, images obtained at night have the same results.



Fig. 9. Low-Cost and Portable Ground Station.



Fig. 10. Image with Cloudy Sky.



Fig. 11. Image with Clear Sky.

## IV. Discussion

In the present research, it was proposed as a primary objective to implement a low-cost and portable ground station system for automatically obtaining satellite images to motivate users in the study of the climatic patterns of its locality regarding the climate changes.

It is concluded that since this portable ground station is automatic, the user no longer has to wait for the satellite to pass through its location, thus facilitating the process of obtaining satellite images.

The portable ground station by its flexibility can be used to do image processing because the execution of applications responsible for the reception and decoding of satellite images work when passing satellites. This task could be used when no application is executed.

The idea is to continue with the research implemented by the system in colleges and universities to encourage students to study the analysis and forecasting of the climate in their region. Also, to prepare them for the future NOAA geostationary and polar orbit satellite systems.

The proposed low-cost and portable ground station could serve for education purpose in universities to do practices and laboratory experiments. It facilitates access to cheap laboratories for universities from developing countries.

As future work, one plans to add an orientation system for the reception of high-resolution images (HRPT), this will allow us to do image processing with high reliability.

### References

[1] CONIDA, "Suministro de imagenes satelitales del centro nacional de operaciones de imagenes satelitales," 2018.

[2] Roman-Gonzalez, Avid, and Natalia Indira Vargas-Cuentas. "Aerospace technology in peru." 66th International Astronautical Congress-IAC 2015. 2015.

[3] J. Gutiérrez Mendoza, O. Muñoz˘ Valverde, and J. Flores Quispe, "Capacidad del satélite´ perusat-1 en el desarrollo de inteligencia de imágenes´ en apoyo a la 3a brigada de artillería del iii ejército de operaciones," 2018.

[4] J. F. Moreno and J. Melia, "A method for accurate geometric correction of noaa avhrr hrpt data," IEEE Transactions on Geoscience and Remote Sensing, vol. 31, no. 1, pp. 204–226, 1993.

[5] Velasco, César, and Christian Tipantuña. "Meteorological picture reception system using software defined radio (SDR)." 2017 IEEE Second Ecuador Technical Chapters Meeting (ETCM). IEEE, 2017.

[6] G. K. Davis, "History of the noaa satellite program," Journal of Applied Remote Sensing, vol. 1, no. 1, p. 012504, 2007.

[7] A. Roman-Gonzalez, and N. I. Vargas-Cuentas; "Analysis of Landslides in Peru Based on Satellite Images to Identify Danger Zones", 69th International Astronautical Congress – IAC 2018; Bremen – Alemania; Octubre 2018.

[8] C. W. Alvarez Busani, "Diseño y construccion de una antena double cross para recepcion de imagenes procedentes de satelites de orbita polar," 2012.

[9] C. Velasco and C. Tipantu˜na, "Meteorological picture reception system using software defined radio (sdr)," in 2017 IEEE Second Ecuador Technical Chapters Meeting (ETCM). IEEE, 2017, pp. 1–6.

[10] B. Farroñan, A. Ivan, and J. L. Zuñiga Peña, "Diseño e implementacion de un timbre inteligente basado en el internet de las cosas (iot) para fortalecer la seguridad contra robos en viviendas sociales."2019.

[11] W. Harrington, Learning Raspbian. Packt Publishing Ltd, 2015.

[12] M. Sruthi, M. Abirami, A. Manikkoth, R. Gandhiraj, and K. Soman, "Low cost digital transceiver design for software defined radio using rtl sdr," in 2013 international mutli-conference on automation, computing, communication, control and compressed sensing (iMac4s). IEEE, 2013, pp. 852–855.

[13] G. T. LASKOSKI and P. NOHAMA, "Um novo software livre para comunicação alternativa".

# Experience in Asteroid Search using Astrometrica Software

Junior Ascencio-Moran[1], Jhon Calero-Juarez[2], Maria del Carmen Pajares-Acuña[3], Avid Roman-Gonzalez[4]

Image Processing Research Laboratory (INTI-Lab)
Universidad de Ciencias y Humanidades
Lima, Peru

*Abstract*—**The present work of this research consists of the analysis of telescopic images provided by the International Astronomical Search Collaboration (IASC) to find asteroids that can be named. The concern in searching for asteroids helps the scientific community that promotes the collaboration of young students and astronomy fans get experience in finding asteroids through campaigns related to these. The Space Generation Advisory Council (SGAC) campaign in partnership with IASC has found around 1500 asteroids since the beginning of October 2006 as each year more than 1000 teams from different countries participate. The Astrometrica software was used which is in charge of receiving the images in FITS format. The configuration of the selected telescope is carried out so that they are later analyzed in greater detail. Finally, a clean and precise Minor Planet Center (MPC) report is made, which is what the campaign requires so that it can go on to a preliminary phase and subsequently be accepted by the international astronomical union. The asteroids that become named will be registered in the catalog of official minor planets of the world. In the campaign related to this study, one finds 28 possible asteroids.**

*Keywords—SGAC; IASC; INTI-Lab; UCH; MPC; asteroids*

## I. INTRODUCTION

Currently, the study and monitoring of asteroids in our Solar System represent a primary need since these celestial bodies are potential threats for the Earth. If they collide with our planet, this will imply, depending on the dimensions and internal structure of the asteroid, disasters of enormous proportions until the extinction of life on Earth [1]. On the other hand, asteroids contain a large amount of information because most of these were formed along with our Solar System and even some precede the formation of our planetary system. So, it is essential to learn from them through observations and research because these minor planets can provide us with information on the origin of evolution and life itself [2][3].

A comet is one minor solar system body, made up of a nucleus "dirty ice" that dirty snowball which we call Comet, which is made up of CH, CN y NH2. When approaching the Sun given the material that makes comets, the evaporation of these compounds is sublimated, so it gives rise to a gaseous sphere called coma and dust particles and ions (triggered by the solar wind) that form the tail. This tail has a length of millions of kilometers [4].

A meteorite is the wake of light left by a meteoroid when it enters into Earth's atmosphere. The purpose of this work is to raise awareness of the importance of participating in an asteroid search campaign, in addition to telling the experience and how the campaign is developed week by week until the culmination of it, presenting real data provided by the International Astronomical Search Collaboration. Particles, scattered throughout space, are grouped in swarms and associated with a comet. Meteors are usually more visible from Earth than asteroids or comets. They are commonly known as shooting stars [4].

Asteroids are small rocky bodies compared to the planets of our solar system. Asteroids have irregular shapes, and only a small part has a spherical shape, their surfaces are full of craters or holes. The vast majority of asteroids are concentrated in the "Asteroid Belt" which is located between the Mars and Jupiter planets, also have an orbit around the Sun [5][6].

Asteroids can be discovered and observed with the help of optical telescopes by astronomers and amateurs who are passionate about asteroids and outer space. Currently, some organizations promote initiative and interest in knowledge and research; also offer the possibility of making their discoveries by participating in asteroid search campaigns.

Prestigious organizations, such as the Space Generation Advisory Council (SGAC) and the International Astronomical Search Collaboration (IASC) works together to promote space science. The Space Generation Advisory Council (SGAC) is a worldwide organization representing students and young professionals involved and interested in the space industry. The International Astronomical Search Collaboration (IASC) is a citizen science program that provides high-quality astronomical data to citizen scientists from around the world.

These citizen scientists can make original astronomical discoveries and participate in practical astronomy, that they regularly sponsor these asteroid search campaigns where only a few teams are selected from multiple schools and participants from various parts of the world. The experience of participating in this campaign also gives participating organizations the privilege of assigning a name to their discoveries.

Since the beginning of October 2006, more than 1500 asteroids have been discovered, of which 52 have been numbered by the International Astronomical Union [7]. The numbered asteroids are recorded in the catalog of official minor planets of the world. This task could be possible with the help of Astrometrica software that allows the superposition and

analysis of images and the ability to easily compare astronomical captures to promote the discovery of objects [8].

Many of these experiences are recorded in documents describing how the process has been, as presented in this document "in search of new asteroids or Astronomy Education Review"[9] [10].

The purpose of this work is to publicize the importance of participating in an asteroid search campaign, also, to tell the experience and how the campaign develops week by week until the culmination of it, presenting real data provided by the International Astronomical Search Collaboration (IASC).

In the second part, the selection process determined by Space Generation Advisory Council (SGAC) to designate the participating teams will be addressed. It will also explain how the preparation process is developed for optimum performance during the period of the campaign, as well as indicating the characteristics of the materials necessary for this process of preparation and development of the campaign. In the third part, it will be presented how the campaign was developed during the 30 days of its duration, the results obtained from experience and how they contribute to the scientific community that is dedicated to working related to the space industry.

## II. METHODOLOGY

To be part of the search asteroid campaign is indispensable, in the first instance, go through a selection process and, subsequently, a period of preparation before the start of the event. This section will be followed by a detailed presentation of how these procedures were developed and how they were tackled, to obtain the privilege of participating in the "SGAC Asteroid Search Campaign.

### A. Selection Process

The announcement for the selection of participants of the campaign "SGAC Asteroid Search Campaign," was launched in early October 2018 by the Space Generation Advisory Council, hence the name of the search campaign. The Space Generation Advisory Council (SGAC) raised four requirements, from which teams would be chosen to form part of the asteroid search campaign. The following criteria are indicated:

- Teams should be created with a minimum of 3 and a maximum of 5 members. In the case of participants who did not have a team, there was the possibility of being associated with other participants who were in the same condition to form a team with the established requirements.

- Each team member should be registered as a Space Generation Advisory Council (SGAC) member.

- Each team member had to answer the following questions, explaining the reasons that encouraged to participate:

- Why do you want to be part of this campaign?

- How do you think that participation in the Asteroid Search Campaign could benefit you?

- Participants ' registrations should be made before the deadline established (21 October 2018).

Teams were selected based on the individual responses of each member and their regional distribution in the globe.

Based on our capacity as an enthusiastic student on research, knowledge, discoveries and countless hidden mysteries, samples of our answers to the questions set by Space Generation Advisory Council (SGAC) are shown below:

- The outer space has always been a subject of my interest because I consider that there lies the response of the origin of life and the universe. I find also exciting the opportunity to discover an asteroid by myself and to know a little more carefully about asteroids in detail. Also, I consider it an excellent opportunity to develop a new academic experience.

- I believe that participating in this campaign is essential because it represents a unique opportunity and will allow me to acquire new knowledge related to the asteroids in general. Also, have access to real data, analyze them and be able to find at least one asteroid and be able to give it a name would be rewarding.

The outer space is a privileged place and therefore deserves dedication and commitment, also, every minimum possibility of being able to get involved and contribute directly or indirectly to issues related to the space industry. It's comfortable.

As seen in the Fig. 1, the selected participants were: Calero Juarez Jhon David, Ascencio Moran Junior Angel, and Pajares Acuña Maria del Carmen under the supervision of Prof. Ph.D. Ing. Avid Roman-Gonzalez.

As seen in Table I, the selected team, which subsequently, they are notified via email on October 26; after that, each team will have four days to perform their training before asteroid hunting.

In this selection process, 17 countries participated, and there were a total of 24 teams selected.



Fig. 1. Students for the Asteroid Search Campaing.

TABLE I.        PARTICIPATING TEAMS

| # | Team | Country |
|---|------|---------|
| 1 | Space Girls | Belarus |
| 2 | Universidad Mayor de San Andres | Bolivia |
| 3 | Municipal center for extracurricular activities | Bulgaria |
| 4 | Candor Chasma Team | Colombia & Peru |
| 5 | Asteroid Caretakers Squad | Colombia-Bolivia |
| 6 | Ryan's Belt | Czech Republic/Slovakia/United Kingdom |
| 7 | YEHA | Ethiopia |
| 8 | The Harwoods | France |
| 9 | Aristotle University of Thessaloniki | Greece |
| 10 | SAKED in Space | Greece |
| 11 | STAR | Greece |
| 12 | Gagan | India |
| 13 | Chaitanya Mandala | Italy |
| 14 | Rock Hunters Trio - Bologna University | Italy |
| 15 | Astro Cypriots - Bekirpasa Lycee | North Cyprus |
| 16 | INTI-Lab - Universidad de Ciencias y Humanidades | Peru |
| 17 | AstropILO | Philippines |
| 18 | Institute Of Astronomy Sri lanka | Sri Lanka |
| 19 | Sky Trailblazers | Sri Lanka |
| 20 | ULKA | Sri Lanka |
| 21 | University of Khartoum | Sudan |
| 22 | Azteroids | USA |
| 23 | Paco to the Stars | USA |
| 24 | HSU | USA |

## B. Training Process

Each team receives a unique set of images that they must then study using the Astrometric software, then prepare a report according to the Minor Planet Center (MPC) that will be sent to Dr. Patrick Miller via email. The report will then be reviewed by the International Astronomical Search Collaboration (IASC) Data Reduction Team (IDaRT) for possible discoveries and reported to the Minor Planet Center (MPC) (Harvard).

The images analyzed are provided by the Institute for Astronomy (IFA) at the University of Hawaii. These images are in FITS (Flexible Image Transport Setting) format and are captured by the 1.8m Pan-STARRS telescope located in Haleakala.

The training begins as soon as the acceptance email is received, this email contains a series of instructions, a user name and a password necessary for the development of the campaign, also includes the email of five experts in the management of the software Astrometrica, who will resolve possible doubts in this regard.

Each team member is assigned a license and a password (well apart from the user password) for unlimited use of the Astrometric software. The software can be downloaded for free from the official IASC website. Also, it is possible to download a quick guide to use the software and a set of images to practice during the four days before the start of the campaign.

The success of each participating team depends on being able to use the software effectively, prepare accurate MPC reports and distinguish between right asteroids and false targets.

## C. Materials and Equipments

To develop the search for asteroids, the Universidad de Ciencias y Humanidades through the Image Processing Research Laboratory (INTI-Lab) provided us with a laboratory with the following characteristics:

A Core i7 6th generation computer with 16RAM DDR 4, so that our team does not have difficulties when installing the program Astrometrica.

Thanks to this collaboration, we were able to carry out the asteroid search without any complications.

## D. Asteroid Search with Astrometrica Software

During the campaign, the Astrometrica software will be used, which can be accessed for free and even provide a manual from its official website, but only the participants of this campaign will be given a user name and password with which they will receive previously after the preparation process.

To develop this experience, it is necessary to download the set of instructions provided by the campaign.

The following steps must be followed:

Step 1: As seen in the Fig. 2 unzip the archive of satellite images and place them in a folder (previously created) with the same name as the downloaded file, then open the Astrometric software and wait a few minutes until it is finished initializing (two windows will open, which must remain open throughout the process of analyzing the images).

Step 2: As seen in the Fig. 3, it is verified that the telescope (Pan-STARrs-2) selected in the software is the same as the one used to capture the images.

Step 3: As seen in the Fig. 4, one proceed to load the Pan-STARRS-2 telescopic images (all images in the folder must be selected).



Fig. 2.    Minor Planet Center (MPC) Database Update Window.

Fig. 3. Window of the Software Astrometrica, Pan-STARRS-2.



Fig. 5. Window of the Astrometrica Software, Stars Present in the Images of the Pan-STARRS-2 Telescope.



Fig. 4. Astrometric Software Window, Pan-STARRS-2 Telescope Images.



Fig. 6. Moving Objects in the Pan-STARRS-2 Telescope Images.

Step 4: As seen in Fig. 5, one proceed to load the Pan-STARRS-2 telescopic images (all images in the folder must be selected).Select in the menu bar the icon "Astrometric Data Reduction" to highlight the stars present in the images.

Step 5: As seen in Fig. 6, select in the menu bar the icon "Known Object Overlay" to highlight the objects that have movements in the image. Then select the "Blink current images" icon so that the software can create a small gif by placing the images in succession and thus identify the moving objects present in the set of images.

Step 6: In the new window a kind of "gif" is shown where we will proceed to look for objects that present movement. In the case of locating a moving object, we should pause the "blinked" and place the pointer as close as possible to the center of the object. After this, a new window will open with the characteristics of the selected object.

As seen in Fig. 7, one must take into account the following considerations to recognize the right asteroid:

- The white dots displayed in the verification window should be distributed very close to the red line.

- The signal to noise ratio (SNR) must be greater than 5.

- The declination and right ascension of the object should be very close to 0.00.

- The selected object must maintain a straight line movement.



Fig. 7. Astrometric Parameters.

```
COD F51
OBS N. Primak, A. Schultz, S. Watters, J. Thiel, T. Goggia
MEA P. Miller, C. Davis, & D. Offner (Hardin-Simmons University, USA)
TEL 1.8-m f/4.4 Ritchey-Chretien + CCD
ACK MPCReport file updated 2016.10.31 13:36:23
NET PPMXL
    HSU0001 C2016 09 12.51415300 08 13.441+11 15 46.81      20.6 R      F51
    HSU0001 C2016 09 12.52755500 08 12.706+11 15 46.34      19.7 R      F51
    HSU0001 C2016 09 12.54105000 08 11.972+11 15 45.44      20.5 R      F51
    HSU0001 C2016 09 12.55458200 08 11.241+11 15 44.56      20.4 R      F51
    U4671       C2016 09 12.51415300 08 13.298+11 11 13.98      20.9 R      F51
    U4671       C2016 09 12.52755500 08 12.536+11 11 13.71      20.9 R      F51
    U4671       C2016 09 12.54105000 08 11.727+11 11 13.62      20.9 R      F51
    U4671       C2016 09 12.55458200 08 10.940+11 11 13.38      20.9 R      F51
    ----- end -----
```

```
COD F51
OBS N. Primak, A. Schultz, S. Watters, J. Thiel, T. Goggia
MEA P. Miller, C. Davis, & D. Offner (Hardin-Simmons University, USA)
TEL 1.8-m f/4.4 Ritchey-Chretien + CCD
ACK MPCReport file updated 2016.10.31 13:36:23
NET PPMXL

Image set: ps1-20161003_2_set084

    HSU0001 C2016 09 12.51415300 08 13.441+11 15 46.81      20.6 R      F51
    HSU0001 C2016 09 12.52755500 08 12.706+11 15 46.34      19.7 R      F51
    HSU0001 C2016 09 12.54105000 08 11.972+11 15 45.44      20.5 R      F51
    HSU0001 C2016 09 12.55458200 08 11.241+11 15 44.56      20.4 R      F51

    U4671       C2016 09 12.51415300 08 13.298+11 11 13.98      20.9 R      F51
    U4671       C2016 09 12.52755500 08 12.536+11 11 13.71      20.9 R      F51
    U4671       C2016 09 12.54105000 08 11.727+11 11 13.62      20.9 R      F51
    U4671       C2016 09 12.55458200 08 10.940+11 11 13.38      20.9 R      F51

    ----- end -----
```

Fig. 8. MPC Report Model.

TABLE II. PARTICIPATING TEAMS AND DISCOVERIES

| # | Team | Country | Discoveries |
|---|------|---------|-------------|
| 1 | Space Girls | Belarus | 1 |
| 2 | Universidad Mayor de San Andres | Bolivia | 7 |
| 3 | Municipal center for extracurricular activities | Bulgaria | 6 |
| 4 | Candor Chasma Team | Colombia & Peru | 12 |
| 5 | Asteroid Caretakers Squad | Colombia-Bolivia | 5 |
| 6 | Ryan's Belt | Czech Republic/Slovakia/United Kingdom | 7 |
| 7 | YEHA | Ethiopia | 20 |
| 8 | The Harwoods | France | 22 |
| 9 | Aristotle University of Thessaloniki | Greece | 10 |
| 10 | SAKED in Space | Greece | 7 |
| 11 | STAR | Greece | 3 |
| 12 | Gagan | India | 5 |
| 13 | Chaitanya Mandala | Italy | 6 |
| 14 | Rock Hunters Trio - Bologna University | Italy | 11 |
| 15 | Astro Cypriots - Bekirpasa Lycee | North Cyprus | 6 |
| 16 | INTI-Lab - Universidad de Ciencias y Humanidades | Peru | 28 |
| 17 | AstropILO | Philippines | 8 |
| 18 | Institute Of Astronomy Sri lanka | Sri Lanka | 8 |
| 19 | Sky Trailblazers | Sri Lanka | 6 |
| 20 | ULKA | Sri Lanka | 16 |
| 21 | University of Khartoum | Sudan | 13 |
| 22 | Azteroids | USA | 9 |
| 23 | Paco to the Stars | USA | 4 |
| 24 | HSU | USA | 1 |

It is clear to take into account these parameters because many times it can be confused with hot pixels, background fluctuations with asteroids. Not all things that seem to move in the images will be asteroids.

Step 7: As seen in Fig. 8, once the analysis of the set of images is finished, a report of the discoveries must be made. This report is quite simple because the software elaborates it automatically. The only thing that we must do is to copy the MPC report in a notepad (it must take as name the set of analyzed images) that later we will send to an e-mail to Dr. Patrick Miller, who will be in charge of confirming or discarding a newly discovered asteroid.

## III. RESULTS

As can be seen in Table II, during the asteroid search campaign October-November 2018, the results of each participating team are shown.

To have found this amount of asteroids in a single team representing Peru is very gratifying because as we observed in other participating countries, there were 2 or 3 teams for each country that means that there is a greater interest in discovering what is in space. When one finds an asteroid brings with it a lot of information, as there is a greater interest in a country, it is possible to get a higher amount of participating teams, and that implies a greater amount of discoveries.

One can see in Table III the number of countries and the total number of asteroids discovered in the preliminary phase, each year, different countries join this campaign, making their interest in exploring space more evident.

In Fig. 9, one can see the report that is sent each week to Dr. Patrick Miller. In the period of the campaign, it was received a set of weekly images, each set had 4 pictures, so in total was obtained a total of 20 images in the search stage. During this period 48 possible asteroids were sent in the report.

TABLE III. COUNTRIES AND THEIR DISCOVERIES

| # | Country | Discoveries |
|---|---------|-------------|
| 1 | Sri Lanka | 30 |
| 2 | Perú | 28 |
| 3 | France | 22 |
| 4 | Greece | 20 |
| 5 | Ethiopia | 20 |
| 6 | Italy | 17 |
| 7 | USA | 14 |
| 8 | Sudan | 13 |
| 9 | Colombia & Perú | 12 |
| 10 | Philippines | 8 |
| 11 | Bolivia | 7 |
| 12 | Czech Republic/Slovakia/United Kingdom | 7 |
| 13 | Bulgaria | 6 |
| 14 | North Cyprus | 6 |
| 15 | India | 5 |
| 16 | Colombia-Bolivia | 5 |
| 17 | Belarus | 1 |

Fig. 9. Minor Planer Center(MPC) Report Model of New Discoveries Report.

According to the results of the International Astronomical Research Collaboration (IASC), as can be seen in the Fig. 10, who had the most considerable number of discoveries was Sri Lanka with its three teams together achieved 30 accepted asteroids in the preliminary phase, then representing Peru our team INTI-Lab - Universidad de Ciencias y Humanidades who obtained 28 accepted asteroids in the initial stage.

These data are beneficial to determine the orbits of objects with greater accuracy or even, as in our case, to make discoveries of new asteroids.



Fig. 10. SGAC Campaing Discoveries.



Fig. 11. Certificate for having Participated in the Asteroid Search Campaign..

At the end of the campaign, participants receive an email from IASC which sends them a certificate via email, as can be seen in the Fig. 11, thanking them for their time, effort and dedication during the search stage.

This opportunity to make discoveries and contribute to the goals of protecting Earth from asteroid collisions and exploring new and/or potentially dangerous asteroids is a unique experience, which can be voluntarily accessed.

## IV. Discussion

To send us the set of images from the Pan-STARRS telescope located in Haleakala, Hawaii has to be completely clear otherwise it will not be able to send the collection of pictures to our team during the first four days of having started the campaign. Due to Hawaii, where the telescope that provided the photos is located, had storms with clouds; it was impossible for the telescope to take images for sending us to analyze them.

During the four weeks in which the report of the 48 possible asteroids found was sent, from which 20 steroids were rejected, 28 have been accepted in the preliminary phase.

The Minor Planet Center requires from 7 to 10 days to analyze the preliminary discovery and to confirm its existence and better establish its orbit [10].

Now the nomination stage takes between three and six years since the Minor Planet Center (MPC) updates the preliminary discovery as an interim discovery and monitors additional discovery observations until the orbit has been determined.

Under the direction of instructors, participating students have successfully discovered new asteroids in the preliminary phase of the Main Belt, this is the result of constant motivation and dedication that any person, student or amateur can achieve, through this work motivates people who can participate in future campaigns.

### References

[1] Anguita, F., & Castilla, G. (2005). Crónicas del sistema solar. Recovered from: https://ebookcentral.proquest.com.

[2] Carrasco, L. E., & Carramiñana, A. A. (2005). Del sol a los confines del sistema solar. Recovered from: https://ebookcentral.proquest.com.

[3] Equipo, S. (Ed.). (2009). Astronomía nº 144. Recovered from : https://ebookcentral.proquest.com.

[4] Jorge, R & Ángel, G. (2008). Astronomía contemporánea (3a. ed.). Recovered from: https://ebookcentral.proquest.com.

[5] Fierro, J., & Herrera, M. Á. (1988). La familia del sol. Recovered from: https://ebookcentral.proquest.com.

[6] Agencia Espacial del Perú CONIDA (2011), Asteroides Recovered from: http://www.conida.gob.pe.

[7] Space Generation Advisory Council. (2018). SGAC Asteroid Search Campaign. Recovered from: https://spacegeneration.org.

[8] IASC International Astronimical Search Collaboration. Astrometrica. Recovered from: http://iasc.hsutx.edu.

[9] J. Patrick Miller, Jeffrey W. Davis, Robert E. Holmes, Jr. Harlan Devore, Herbert Raab, Carlton R. Pennypacker. (2008). Internet-Based Hands-On Research Program for High Schools andColleges, in Collaboration with the Hands-On Universe Project en Astronomy Education Review vol 7.

[10] USA HOU (2018). IASC participate. Recovered from: http://handsonuniverse.org.

# Size Reduction and Performance Enhancement of Pi Shaped Patch Antenna using Superstrate Configuration

Pir Saadullah Shah[1], Shahryar Shafique Qureshi[2], Muhammad Haneef[3], Sohail Imran Saeed[4]

Electrical Engineering Department, Iqra National University, Peshawar, Pakistan[1, 2, 4]
Electrical Engineering Department, Foundation University Islamabad, Islamabad, Pakistan[3]

*Abstract*—Patch antennas are modern elements of today's world communication technology. They appear to have unique characteristics and features with their unique power with handling capabilities and lighted structure. This paper focuses on superstrate configuration of patch antenna with defected ground plane and Pi Slotted radiating patch. The three different cases were taken in terms of wavelength distance to observe the performance characteristics of patch structure. The antenna designed in this study can be used for S and C band applications.

*Keywords*—*Radiating patch; ground plane; slotted; bandwidth; gain; directivity*

## I. INTRODUCTION

Modern technology has seen a rapid growth in communication technology sector. Since from beginning of wireless communication the wireless spectrum technology has spread worldwide. Antenna technology has always been taken as keen interest by engineers and researchers of industry. Parabolic reflectors, dipole, monopole antennas almost in every field antenna technology has seen marvelous growth in communication sector. IEEE protocols have set almost bands for every communication of application purposes. ISM band deals with wearable antenna characteristics and with on its boundary S band connected which deals with GSM Wimax and Wireless technologies also serving Wi-Fi and space suit communication for astronauts and deep space missions [1-3].

This research presents a novel research of antenna on its size reduction. Patch antenna have been a prominent research figure in antennas researchers due to its unique capabilities. Patch antennas are low profile antennas that have power handling capabilities. When assembled together to form an array, these small elements can be boosted in terms of performance parameters and have also been used for deep space missions although too little spacing among these elements can give rise to coupling issue that's interaction of their magnetic fields which is improved through proper isolation techniques [4-6]. Patch antenna size can be reduced through number of techniques like Artificial Magnetic Conductors [7-8] use of Electromagnetic Band Gap Structures, EBGs [9-10] and Meta materials. However designing such antennas is little tough as they give relatively low performance parameter results and also miniaturization is not extended to big extent. Meta materials show goof approach but as they are designed individually they are extremely rare and cost of

making them can differ up to great extent [11-12]. Slots can be made on the ground plane and radiating patch and have no design mechanisms and have been reported well to perform in some cases [13-15]. As some times antenna radiation and performance parameters may get distorted, they can be changed to satisfactory levels by stack and superstrate configuration [16-17].

In this research an antenna size is reduced through introduction of slots in its radiating patch and its full ground plane. After miniaturization, antenna main performance parameter which is reflection co efficient, is improved with help of Superstrate configuration alongside its other performance parameters. The superstrate configuration consists of a radiating patch faced parallel to another patch also known as parasitic patch. These patches are coupled magnetically and then are enhanced in terms of their performance characteristics.

The paper is organized as follow. Introduction comes first then Antenna design. Results and discussions covers all the performance parameters and their analysis and in the last comes Conclusion and future work

## II. ANTENNA DESIGN

The resonance frequency for any mn of a rectangular MSA, is given by as:

$$f_0 = \frac{c}{2\sqrt{\varepsilon_{edc}}} \left[ \left(\frac{m}{L}\right)^2 + \left(\frac{n}{W}\right)^2 \right]^{\frac{1}{2}} \qquad (1)$$

The m and n are modes along L and W. For efficient radiation of patch, the width W is given by:

$$W = \frac{c}{2f_0 \sqrt{\frac{\varepsilon_r+1}{2}}} \qquad (2)$$

The antenna permittivity plays an important role in its design and also in performance parameter terms. Substrate thickness also has keen role in antenna design as thick substrate confines field in it.

Antenna with higher permittivity substrate always has low gain but bandwidth is better and antenna with lower permittivity substrate and with higher thickness has good gain results but bandwidth is less. In our design FR4 substrate is taken with relative permittivity of 4.4. The basic design of patch antenna is shown in Fig. 1 with co axial cable feed.

461 | P a g e

Fig. 1. Patch Antenna Layout.

As discussed earlier this research focuses on miniaturization and performance characteristics enhancement using superstrate configuration, the basic layout of superstarte configuration is shown in Fig. 2.



Fig. 2. Superstare Configurated Patch Antenna.

The patch antenna of 4.5GHz was designed and was introduced to slots and ground irregularities in both radiating patch and ground plane. The patch dimensions with help of equations mentioned in [18] were taken to be 20mm width and 15mm in length. The ground plane is twice the size of the calculated patch.



Fig. 3. Proposed Antenna (a) Radiating Patch (b) Ground Plane.

On the radiating patch, the fractional shape is made and length of $SL_1$ and $SL_2$ are taken to be 12mm and $SW_1$ and $SW_2$ are taken to be 5mm. $SL_3$ and $SL_4$ are taken to be 5mm and $SW_3$ and $SW_4$ are taken to be 12mm. The Pi slot is made with lengths of 5 mm and thickness of 2mm each side. On the ground plane, the U and L slots are made the length and width of the both slots are taken to be 10 mm and width of each and 4mm wide. The proposed geometry can be seen in Fig. 3(a) and (b) which represents patch and ground plane respectively. Through the slots, 72% miniaturization is achieved as proposed

antenna current patterns are altered through slots which resulted in size reduction and multiband response.

These results are further then enhanced by implementing superstarte configuration through distance of 12, 24 and 36mm apart. The basic distance is taken to be 12mm distance which represents half wavelength distance.

III. RESULTS AND DISCUSSIONS

The reflection co efficent parameters are shown in Fig. 4. The part (a) shows the basic graph in which antenna is clearly seen to be operating at 4.5GHz. The antenna after miniaturization showed multiband response and when it was introduced to superstrate configuration, the proposed antenna results were enhanced and s parameter response from tri band was changed to quad band.



Fig. 4. Proposed Antenna S Parameters(a) Conventional Patch (b) 12mm Apart (c) 24mm Apart (d) 36mm Aprt.

It can be seen that the basic S parameters were boosted at 24mm apart distance which summed up to be clearly one wavelength distance between radiating and parasitic patch. The S parameters were enhanced down to -10db further and bandwidth of the proposed antenna was also increased. The max gain over frequency is shown in Fig. 5.

It can be clearly seen that as duplicated in part (a) the max gain of the miniaturized patch showed low gain outputs as compared to superstrate composed gain output results. For the sake of simplicity, the gain of 24mm is shown as it offered highest enhanced results.

Fig. 6 and Fig. 7 show 3 dimensional gain and 1 Dimensional polar plots of case 2 which is 24mm apart configuration. From the patterns it can be seen that the proposed antenna is operating at nearly all directions in quad band response.

The detailed analysis of superstrate configuration is mentioned in Tables I to III. From the table it can clearly be seen that 24mm apart distance showed better results. In all the cases the proposed antenna is well matched and is delivering maximum power as accordingly to maximum power transfer power theorem. Table IV shows the size reduction comparison at resonating frequencies.



Fig. 6. Gain 3D Radiation Patterns. (a) 3.06GHz (b) 3.60GHz (c) 3.70GHz (d) 4.5GHz.



Fig. 7. Gain 1D Radiation Patterns. (a) 3.06GHz (b) 3.60GHz (c) 3.70GHz (d) 4.5GHz.



Fig. 5. Max Gain Over Frequency(a) Minituraized Patch (b) Miniaturaized Patch with Superstrate Configuration.

TABLE I. PERFORMANCE CHARACTEISTICS OF CASE 1

| Antenna Parameter | Simulated Value | Simulated Value | Simulated Value | Simulated Value |
|---|---|---|---|---|
| Resonant Frequency | 3.06 GHz | 3.60 GHz | 3.7 GHz | 4.5 GHz |
| Return Loss (S11) | -14.33dB | -15dB | -15.3dB | -13.88dB |
| VSWR | 1.12 | 1.22 | 1.25 | 1.30 |
| Gain | 4.12dB | 2.12dB | 2.39dB | 2.39dB |
| Directivity | 6.46dBi | 4.3dBi | 3.3dBi | 3.2dBi |
| Bandwidth | 110MHz | 60MHz | 175MHz | 103MHz |

TABLE II.  PERFORMANCE CHARACTEISTICS OF CASE 2

| Antenna Parameter | Simulated Value | Simulated Value | Simulated Value | Simulated Value |
|---|---|---|---|---|
| Resonant Frequency | 3.06 GHz | 3.60 GHz | 3.7 GHz | 4.07 GHz |
| Return Loss (S11) | -34.55dB | -36.00dB | -33.90dB | -11.30dB |
| VSWR | 1.22 | 1.122 | 1.325 | 1.40 |
| Gain | 5.12dB | 2.42dB | 4.19dB | 2.69dB |
| Directivity | 6.46dBi | 4.3dBi | 6.3dBi | 5.2dBi |
| Bandwidth | 130MHz | 250MHz | | 110MJz |

TABLE III.  PERFORMANCE CHARACTEISTICS OF CASE 3

| Antenna Parameter | Simulated Value | Simulated Value | Simulated Value | Simulated Value |
|---|---|---|---|---|
| Resonant Frequency | 3.06 GHz | 3.60 GHz | 3.7 GHz | 4.5 GHz |
| Return Loss (S11) | -26.33dB | -27dB | -25.3dB | -13.88dB |
| VSWR | 1.12 | 1.22 | 1.25 | 1.30 |
| Gain | 4.12dB | 1.42dB | 3.19dB | 2.00dB |
| Directivity | 5.20dBi | 1.98dBi | 4.3dBi | 2.2dBi |
| Bandwidth | 108MHz | 252MHz | | 108MHz |

TABLE IV.  SIZE REDUCTION OF PROPOSED PATCH

| S.No | Antenna (GHz) | Dimensions mm² | Size Reduction |
|---|---|---|---|
| 1 | 3.06 | 1131 | 72% |
| 2 | 3.6 | 850 | 52% |
| 3 | 3.7 | 780 | 47% |

## IV. CONCLUSION

Patch Antennas play an important role in up to date communication technology. These minute elements can be designed at much ease as compared to other traditional antennas. They are frequency reconfigurable and easily adjustable elements. Patch Antennas offer a unique characteristics of multi band response when introduced with slots. However at every resonant frequency, their size varies and so the slots have no specific criteria either they should be placed on top or right or any other side. In this research a miniature antenna was presented with multi band response. Different techniques were considered in reducing the size of antenna and size reduction through slots was implemented. It was seen that with adding different slots on the large ground plane and radiating patch antenna exhibited Quad band response with resulting the size reduction of 72%. Antenna with slots can offer low results as slots can effect antenna performing parameters through which antenna efficiency is affected. Their VSWR plots or impedance matching values are effected if slots are not planed and examined at each state and with higher VSWR values can horribly affect antenna performance. Superstrate configuration was implemented after reducing the size of antenna and a parasitic patch was placed at top of radiating patch. The parasitic patch was studied on different distance levels and at half wavelength distance antenna performance parameters showed satisfactory levels of performance. The reflection co efficient is main parameter of antenna and it showed enhanced results as compared to single radiating patch. These results can further be improved in future if this set up can be performed in a cavity box. By placing it in a cavity box we can improve its radiating patterns further and also can examine its gain and efficiency values. , this technique has been previous proved helpful for narrow-band antennas. The technique can be validated by applying it to band notch version of the same antenna. Also this results can be checked for their behaviour response if the number of parasitic patches can be increased and their structure as it can lead to wider band response if the parasitic patch is used without addition of slots and can be used as a reflector. This set up was performed in a professional antenna simulator software Computer Simulation Technology 2014 and mesh size per wavelength was kept at 5 levels. The proposed antenna is miniaturized as compared to conventional antenna in size and has shown good bandwidth and gain alongside directivity and other performance parameters. The proposed antenna can be used in S band applications and Wi-Max and other small wireless technologies.

## REFERENCES

[1] Chen, Zhi Ning, and Xianming Qing. "Dual-Band Circularly Polarized $ S $-Shaped Slotted Patch Antenna With a Small Frequency-Ratio." IEEE Transactions on Antennas and Propagation 58.6 (2010): 2112-2115.

[2] Saad Hassan Kiani, Khalid Mahmood and Ahsan Altaf, "A Linear Array for Short Range Radio Location and Application Systems" International Journal of Advanced Computer Science and Applications(IJACSA), 9(4), 2018. http://dx.doi.org/10.14569/IJACSA.2018.090420

[3] Christodoulou, Christos G., et al. "Reconfigurable antennas for wireless and space applications." Proceedings of the IEEE 100.7 (2012): 2250-2261.

[4] Saad Hassan Kiani, Khalid Mahmood, Ahsan Altaf and Alex J. Cole, "Mutual Coupling Reduction of MIMO Antenna for Satellite Services and Radio Altimeter Applications" International Journal of Advanced Computer Science and Applications(IJACSA), 9(4), 2018. http://dx.doi.org/10.14569/IJACSA.2018.090405

[5] Qian, Jian-Feng, et al. "A Wide Stopband Filtering Patch Antenna and Its Application in MIMO System." IEEE Transactions on Antennas and Propagation 67.1 (2019): 654-658.

[6] Zhu, Jianfeng, et al. "Wideband low-profile highly isolated MIMO antenna with artificial magnetic conductor." IEEE Antennas and Wireless Propagation Letters 17.3 (2018): 458-462.

[7] Wen, Dingliang, et al. "Design of a MIMO Antenna With High Isolation for Smartwatch Applications Using the Theory of Characteristic Modes." IEEE Transactions on Antennas and Propagation 67.3 (2019): 1437-1447.

[8] Zhang, C., Gao, J., Cao, X., Xu, L., & Han, J. (2018). Low Scattering Microstrip Antenna Array Using Coding Artificial Magnetic Conductor Ground. IEEE Antennas and Wireless Propagation Letters, 17(5), 869-872.

[9] Han, Z. J., Song, W., Zhu, Y. Q., & Sheng, X. Q. (2018, December). RCS Reduction and Gain Enhancement for Patch Antenna by Using

Low Profile EBG. In 2018 12th International Symposium on Antennas, Propagation and EM Theory (ISAPE) (pp. 1-2). IEEE.

[10] Jam, S., & Simruni, M. (2018). Performance enhancement of a compact wideband patch antenna array using EBG structures. AEU-International Journal of Electronics and Communications, 89, 42-55.

[11] Islam, M. M., Islam, M. T., Samsuzzaman, M., & Faruque, M. R. I. (2015). Compact metamaterial antenna for UWB applications. Electronics Letters, 51(16), 1222-1224.

[12] Abdalla, M. A., & Ibrahim, A. A. (2013). Compact and closely spaced metamaterial MIMO antenna with high isolation for wireless applications. IEEE Antennas and Wireless Propagation Letters, 12, 1452-1455.

[13] Huang, C. Y., & Yu, E. Z. (2011). A slot-monopole antenna for dual-band WLAN applications. IEEE Antennas and Wireless Propagation Letters, 10, 500-502.

[14] Qin, P. Y., Weily, A. R., Guo, Y. J., & Liang, C. H. (2010). Polarization reconfigurable U-slot patch antenna. IEEE Transactions on Antennas and Propagation, 58(10), 3383-3388.

[15] Dang, L., Lei, Z. Y., Xie, Y. J., Ning, G. L., & Fan, J. (2010). A compact microstrip slot triple-band antenna for WLAN/WiMAX applications. IEEE Antennas and Wireless Propagation Letters, 9, 1178-1181.

[16] Hu, W., Yin, Y. Z., Fei, P., & Yang, X. (2011). Compact triband square-slot antenna with symmetrical L-strips for WLAN/WiMAX applications. IEEE Antennas and Wireless Propagation Letters, 10, 462-465.

[17] Cetiner, B. A., Crusats, G. R., Jofre, L., & Biyikli, N. (2010). RF MEMS integrated frequency reconfigurable annular slot antenna. IEEE Transactions on Antennas and Propagation, 58(3), 626-632.

[18] Balanis, Constantine A. Antenna theory: analysis and design. John wiley & sons, 2016.

# Speech Recognition System based on Discrete Wave Atoms Transform Partial Noisy Environment

Mohamed Walid[1], Bousselmi Souha[2], Cherif Adnen[3]

Laboratory Analysis and processing of electrical and energy systems
Faculty of Sciences of Tunis, FST, Tunis, Tunisia

*Abstract*—**Automatic speech recognition is one of the most active research areas as it offers a dynamic platform for human-machine interaction. The robustness of speech recognition systems is often degraded in real time applications, which are often accompanied by environmental noises. In this work, we have investigated the efficiency of combining wave atoms transform (WAT) with Mel-Frequency Cepstral Coefficients (MFCC) using Support Vector Machine (SVM) as classifier in different noisy conditions. A full experimental evaluation of the proposed model has been conducted using Arabic speech database (ARADIGIT) and corrupted with "NOISEUS database" noises at different levels of SNR ranging from -5 to 15dB. The results of Simulation have indicated that the proposed algorithm has improved the recognition rate (99.9%) at 15 dB of SNR. A comparative study was conducted by applying the proposed WAT-MFCC features to multilayer perceptron (MLP) and hidden Markov model (HMM) in order to prove the efficiency and the robustness of the proposed system.**

*Keywords*—*WAT; SVM; HMM; thresholding; noise; MFCC; MLP*

## I. INTRODUCTION

Automatic speech recognition allows the machine to understand and process information provided orally by a human user. It consists of using matching techniques to compare a sound wave to a set of samples, usually composed of words but also, more recently, phonemes (minimum sound unit). Speech recognition is based on the knowledge of several sciences: anatomy (The functions of the phonatory apparatus and the ear), the signals emitted by speech, phonetics, signal processing, linguistics, computer science, artificial intelligence and statistics. Automatic speech recognition opens new perspectives, given the considerable difference between manual and voice control. The use of natural language in the human-machine dialogue puts technology within the reach of all and leads to its popularization, reducing the constraints of the use of keyboards, mice and command codes to control. By simplifying the human-machine dialogue protocol, the automatic speech processing also aims to gain productivity since it is the machine that adapts to humans to communicate, not the other way around. In addition, it makes possible the simultaneous use of the eyes or hands to another task. It helps to humanize information management systems by focusing their design on users. A good speech recognition rates have been mostly reached using small vocabularies. In fact, this result is considered to be sufficient for the implementation of the most voice control devices. Thus, the error rate and learning time are steadily decreased. Also, this rate is

obviously variable and depends on vocabulary and language dialect. Sometimes, the ASR system has some troubles making it enables to avoid some linguistic traps. However, the development of sophisticated ASR systems with an improved speech recognition rates has become an interesting research subject for all scientific researchers in speech recognition domain. Indeed, several classification and parameterization methods have been emerged to accomplish this task. Among the most commonly used classification methods, we can cite HMM, SVM, Artificial neural networks (ANN), and Gaussian Mixture Model (GMM) [1]-[4]. Nevertheless, MFCC, Linear Predictive Coding (LPC), and Perceptual Linear Prediction (PLP) [5] constitute the well-known parameterization approaches. Among the most recent researches in the field of automatic speech recognition, attention may be attracted to [6] in which a deep Belief Networks (DBNs) has been presented at the aim to extract discriminative information using frames with a larger size in a speech signal. The desired objective in that work was to explore efficiency the DBNs in learning features which are more invariant to deep fluctuation in speech signal. Relying on other recognition methods, the author has tried to prove the reliability of his adopted technique which has significantly contributed to reduce the recognition error rate. However, in [7] an analysis of the impact of database size has been performed and that one of the impact of dialect in the context of independent-speaker text using SVM alone and then hybridized with GMM has been carried out too.

In this paper, an efficient isolated-words recognition system has been developed using SVM as classifier, and Mel-frequency Cepstral coefficient (MFCC) combined with discrete wave atoms transform (WAT) as feature extraction methods. Indeed, the adopted approach has been tested on Arabic language database in both clean and noisy conditions.

This manuscript is structured as follows: In Section II, a brief literature review of ASR Systems is presented. In Section III, the proposed speech recognition system is exhibited. However, an analysis and experimental results are given in Section IV. Concluding remarks and perspectives are presented in Section V.

## II. RELATED WORKS

Several works have been applied for isolated-words recognition systems using different classification algorithms, such as machine learning via ANN, GMM, Dynamic Time Warping (DTW), K-Nearest Neighbor (KNN), SVM, and HMM. A summarize of the performances of some ASR systems are given in Table I.

TABLE I.        LITERATURE SURVEY OF SPEECH RECOGNITION

| REF | parameterization | classifier | database | Accuracy (%) |
|-----|------------------|------------|----------|--------------|
| [8] | MFCC | ANN | Language Indian | 99.84 |
|     |      | SVM |                | 94.25 |
| [9] | MFCC | DTW KNN | Language Quechua | 91.1 |
| [10] | MFCC | GMM KNN | Language Hindi | 94.31 |
| [11] | MFCC | DTW KNN | Language English | 98.4 |
| [12] | CMUSphinx toolkit | HMM | Language Indonesian | 80 |
| [13] | MFCC | HMM | Language Maly | 80 |
| [14] | MFCC | HMM | Language Maithili | 95 |

## III. THE PROPOSED SPEECH RECOGNITION SYSTEM

A new speech recognition system based on WAT-MFCC and SVM was developed in this paper to improve the accuracy of recognition. The general scheme of the proposed speech recognizer is depicted in Fig. 1. It contains three parts: the preprocessing step, the features extraction step and SVM classification. Each part is explained in the following sections of this manuscript.

### A. Feature Extraction Stage

#### 1) Pre-processing speech stage

*a) Pre-emphasize:* This step consists of signal filtering and the determination of its first order finite impulse response such as.

$$H(z) = 1 - \alpha Z^{-1} \quad 0.9 \leq \alpha \leq 1 \tag{1}$$

The signal is sampled at 16 kHz and then pre-emphasized to determine high frequencies which are less energetic than low ones.

*b) Windowing:* In this step, the signal is split into frames by multiplying the frame samples by a Hamming window of length 20 ms. This latter is given by the following expression:

$$w(n) = 0.54 - 0.46 \cos\left(\frac{2\pi n}{N-1}\right) \quad 0 \leq n \leq N-1 \tag{2}$$

*2) Discrete wave atoms transform:* The wave atoms transform is an orthogonal transformation that was developed by Demanet. This transformation offers a parsimonious representation of a signal. In [15], a signal is assumed as an oscillatory model that can be described by the following function:

$$f(x) = \sin(Ng(x))h(x) \tag{3}$$

Where $x$ represents the coordinate, $g$ and $h$ are $C^{\infty}$ scale function. $h$ Has a compact support in the closed set $[0,1]^2$ and N is a constant. In order to get sparse solution of signal f, wave atoms were proposed by Demanet and Ying in [15] and Demanet in [16].



Fig. 1.   Block Diagram of the Proposed Speech Recognition System.

The wave transform is a variant of wavelet packets. However, it allows a high-frequency localization that cannot be obtained using a filter bank as the case of the wavelet packets transform. To situate the wave atoms transform among the existing geometrical transforms, according to the analysis given by Demanet, two parameters are involved. The parameter, α represents the multi-scale aspect of a transformation and the parameter β, represents its directionality aspect. A transformation is uniform for α= 0, dyadic for (α = 1).

Fig. 2 shows the distribution of various transforms according to their multi-scale and multidirectional aspects.

When α = 0, the transform is not multi-scale, such as the case of the Gabor transform. When α = 1 the transform is multi-scale, this is the case of wavelets, ridgelets and curvelets. When β = 0 the transform has a maximum selectivity as the Gabor transform where the decomposition can be done in all possible directions. Whereas β = 1 represents a minimal selectivity, this is the case of wavelets where there is very little directional information.

In this context, the transformation into wave atoms is applicable in the case where α = β = 1/2. This transformation is as multi-scale as multidirectional. Therefore, it represents the compromise between the two properties.

Wave atoms 1D family function is defined as $\varphi_\mu(x)$, with subscript $\mu = (j,\ m,\ n)$. The indexed point $(x_\mu, \omega_\mu)$ in phase-space is defined as follows.

$$x_\mu = 2^{-j}m, \omega_\mu = \pi 2^j n$$
$$C_1 2^j \leq \max |m_i| \leq C_2 2^j \tag{4}$$

The elements of frame $\varphi_U(x)$ are named wave atoms when:

$$|\hat{\varphi}(\omega)| \leq C_M 2^{-J}(1+2^{-j}|\omega-\omega_\mu|)^{-M} + C_M 2^{-J}(1+2^{-j}|\omega+\omega_\mu|)^{-M}$$
$$|\varphi_\mu(x)| \leq C_M 2^j (1+2^j |x-x_y|)$$
$$M > 0 \tag{5}$$



Fig. 2. Identification of Various Transforms as (α, β) Families of Wave Packets [16].

Practically, wave atoms are constructed from tensor products of a wavelet packet obeying the parabolic scaling law that is performed using decomposition as incomplete wavelet packets, as shown in Fig. 3. This figure describes the decomposition tree corresponding to the transformation into wave atoms defining the partitioning of the frequency axis in 1D. LH and RH respectively designate the left and right windows.

*3) Thresholding:* Thresholding is considered as a crucial step in speech recognition in tumultuous environments. It allows the rejection of the coefficients in which the WAT transformation is lower than a given threshold. There are several Thresholding methods, like the hard threshold and the soft threshold, which are the most used methods. In this article, we will apply the hard threshold given by the following equation:

$$C_{\text{Re}} = \begin{cases} C_{\text{Re}} & if\ |C_{\text{Re}}| \geq T \\ 0 & otherwise \end{cases} \tag{6}$$

*4) Mel-Scaled Frequency Cepstral Coefficients (MFCCs):* Owing to their low complexity estimation and their good performance, MFCCs are commonly integrated in automatic speech recognition systems [17]. Thus, the MFCC representation appears to better approximate the structure of human auditory system than the traditional linear and predictive features. Although, MFCC coefficients are easily affected by common perturbations of random localized-frequency, to which human perception is largely insensitive. For each speech frame, a MFCC vector is computed as follows:

- The spectrum power of the windowed signal block is mapped onto the Mel scale using triangular filters.

- The Discrete Cosine Transform (DCT) is then applied to filter bank output logarithm.

The relationship between scale Mel and frequency is given by the following expression

$$f_{mel} = 2595 \times \log_{10}\left(1+\frac{f}{700}\right) \tag{7}$$

An illustration of the computing of MFCCs coefficients is shown in Fig. 4.



Fig. 3. Strategy of Wave Atoms and Corresponding Set of Sub-Bands [16].

Fig. 4. Process to Create WAT- MFCC.

### B. Support Vector Machine (SVM) Classification

SVM is considered as a segregating classifier associated with a few real-world applications. By maximizing the margin between boundary points of the classes and the separating hyper plane, SVM represents the "best" separating hyper plan (Fig. 5). These end-focuses are known as support vectors. SVMs apply linear and nonlinear separation hyper plans to classify data.

By specifying how SVM finds an ideal hyper plan that ranks the new models. The feature vector is noted by $x_i \in R^n$ , $i \in \{1,\dots,M\}$ , where M is a number of training samples, n is the number of speech signal characteristics. The aim is to classify a speech unit into two classes $y_i = -1$ and $y_i = +1$, which design respectively the out vocabulary units and the hyper plane vocabulary unit.

The hyper plan form is described by the following equation:

$$\omega^t x + b = 0 \tag{8}$$

Where $\omega$ is a $M$ dimensional vector normal to the plane and $b$ is a scalar.



Fig. 5. Binary SVM Classification [18].

SVMs construct this hyper plan ($\omega$ and b) to be far from the boundary of both classes.

$$\omega^t x_i + b \geq 1 - \xi_i \ \text{if} \ y_i = +1 \tag{9}$$

$$\omega^t x_i + b \leq -1 + \xi_i \ \text{if} \ y_i = -1 \tag{10}$$

Where $\xi = (\xi_1,\dots,\xi_M)$ is slack variable, it controls the further processing of outliers, called "Sof-margin SVM".

In order to find best hyper plan we have to maximize the margin with constraints combining in constraints (9) and (10). In such case, the margin is equal to $\dfrac{1}{\|\omega\|}$ ; hence, this problem is equivalent to minimizing $\dfrac{1}{2}\|\omega\|^2 + c\sum_{i=1}^{M}\xi_i$ . This optimization problem can be solved by Lagrange Multipliers method [18].

$$L_p \equiv \frac{1}{2}\|\omega\|^2 + c\sum_{i=1}^{M}\xi_i - \sum_{i=1}^{M}\alpha_i y_i (x_i.\omega + b) + \sum_{i=1}^{M}\alpha_i \tag{11}$$

It is a convex optimization problem and can be solved by Quadratic Programming method which returns the optimal $\alpha$ and it permits us to find $\omega$ by applying the formula given by (12).

In the next step, we have to describe the set of Support Vectors S which indices hold $\alpha \rangle 0$ condition and calculate b:

$$\frac{\partial Lp}{\partial \omega} = 0 \Rightarrow \omega = \sum_{i=1}^{M}\alpha_i y_i x_i \tag{12}$$

$$\frac{\partial Lp}{\partial b} = 0 \Rightarrow \sum_{i=1}^{M}\alpha_i y_i = 0 \tag{13}$$

Substituting (12) and (13) into (11) gives Dual form $L_D$ of the Primary $L_P$:

$$L_D \equiv \sum_{i=1}^{M}\alpha_i - \frac{1}{2}\sum_{i,j}^{M}\alpha_i H_{ij}\alpha_j \tag{14}$$

Where

$$H_{ij} = y_i y_j x_i x_j$$

Minimizing $L_P$ equivalent to maximizing $L_D$:

$$\sum_{i=1}^{M}\alpha_i y_i = 0, \alpha_i \geq 0, i = 1,\dots\dots\dots,M$$

It is a convex optimization problem and can be solved by Quadratic Programming method which returns "best" $\alpha$ and it is allow us to calculate w by the formula (12).

In the next step we need to define the set of Support Vectors S which indices hold $\alpha_i > 0$ condition and calculate b:

$$b = \frac{1}{N_s} \sum_{s \in S} \left( y_s - \sum_{m \in S} \alpha_m y_m x_m . x_s \right) \tag{15}$$

Each new sample $x^*$ is classified by $y^* = sgn(\omega.x^* + b)$ in the test process.

SVMs transform data (feature vectors) into a high dimensional space, whose training data become linearly separable, i.e. SVMs transform n-dimensional feature vector x into a N-dimensional feature vectors [19]:

$$\Phi : R^n \rightarrow R^N$$

This transformation is realized via Kernel functions. There are different types of Kernel functions in SVM method such as:

Linear Kernel: $K(x_i, x_j) = x_i^T x_j \tag{16}$

Sigmoid Kernel: $K(x_i, x_j) = \tanh(ax_i \cdot x_j - b) \tag{17}$

Radial Basis Kernel: $K(x_i, x_j) = e^{\left( -\frac{\|x_i - x_j\|^2}{2\sigma^2} \right)} \tag{18}$

Polynomial Kernel: $K(x_i, x_j) = (x_i . x_j + a)^b \tag{19}$

Where a and b are kernel's parameters.

Each new sample $x^*$ is classified by the following formula for the non-linear Kernel based SVM

$$y^* = sgn \left( \sum_{s \in S} y_s a_s K(x^*, x_s) + b \right) \tag{20}$$

SVM classification is a binary classification machine-learning algorithm. For this reason, it is required to modify decision-making part of the algorithm for the multiclass tasks in real world situations. Two most popular methods are used for multiclass classification, which are "one against one" (OAO) and "one against all" (OAA) techniques. The OAA SVM multiclass approach involves the division of N class dataset into N two-class cases, whereas OAO approach involves training separate classifier for each pair of classes. This leads to $\frac{N(N-1)}{2}$ classifiers. This method is less sensitive to the imbalanced datasets but it is more computationally expensive. In this paper, we use "one-against-all" SVM multiclass technique.

## IV. ANALYSIS AND RESULTS

### A. Speech Database

The database used in this work is the ARADIGIT database [20], which is a local Arabic speech database, collected from 110 Algerian speakers who live in different regions of their country and speaks different dialects. The individuals chosen to record the words used in the database are aged between 18 and 50 years old. The database was recorded in a very quiet environment, at sample frequency 22.050 kHz and down-sampled at 16 kHz. In the training phase, they have used 1800 utterances pronounced by 60 speakers of both genders (30 male speakers and 30 female speakers) where each speaker repeats the same digit 3 times. In the testing phase, they have used 1000 utterances pronounced by 50 other speakers of both gender (25 males and 25 females) where each speaker repeats the same digit 2 times.

### B. Experimental Results

Performances of the proposed system were tested on both clean and noisy signals. For noisy signals, we have corrupted data with three noises (white, car and babble) extracted from the NOISEX92 database (Varga et al. 1992). Noises are added to the speech signal in a SNR range from -5 to 15 dB with step size 5 dB. For training phase, clean speech database are explored. Indeed, noises are only added for testing the recognition performance. The robustness of the system is evaluated according to the recognition rate defined below.

$$Accur = \left[ \frac{Correctly\ recognised\ samples}{Total\ number\ of\ test\ samples} * 100 \right] \tag{21}$$

Using three classifiers (SVM, HMM and MLP) in clean condition, the obtained recognition rates with MFCC only and WAT-MFCC parameterizations are presented in Fig. 6. Obviously, we can note that SVM combining MFCC soft and WAT-MFCC classifier has reached best performances in terms of recognition accuracy in comparison with other classifier as HMM and MLP. Thus, SVM based MFCC only reached 94.6% whereas using WAT-MFCC; we obtained a rate of 100%. In addition, it is remarkable that MFCC soft has contributed to reach competitive performances compared to WAT-MFCC. Despite worst performances have been obtained using MLP based MFCC soft with an achieved rate of 84.2%; the use of WAT-MFCC has registered an acceptable accuracy 92.4%.

Recognition results obtained from comparison between MFCC and WAT-MFCC using three speech recognition algorithms are illustrated in the following tables. Based on the displayed numbers in Tables II to IV, it has to note that WAT-MFCC has the ability to achieve the best recognition accuracy with respect to MFCC for all SNR levels and in all noisy conditions using SVM classifier. These results were achieved in comparison with HMM and MLP algorithms (see Tables II to zizv). Indeed, the best recognition rate (99.9%) was recorded at 15 dB of SNR. As shown, HMM

and MLP algorithms have contributed to give the second and third best recognition accuracies, respectively, at different levels of SNR under different noisy conditions. Indeed, these algorithms have recorded respectively 84.8% and 81.3% rates in "white" noise condition at 15 dB of SNR. Moreover, it can be seen that more the level of SNR increases, more the recognition accuracy increases too. In all tests, worse results were recorded with "babbel" noise compared to those obtained with "white" and "car" noises.



Fig. 6. Recognition Accuracy and Training Time using Three Classifiers with MFCC only and WAT-MFCC in Clean Condition.

TABLE II. COMPARISON OF THE OBTAINED RECOGNITION ACCURACY FOR MFCC ONLY AND WAT-MFCC USING DIFFERENT CLASSIFIERS IN CAR NOISE CONDITION

| Noise | Speech recognition algorithm | | SNR | | | | |
|---|---|---|---|---|---|---|---|
| | | | -5db | 0db | 5db | 10db | 15db |
| Car | SVM | WAT-MFCC | **67.4** | **72.7** | **84.8** | **95.4** | **96.7** |
| | | MFCC | 60.60 | 65.4 | 70.4 | 75.7 | 90.90 |
| | HMM | WAT-MFCC | 63.6 | 65.2 | 71.2 | 75.1 | 83.8 |
| | | MFCC | 57.5 | 60.60 | 62.1 | 74.2 | 80.30 |
| | MLP | WAT-MFCC | 60.60 | 62.3 | 66.6 | 73.4 | 81.2 |
| | | MFCC | 54.5 | 54.5 | 59.6 | 64.4 | 71.9 |

TABLE III. COMPARISON OF THE OBTAINED RECOGNITION ACCURACY FOR MFCC ONLY AND WAT-MFCC USING DIFFERENT CLASSIFIERS IN WHITE NOISE CONDITION

| Noise | Speech recognition algorithm | | SNR | | | | |
|---|---|---|---|---|---|---|---|
| | | | -5db | 0db | 5db | 10db | 15db |
| White | SVM | WAT-MFCC | **68.7** | **86.3** | **90.90** | **92.4** | **99,9** |
| | | MFCC | 64.3 | 68.9 | 75.3 | 84.8 | 90,90 |
| | HMM | WAT-MFCC | 60.60 | 65.9 | 68.9 | 73.4 | 84.8 |
| | | MFCC | 56.8 | 59.1 | 63.6 | 67.4 | 73.4 |
| | MLP | WAT-MFCC | 57.5 | 60.60 | 65.1 | 69.6 | 81.3 |
| | | MFCC | 53.6 | 58.3 | 62.8 | 65.90 | 76.9 |

TABLE IV. COMPARISON OF THE OBTAINED RECOGNITION ACCURACY FOR MFCC ONLY AND WAT-MFCC USING DIFFERENT CLASSIFIERS IN BABBLE NOISE CONDITION

| Noise | Speech recognition algorithm | | SNR | | | | |
|---|---|---|---|---|---|---|---|
| | | | -5db | 0db | 5db | 10db | 15db |
| Babble | SVM | WAT-MFCC | **66.6** | **71.8** | **84.6** | **91.90** | **96.9** |
| | | MFCC | 62.8 | 64.4 | 66.8 | 77.2 | 80.30 |
| | HMM | WAT-MFCC | 63.6 | 65.90 | 73.4 | 75.7 | 81.5 |
| | | MFCC | 59.6 | 61.1 | 65.90 | 70.6 | 73.4 |
| | MLP | WAT-MFCC | 60.60 | 63.6 | 71.2 | 73.4 | 77.9 |
| | | MFCC | 56.8 | 60.60 | 65.8 | 71.2 | 71.6 |

From the obtained results, we can observe that SVM classifier with WAT-MFCC extraction features has proved its efficiency compared to HMM and MLP classifiers. Also, the adopted speech recognition method seems to be promising as it can greatly contribute to reach the best performances in noisy conditions.

## V. CONCLUSION

In this paper, a new model for Arabic speech recognition using a combination of WAT with MFCC feature has been presented. The obtained results of the proposed method have shown that WAT is a powerful and effective technique in optimizing MFCC parameters. In fact, it has enhanced the learning ability of SVM which has given us a high recognition rate (100%) in clean environment. Furthermore, the assessment of the proposed method has been performed in noisy conditions without any speech enhancement algorithm. Also, it has been compared to other speech recognition approaches such as HMM and MLP. Moreover, the evaluation has shown that the proposed technique have reached the best performance (99.9%) in term of recognition rate in noisy condition at 15dB of SNR.

As a further work, we would like to test the proposed model on electronic architectures such as raspberry pi3, FPGA, and STM32 in order to more follow its evolution in real time applications.

REFERENCES

[1] J. Ming, D. Crookes, Speech Enhancement Based on Full-Sentence Correlation and Clean Speech Recognition, IEEE/ACM Transactions on Audio, Speech, and Language Processing.25 (2017) 531–543.

[2] Z-Q. Wang, D. Wang, A Joint Training Framework for Robust Automatic Speech Recognition, IEEE/ACM Transactions on Audio, Speech, and Language Processing. 24 (2016) 796–806.

[3] S. Chandrakala, N. Rajeswari, Representation Learning Based Speech Assistive System for Persons With Dysarthria, IEEE Transactions on Neural Systems and Rehabilitation Engineering. 25 (2017) 1510–1517.

[4] Z.Tang, L. Li, D. Wang, R.Vipperla, Collaborative Joint Training With Multitask Recurrent Model for Speech and Speaker Recognition. IEEE/ACM Transactions on Audio, Speech, and Language Processing. 25 (2017) 493–504.

[5] D. Namrata. Feature extraction methods LPC, PLP and MFCC in speech recognition. International journal for advance research in engineering and technology, 2013, vol. 1, no 6, p. 1-4.

[6]   M. Farahat,"Noise Robust Speech Recognition Using Deep Belief Networks", International Journal of Computational Intelligence and Applications, Vol. 15, No. 1 (2016) 1650005 (17 pages).

[7]   K. Yasmine, Z. and A. Amrouche, "SVM against GMM/SVM for Dialect Influence on Automatic Speaker Recognition Task", International Journal of Computational Intelligence and Applications, Vol. 13, No. 2 (2014) 1450012 (10 pages).

[8]   Londhe, N. D., & Kshirsagar, G. B. (2017). Speaker independent isolated words recognition system for Chhattisgarhi dialect. 2017 International Conference on Innovations in Information, Embedded and Communication Systems (ICIIECS).

[9]   Hernan Faustino Chacca Chuctaya, Rolfy Nixon Montufar Mercado and Jeyson Jesus Gonzales Gaona, "Isolated Automatic Speech Recognition of Quechua Numbers using MFCC, DTW and KNN" International Journal of Advanced Computer Science and Applications(ijacsa), 9(10), 2018.

[10]  P. Wani, U. G. Patil, D. S. Bormane, and S. D. Shirbahadurkar, "Automatic speech recognition of isolated words in Hindi language," in Proceedings - 2nd International Conference on Computing, Communication, Control and Automation, ICCUBEA 2016, 2017.

[11]  I.,Muhammad, Atif and R. Gulistan. Isolated Word Automatic Speech Recognition (ASR) System using MFCC, DTW & KNN. In : 2016 Asia Pacific Conference on Multimedia and Broadcasting (APMediaCast). IEEE, 2016. p. 106-110.

[12]  P. Hamdan, F. Ridi, and H.   Rudy. Indonesian Automatic Speech Recognition system using CMUSphinx toolkit and limited dataset. In : 2016 International Symposium on Electronics and Smart Devices (ISESD). IEEE, 2016. p. 283-286.

[13]  ANAND, Anu V., DEVI, P. Shobana, STEPHEN, Jose, et al.Malayalam Speech Recognition system and its application for visually impaired people. In : 2012 Annual IEEE India Conference (INDICON). IEEE, 2012. p. 619-624.

[14]  R. Rajeev, and D. R.  Kumar. Isolated word recognition using HMM for Maithili dialect. In : 2016 International Conference on Signal Processing and Communication (ICSC). IEEE, 2016. p. 323-327.

[15]  L. Demanet, L. Ying, Wave atoms and sparsity of oscillatory patterns, Applied and Computational Harmonic Analysis. 23 (2007) 368–387.

[16]  L. Demanet, "Curvelets, wave atoms, and wave equations," PHD Thesis in California Institute of Technology, 2006.

[17]  S. Mada, and S. Zabidin. Implementasi Pengenalan Pola Suara Menggunakan," Mel-Frequency Cepstrum Coefficients (MFCC) dan Adaptive Neuro-Fuzzy Inferense System (ANFIS), '' sebagai Kontrol Lampu Otomatis. ALHAZEN, 2014, vol. 1, no 1, p. 43-54.

[18]  A. Ganapathiraju, J. Hamaker and J. Picone, Applications of Support Vector Machines to Speech Recognition. IEEE Transactions on Signal Processing. 52 (2004) 2348–2355.

[19]  M. Elleuch, R. Mokni, M. Kherallah, Offline Arabic Handwritten Recognition System with Dropout applied in Deep Networks based-SVMs.International Joint Conference on Neural Networks,Vancouver, BC, Canada 2016,pp.3241-3248.

[20]  A. Amrouche "Reconnaissance automatique de la parole par les modèles connexionnistes" .Thèse de doctorat, facultéd' électroniqueetd' informatique, USTHB. 2007.

# Intelligent Scheduling of Bag-of-Tasks Applications in the Cloud

Preethi Sheba Hepsiba[1], Grace Mary Kanaga E[2]

Department of Computer Science and Engineering, Karunya Institute of Technology and Sciences, Coimbatore, India[1, 2]
Department of Computer Science and Engineering, CMR Institute of Technology, Bangalore, India[1]

*Abstract*—**The need of efficient provision resources in cloud computing is imperative in meeting the performance requirements. The design of any resource allocation algorithm is dependent on the type of workload. BoT (Bag-of-Tasks) which is made up of batches of independent tasks are predominant in large scale distributed systems such as the cloud and efficiently scheduling BoTs in heterogeneous resources is a known NP-Complete problem. In this work, the intelligent agent uses reinforcement learning to learn the best scheduling heuristic to use in a state. The primary objective of BISA (BoT Intelligent Scheduling Agent) is to minimize makespan. BISA is deployed as an agent in a cloud testbed and synthetic workload and different configurations of a private cloud are used to test the effectiveness of BISA. The normalized makespan is compared against 15 batch mode and immediate mode scheduling heuristics. At its best, BISA produces a 72% lower average normalized makespan than the traditional heuristics and in most cases comparable to the best traditional scheduling heuristic.**

*Keywords—Bag-of-tasks applications; intelligent agent; reinforcement learning; scheduling*

## I. INTRODUCTION

A formidable challenge in cloud computing is the effective allocation of resources. In traditional scheduling, the efficacy of any scheduling heuristic largely depends on the characteristics of the workload (tasks to be scheduled) and resources. For a cloud service provider, maximum resource utilization and minimum power consumption result in the most profitability. For the customer, the scheduling of tasks in the cloud is also affected by network latency and cost. In addition, resource providers also strive to meet QoS (Quality of Service) requirements [1] from the customer in addition to minimizing energy consumption and maximizing resource utilization. Several VMs are hosted in a multi-tenant architecture, migrated to different hosts with different computing capacities in a datacenter. The principles that aid in attaining these objectives often result in a dynamically changing environment on the available cloud resources. This inadvertently leads to performance unpredictability. Andreadis et al [2] when proposing a reference architecture for datacenter scheduling pointed out that because of the complex nature of the environment in a datacenter, comparing scheduling heuristics and improving performance is challenging.

Self-learning systems [3] are the future of cloud computing. Reinforcement learning has been used for energy-aware resource scheduling [4] and results show that it can enhance energy efficiency in data centers. Adaptive scheduling for resource provisioning [5] also uses reinforcement learning.

This work differs from adaptive scheduling in that the heuristics or meta-heuristics are not modified. However, by incorporating reinforcement learning, the state is observed and based on the information, the agent is trained to choose the best heuristic from its repository.

The bag-of tasks applications are a popular type of workload in the cloud and challenging to schedule in heterogeneous machines as opposed to homogeneous machines[6] for which an optimal schedule can be produced. Once tasks are allotted to a private cloud or any set of resources on a public cloud, there is a need for an intelligent agent in scheduling that is adaptable to a dynamic environment prevalent in the cloud. The hypothesis is that if an agent can sense the environment, apply learned best scheduling heuristic and simultaneously explore other options, over time the agent will choose the best heuristic for a state. In this work, an intelligent agent, BISA (BoT Intelligent Scheduling Agent) is proposed that uses reinforcement learning to choose the best scheduling heuristic in a state. BISA recognizes the current state based on the characteristics of the BoT workload and the available resources. It uses reinforcement learning to choose the best-known scheduling heuristic for the given state. The learning parameters in BISA, namely α and β control the exploration vs. exploitation strategy of BISA.

In previous work [7], the BIS agent was presented and the preliminary results from the training phase of BISA were divulged and discussed. In this work, BoT workload is generated synthetically. The testing is carried out rigorously for different configurations, and comparison of the normalized makespan over a series of runs of each cycle is presented for various learning parameters.

The contributions for scheduling bag-of-tasks using agents are:

- To develop BISA with an objective to minimize makespan.

- To test and present the results of the BISA agent on synthetic workload for different sets of learning parameters

- To test the hypothesis that BISA works in a dynamic environment. (simulated by changing configurations of Hosts and VMs)

The rest of the paper is organized as follows. In Section II, an overview of the BoT Workload and several heuristics in literature and the agent-based paradigm is presented. The

framework of the intelligent agent is presented in Section III. The results and discussion are elaborated in Section IV. The conclusion and future work are described in Section V.

## II. RELATED WORK AND MOTIVATION

In cloud computing which from an architectural point of view is a large-scale distributed system, traditional scheduling algorithms are not enough to schedule tasks efficiently. The type of workload also affects the design of a scheduling algorithm.

The bag-of-tasks workload is examined and the motivation for generating a synthetic workload is presented. Traditional heuristics and meta-heuristics that have been used to solve the resource allocation problem is presented followed by the motivation for the agent-based approach.

### A. BoT Workload

A type of workload prevalent in large scale distributed systems such as the grid or cloud is the bag-of-tasks (BoT) workload. BoT workload comprises of independent tasks typically submitted as part of a larger application by the same user. The scheduling of independent tasks (bag-of-tasks) in heterogeneous systems is known to be an NP-Complete problem[8]. The following rules were used to identify BoTs from a trace [9]:

- BoT Size ranges from 2 to 64 parallel tasks.

- There are no serial jobs

- The run times are not extremely small and range from a few minutes to a day.

A synthetic workload [10] can also be generated based on the characteristics of a BoT. The model used was developed to test the performance of BoT workload in large scale distributed systems[11]. A synthetic workload offers more control in varying parameters. In this work, a synthetic BoT workload is generated, and the parameters used for the distribution are elaborated in Section III.

### B. Traditional Heuristics

Schedulers typically use either an immediate mode or batch mode heuristic[8], [12] to schedule tasks in the cloud. The immediate mode heuristics described are (i) MCT (Minimum Completion Time), (ii) MET (Minimum Execution Time), (iii) Switching Algorithm, (iv) K-percent best, (v) Opportunistic Load Balancing. The Batch Mode Heuristics described are Min-Min, Max-Min Heuristic, and Sufferage. Some of these heuristics are based on concepts described in prior research work [8]. Enhanced versions of the Fastest Processor to Largest Task First (FPLTF) and sufferage are also described in [13].

An extensive list of 14 heuristics derived from 3 types of ordering of tasks within the BoT (Uniform, Large to Small and Small to Large) and the mapping of the task (Random Mapping, Maximum Expected Remaining Allocation Time Mapping, Maximum Current Remaining Allocation Time Mapping, Minimum Expected Remaining Allocation Time Mapping and Minimum Current Remaining Allocation Time Mapping) is presented by Garcia & Sim [14].

In many schedulers used in datacenters, traditional scheduling heuristics are used, but the pitfall is that none of the traditional heuristics are optimal in a heterogeneous system such as the cloud.

### C. Meta-Heuristics

Several metaheuristics have also been applied to resource scheduling in the cloud to reduce the makespan and increase the throughput. These heuristics aim to produce optimal schedules for a multi-objective scheduling problem.

Tabu Search (TS) and Simulated Annealing (SA) were compared against Fastest Processor Largest Task (FPLT) [9] and the results show that Tabu Search and Simulated Annealing performed consistently better than FPLT even as BoT sizes increased. SA and TS had an 8% to 9% smaller makespan as compared to FPLT for globally arriving BoTs. The only shortcoming is the performance overhead incurred by TS and SA. A variation of SA, Thermodynamic Simulated Annealing (TSA) [9] is presented that also performs considerably better than SA.

A parallel Genetic algorithm [15] has been used in cloud scheduling which also tries to improve the VM utilization rate instead of just concentrating on a good resource scheduling algorithm. An IGA (Improved Genetic Algorithm) [16] to improve the utilization rate of VM's has also been proposed and the results show that the performance of IGA is twice that of TGA (Traditional Genetic Algorithm). A hybrid algorithm that used Genetic Algorithm, round-robin scheduling in deep neural learning works effectively in minimizing the makespan and other parameters in a cloud workflow. A software framework was developed for rapid prototyping of hybrid meta-heuristic schemes [17] as hybrid meta-heuristic techniques have proved to be more effective and adaptable.

A common disadvantage of most metaheuristics is the complexity but for a large-scale scheduling problem, the performance overhead is compensated with improvement in performance.

### D. Agent-Based Approach

In a large-scale system like the cloud, accurate system requirement is difficult to obtain. An estimate of the execution time of the task will not be available. Meta-heuristics may sometimes provide near-optimal schedules but the time complexity for producing the schedule may be extremely high when the magnitude of the system increases which is the case in cloud.

A concise set of challenges that result from these heuristics are that:

- Heuristics that solve this problem require accurate information about the environment.

- The execution time of each task should be known a priori.

- The complexity (although polynomial) may be unacceptably high [13] when there are many tasks or resources.

A vital observation based on applying 14 scheduling heuristics based on ordering of tasks within a BoT and the choice of mapping [14] is that due to the *NP-complete* nature of the scheduling problem, there was not a dominant scheduling heuristic from among the proposed heuristics or the benchmark scheduling heuristics . In a configuration where there are some powerful virtual machines (VMs) that are free, a batch mode heuristic that sorts a set of tasks, largest to smallest (LtoS) and assigns the largest tasks to a machine that is more powerful first would reduce the makespan. In a different configuration where all the powerful machines are already overloaded, assigning tasks sorted smallest to largest (StoL) and assigning the smallest tasks to the overloaded powerful machines and the larger tasks to a less powerful machine would result in a reduction in makespan.

### E. The Model of BISA

The BISA (BoT Intelligent Scheduler Agent) is an intelligent agent [18] that senses the environment and based on the current state chooses a scheduling policy, calculates the utility of that scheduling policy using a reward function and chooses to execute the schedule produced or choose another scheduling policy that gives a better utility. The result of the BIS agent is a schedule for the set of tasks in a BoT on a set of resources with a goal to minimize the makespan. Fig. 1 shows the abstract working of BISA and the context it works in a hybrid cloud environment. BISA works within a cloud (private or public) to effectively schedule the BoTs.

The following variables are be used to describe the working of the BIS agent.

Let $S$ be the set of states that the agent has come across where

$S = \{s_1, s_2 \dots s_l\}$

Let $BoT_t$ be a *BoT* submitted at time $t$.

Let U be the set of users who submit BoTs where U $= \{u_1, u_2, \dots u_n\}$

Let $N_t$ denote the number of tasks in a $BoT_t$ submitted at time $t$.

Let $I_t$ be the ideal makespan a *BoT* submitted at time $t$.

Let $M(BoT_t)$ be the makespan of BoT submitted at time $t$ after it has finished execution.

Let $ST_t$ be the submission time of $BoT_t$.

Let $st_{i,t}$ be the execution start time of $task_i$ in $BoT_t$ where $i=1, \dots N_t$.

Let $wt_{i,t}$ be the waiting time of $task_i$ in $BoT_t$ where $i=1, \dots N_t$.

Let $ft_{i,t}$ be the finish time of $task_i$ in $BoT_t$ where $i=1, \dots N_t$.

Let $size\_mi_{i,t}$ be the size of task $i$ in $BoT_t$ where $i=1, \dots N_t$ in Million Instructions(MI).

Let $H$ be the set of hosts in a cloud where $H = \{h_1, h_2 \dots h_n\}$.

Let $VM$ be the set of virtual machine's in a cloud where $VM = \{vm_{1,1}, vm_{1,2}, \dots vm_{r,p}\}$ where $vm_{r,p}$ denotes a *VM p* assigned to host $r$.

Let $HN_r$ be the number of VMs assigned to host $r$.

Let $mips_{r,p}$ denotes the MIPS rating of a *VM p* assigned to host $r$.

Let $A$ be the set of actions that is available to the agent. $A = \{a_1, a_2 \dots a_m\}$. Table I shows the ordering and mapping policy that are used to produce the 15 possible heuristics as in [14]. Each action is produced by combining a task ordering with the task mapping policy. CT and ET are interpreted based on definitions in [12].

Let $Q_{private}$ and $Q_{public}$ denote the set of tasks that are assigned to the private or public Queue.

Let $r_{l,j,k}$ be the reward for the $k^{th}$ time an action $j$ applied on state $l$.

Let $Q_{l,j,k}$ be the average of the first $k$ rewards on action $j$ in state $l$.

The working of the BIS intelligent agent is based on reinforcement learning [19] and the various stages are outlined below:



Fig. 1. The BISA Agent in Hybrid Cloud.

TABLE I. ORDERING AND MAPPING POLICY OF THE HEURISTICS

| Ordering | |
|---|---|
| **Unordered (U)** | Random ordering of tasks within each BoT. |
| **LtoS** | Tasks are ordered Large-to-Small within the BoT. |
| **StoL** | Tasks are ordered Small-to-Large within the BoT. |
| **Mapping Policy** | |
| *CT: Expected completion Time on a machine after completing all the tasks the machine that was previously allotted.* | |
| *ET: Expected completion time on a machine when no load is assigned to it.* | |
| **MinCT** | Allocate task to the machine with minimum CT |
| **MaxCT** | Allocate task to a machine with maximum CT |
| **MinET** | Allocate task to a machine with minimum ET |
| **MaxET** | Allocate task to a machine with maximum ET |
| **Random (R)** | Allocate task to a random machine |

*1) Recognize the state:* The size ratio of a BoT and the VM utilization for each classification of VMs form a state.

Let $geo\_mean(BoT_i) = \left(\prod_{n=1}^{N_t} size\_mi_{n,t}\right)^{1/N_t}$ be the geometric mean of $BOT_t$.

Let #small_tasks $(BoT_t)$ be the count number of tasks in $BoT_t$ whose $size\_mi_{i,t}$ is lesser than $geo\_mean(BoT_i)$

Let #large_tasks $(BoT_t)$ be the count number of tasks in $BoT_t$ whose $size\_mi_{i,t}$ is greater than $geo\_mean(BoT_i)$

Let $size\_ratio(BoT_t)$ be the size ratio of BoT submitted at time $t$ calculated in Equation (1)

$$size\_ratio(BoT_t) = \frac{\#small\_tasks\ (BoT_t)}{\#large\_tasks\ (BoT_t)} \qquad (1)$$

The size ratio is calculated as the number of small tasks divided by the number of large tasks. It determines whether the BoT is left normally distributed (a greater number of small tasks) or right normally distributed (a greater number of large tasks) or an equal number of large and small tasks. The size ratio is essential when the agent runs in a real-time environment where the parameters for the normal distribution that are typically used in a simulation are not known a priori.

In order to calculate VM utilization, VMs are first classified into three categories based on MIPS following which the average percentage of computing utilization in each of the three categories are obtained.

Let $C=$ *{LG, ST, SM}* denote the classifications of VMs. Large VMs (*LG*) having processing capacity of MIPS > 80000, Standard VMs (*ST*) having MIPS between 40,000 and 80,000 and Small VMs (*SM*) having MIPS below 40,000.

Let $VM\_util_c$, = *{H, M, L}* denote the average computing utilization of all the cloudlets/tasks running on the *VM* where $c$ = *{LG, ST, SM}*. *H* denotes high CPU utilization (>80%). *M* denotes medium CPU utilization (40% to 80%). *L* denotes low CPU utilization (<40%).

*2) Decide on the action that is matching with the goal:* In this stage, the agent will choose the appropriate scheduling heuristic that minimizes the makespan of the BoT. The policy is chosen using a roulette wheel selector where all the probabilities for a set of actions in each state is on the scale of one. The values chosen for α and β allow us to control the tradeoff between exploration vs. exploitation. Exploration allows the agent to choose a new scheduling heuristic whereas exploitation acts like a greedy heuristic that choose the scheduling policy known to perform well.

*3) Execute action:* Schedule the BoT according to the scheduling policy decided upon in the previous step.

*4) Assess chosen action:* After the state is recognized, the ideal makespan is calculated for the given BoT. The ideal makespan, $I_t$ is taken as the time taken for the largest task in a BoT to execute on the fastest available machine. It gives the agent a reference value to calculate the reward.

The reward is calculated by how well a policy performs with respect to the ideal makespan. The ideal makespan is calculated using Equation (2) which is the time is taken for the largest task i.e. the task with the longest instruction length (MI) to execute in the fastest machine (the VM with the largest MIPS rating).

$$I_t = \frac{max \sum_{i=1}^{N_t}(size\_MI_{i,t})}{max \sum_{r=1}^{H} \sum_{p=1}^{HN_r} mips_{r,p}} \qquad (2)$$

The makespan of a $BoT_t$ on action $a_j$ is calculated using Equation (3) which is the difference between the maximum finish time of a task in $BoT_t$ and the submission time $BoT_t$.

$$M(BoT_t) = max \sum_{i=0}^{N_t}(ft_{i,t}) - ST_t \qquad (3)$$

This value of the makespan divided by the ideal makespan will be the reward, $r_{l,j,k}$ at the $k^{th}$ time the policy $a_j$ was used on a given state $l$ using Equation (4). Initially, the reward, $r_{l,j,0}$ is 0.

$$r_{l,j,k} = \ln(1/(M(BoT_t, a_j) - I_t) \qquad (4)$$

An incremental implementation is used to calculate the cumulative reward at step $k+1$ in Equation (5) and this value is updated in the table of rewards.

$$Q_{l,j,k+1} = Q_{k+1,j} + \frac{1}{r_{l,j,k}}\left[r_{l,j,k+1} - Q_{l,j,k}\right] \qquad (5)$$

In order to decide which action to choose at $k^{th}$ time to schedule $BoT_t$ in state $l$, the probability of choosing each action or policy is calculated. This probability is calculated by looking up a table of preferences for each action based on the reference reward.

Let $\pi_{l,\tau}(a_j)$ be the probability of choosing action $a_j$ at play $\tau$ for state $l$, i.e. the $t^{th}$ time the state $l$ is encountered.

Let $P_{l,\tau}(a_j)$ be the preference of an action selected at play $\tau$ of state $l$. The preference for an action determines the likelihood of an action being selected for a state. The initial action preferences to 0. The initial reference reward, $\bar{r}_{l,0}$ for every state $l$ where play $\tau$ is 0 is set to 0.1 meaning if the makespan of the scheduling policy is more than 10% as efficient as the ideal makespan, the reward is incremented and the preference for that action is set to a positive value. The reference reward following the first run, $\bar{r}_{l,\tau+1}$ is calculated as in Equation (6) where $0 < \alpha \le 1$.

$$\bar{r}_{l,\tau+1} = \bar{r}_{l,\tau} + \alpha\left[r_{l,\tau} - \bar{r}_{l,\tau}\right] \qquad (6)$$

During the first run, all the scheduling policies will have an equal probability of being chosen. Once a scheduling policy is chosen in play $\tau$ in state $l$, the preference for that action is updated. The preference is the difference between the reward $r_{l,\tau}$ and the reference reward $\bar{r}_t$ as given in Equation (7). The initial action preferences are set to 0.

$$P_{l,\tau+1}(a_j) = P_{l,\tau}(a_j) + \beta\left[r_{l,\tau} - \bar{r}_{l,\tau}\right] \qquad (7)$$

The constant $\beta$ is a positive step size parameter. A high reward will increase the probability of reselecting the action and a low reward should decrease the probability.

The probability $\pi_{l,\tau}(a_j)$ of selecting an action $j$ at play $\tau$ in state $l$ is calculated using Equation (8).

$$\pi_{l,\tau}(a_j) = \frac{e^{P_{l,\tau}(a_j)}}{\sum_{b=1}^{A} e^{P_{l,\tau}(a_b)}} \qquad (8)$$

The performance metrics used to assess BISA are presented below:

- Overall makespan:

The overall makespan of each run is calculated by first calculating the makespan of individual BoTs and then taking the sum of the makespan. The individual makespan is calculated as shown in Eq (2) and the overall makespan is computed according to Equation (9).

$$Overall\ Makespan = \sum_{t=0}^{BoT} M(BoT_t) \qquad (9)$$

- Normalized makespan:

The overall makespan can be a misleading metric because during each run of the simulation the load varies. The normalized makespan is based on the overall Makespan for all the BoTs generated within each run (48 intervals) calculated from Equation (9), the total runtime that denotes the actual runtime in the allotted VM of each task in a BoT during simulation and the number of tasks are those generated during each run (48 intervals). The normalized makespan is calculated by taking the sum of the tasks in all the BoTs generated during each run as shown in Equation (10).

$$Normalized\ Makespan = \frac{Overall\ Makespan}{Total\ Runtime * Number\ of\ Tasks} \qquad (10)$$

## III. EXPERIMENTAL SETUP

The framework provided by cloudsim [20] to simulate the cloud and incorporate the intelligent agent as a thread. The VM's in the cloud run on the assigned host in the datacenter the allocation policy is time-shared. Each VM runs on a Processing Element (PE) which is 1 core of the host machine.

To test the reinforcement learning agent, 7 Virtual Machines (VMs) are set up in various configurations. Configuration 1 & 2 are given in Table II. In configuration 1, the ratio of LG:ST:SM is 1:3:3 with a total capacity of 30, 7000 MIPS. In configuration 2, the ratio of LG:ST:SM is 3:2:2 with a total capacity of 39, 6000 MIPS.

The Bag-of-tasks workload is generated synthetically based on characteristics outlined in [11]. Table III gives the parameters used to generate a synthetic workload. In each run of the simulation, a daily cycle of 48 intervals is simulated. The user submitting the BoT is generated using a Zipf distribution. The inter-arrival time between each BoT in a given interval is modeled using a Weibull Distribution. The size of BoTs is assigned to be powers of 2 between 4 to 512 also modeled using a Weibull Distribution. The average task runtime within a BoT is given using a Normal distribution and the task runtime variability of the tasks is given using a Weibull Distribution. The runtime is specified in MI (Million

Instructions. The arrival rate is used to control the system load, and for this simulation, it is set to 1 The VMs operate on a time-shared basis on the hosts. The tasks (cloudlets) allotted to the VMs are also time-shared and hence, context switching takes place when multiple tasks take turns in using the processing element.

TABLE II. CONFIGURATION OF THE CLOUD

| Configuration 1 | | Configuration 2 | |
|---|---|---|---|
| | MIPS | | MIPS |
| VM #0 | 82000 | VM #0 | 82000 |
| VM #1 | 49000 | VM #1 | 82000 |
| VM #2 | 49000 | VM #2 | 82000 |
| VM #3 | 49000 | VM #3 | 49000 |
| VM #4 | 26000 | VM #4 | 49000 |
| VM #5 | 26000 | VM #5 | 26000 |
| VM #6 | 26000 | VM #6 | 26000 |
| **Total Capacity** | **307000** | **Total Capacity** | **396000** |
| **LG:ST:SM = 1:3:3** | | **LG:ST:SM = 3:2:2** | |

TABLE III. WORKLOAD CHARACTERISTICS OF BOT WORKLOAD

| | |
|---|---|
| **Number of Intervals in each Run** | **48** |
| **User (Zipf Distribution)** | (368,1.31) |
| **Interarrival Time (IAT) Weibull Distribution** | (4.25,7.86) |
| **Size of BoT (Weibull Distribution)** | pow ((1.76,2.11),2) |
| **Class of Bots** | {4, 8, 16, 32, 64, 128, 256, 512} |
| **Average Task Runtime within BoT (Normal Distribution)** | (2.73,6.1) |
| **Task Runtime Variability (Weibull Distribution)** | (2.05,12.05) |
| **Arrival rate** | 1 |
| **Runtime factor** | 1000,000 |

## IV. RESULTS AND DISCUSSION

The private cloud is set up in cloudsim and the BIS agent runs like a thread. BISA recognizes the state, decides on the action and updates the rewards and number of plays on a state. Fig. 2 shows a sample set of states observed after 10 runs comprising of a daily cycle of 48 cycles per run. The size ratio, along with the utilization of processing capacity in the various categories of VMs forms a state. State 30 with size_ratio as 1 and all the categories of VM utilization being "High" has the highest number of plays followed by State 1 and 18 that also have size ratio as 1 but with varying VM utilization.

Fig. 3 presents the normalized makespan after every run for a succession of 10 runs with different values for $\alpha$ and $\beta$. Each run consists of 48 intervals with BoTs generated using characteristics in Table II. Every marker represents the normalized makespan after a run. The normalized makespan is plotted against the number of tasks executed after ever run rather than run 1 to 10 so that the trend can be observed accurately.

Fig. 2.   The Number of Plays, $\tau$ on Various States at the end of 10 Runs (Daily Cycle of 48 Intervals Each).



Fig. 3.   Normalized Makespan Over 10 Runs for 6 Sets of Learning Parameters (Configuration 1-Tasks Time-Shared in VM).

For all learning parameters, a logarithmic decrease in the makespan can be observed. In some cases, the decrease is steeper. For $\alpha$=0.3 and $\beta$=0.5 the likelihood that the agent will rely on a previously learned policy that produces a better schedule is slightly more than its behavior to explore other possibilities. The normalized makespan fluctuates with every run and even after 33720 tasks at the end of the $10^{th}$ run. If a high value is set for both the exploration and exploitation parameters at $\alpha$=0.8 and $\beta$=0.8, this results in the slightly higher normalized makespan in most runs. If an exploration value is comparatively much higher than exploitation value in the case of $\alpha$=0.8 and $\beta$=0.2, this prevents the agent from getting stuck in any local optima and ensures all possibilities are explored. It can be observed that there is a steep learning curve in the initial 3 runs, but because of the high exploration value, the agent keeps choosing different heuristics and the makespan could not be reduced further. The agent produces surprisingly good results for a comparatively high exploitation value as compared to exploration value with $\alpha$=0.2 and $\beta$=0.8. The normalized makespan in the first run itself is 23% better than the lowest recorded value when $\alpha$=0.3 and $\beta$=0.3. Likewise, the learning curve is also steepest among all the sets. The lowest recorded makespan is 7 times lower than all other normalized makespan produced for any other set. Lastly, a very high exploitation value was set with $\alpha$=0.12 and $\beta$=0.98. The agent stabilized after 4 runs at a suboptimal normalized makespan but produces the largest recorded value for normalized makespan in the second run.

The average normalized makespan of the last 5 runs in a total of 10 runs on BISA with the different set of learning parameters is presented in Fig. 4. The average of the last 5 runs is taken as the agent is expected to stabilize. It can be observed that BISA with $\alpha$=0.2 and $\beta$=0.8 produced the best results. All other sets of learning parameters produce suboptimal results. For this configuration, U_MinCT produces the best results among the 15 scheduling heuristics. BISA with $\alpha$=0.2 and $\beta$=0.8 has an average normalized makespan that is 72% smaller than U_MinCT.

BISA with $\alpha$=0.8 and $\beta$=0.8 which has the largest average normalized makespan has an average normalized makespan which is 1.8 times larger than U_MinCT and performs better than 9 out of the 15 scheduling heuristics.

BISA was also tested in a second configuration as outlined in Table II. As can be seen in Fig. 5, the steepest learning curve is observed for $\alpha$=0.2 and $\beta$=0.8. The best scheduling heuristic is U_MinCT and BISA ($\alpha$=0.3 and $\beta$=0.5) produces a 39% larger average normalized makespan as compared to U_MinCT. On average, BISA produces suboptimal results and results in average normalized makespan that is less than 50% higher than the best (Fig. 6).

To test the hypothesis that BISA will be able to adapt to dynamic changes to the resources, BISA is run with the best performance learned from configuration 1 ($\alpha$=0.2 and $\beta$=0.8) to configuration 2 (Fig. 7).

Fig. 4.    Comparison of Average Normalized Makespan of Last 5 Runs (Configuration 1–Tasks Time-Shared in VM).



Fig. 5.    Normalized Makespan Over 10 Runs for 6 Sets of Learning Parameters (Configuration 2-Tasks Time-Shared in VM).



Fig. 6:   Comparison of Average Normalized Makespan of last 5 runs (Configuration 2 –tasks time-shared in VM)

Fig. 6.    Comparison of Average Normalized Makespan of Last 5 Runs (Configuration 2-Tasks Time-Shared in VM).



Fig. 7.   Comparison of Normalized Makespan of the First Run of BISA Run aFresh on Configuration 2 (C2) with Learned Rewards and Preferences from Configuration 1(C1) (α=0.2 and β=0.8) on Configuration 2.

The first normalized makespan learned afresh in configuration 2 is compared with the first normalized makespan in configuration 2 learned from configuration 1 on $\alpha=0.2$, $\beta=0.8$. The normalized makespan for the set ($\alpha=0.3$ and $\beta=0.5$) operating on learned heuristics is half of the first normalized makespan when BISA is run on Configuration 2 afresh. The other sets do not show a comparable difference. It can also be observed that when using the same learning parameters $\alpha=0.2$, $\beta=0.8$ when transitioning from configuration 1 to configuration 2 a negative result of increased makespan is observed. This indicates the state variables used in this simulation is not enough to differentiate. Hence in further work, states will be defined in further details so that learned best heuristics on a configuration can seamlessly be applied to other underlying configurations.

## V. Conclusions and Future Work

An intelligent agent, BISA (BoT Intelligent Scheduling Agent) is proposed that learns the best scheduling heuristic to use in a state. The goal of BISA is to schedule each BoT such that the makespan is minimal. It is tested in two underlying configurations of cloud by allotting both tasks to VMs in a time-shared manner. It produces sub-optimal results comparable to the best scheduling heuristic for a given configuration. It can adapt to different underlying configurations and re-learn the best heuristic to use in a configuration. In configuration 1, for one set of learning parameter ($\alpha=0.2$ and $\beta=0.8$) BISA outperforms all the traditional scheduling heuristics. In configuration 2 all the sets of learning parameters result in sub-optimal results with the best results on the set ($\alpha=0.3$ and $\beta=0.5$). A 50% decrease in normalized makespan can be observed when transitioning from the best run on configuration 1 to configuration 2 on ($\alpha=0.3$ and $\beta=0.5$) as compared to running BISA afresh in configuration 2.

It can be inferred that the choice of learning parameters determines the effectiveness of BISA. A good tradeoff between exploration and exploitation produces a lower makespan and converges to a near-best heuristic for a state rapidly. With BISA there is always the possibility of agents settling for a sub-optimal heuristic and it is crucial to choose learning parameters to minimize this occurrence.

BISA could be improved by re-examining the parameters that comprise the state and specifying them in a fine-grained manner. Additional objectives of minimizing cost, latency, and energy consumption, if factored into the reward obtained in BISA, solves a multi-objective problem. BISA trained to learn the best learning parameters that result in optimal results would also improve the performance of BISA in converging to and choosing the optimal scheduling heuristic in a state.

### References

[1] R. Buyya, S. K. Garg, and R. N. Calheiros, "SLA-oriented resource provisioning for cloud computing: Challenges, architecture, and solutions," Proc. - 2011 Int. Conf. Cloud Serv. Comput. CSC 2011, no. Figure 1, pp. 1–10, 2011.

[2] G. Andreadis, L. Versluis, F. Mastenbroek, and A. Iosup, "A Reference Architecture for Datacenter Scheduling: Design, Validation, and Experiments," in Proceedings of the International Conference for High Performance Computing, Networking, Storage, and Analysis, 2018, p. 37.

[3] B. Varghese and R. Buyya, "Next generation cloud computing: New trends and research directions," Futur. Gener. Comput. Syst., vol. 79, pp. 849–861, 2018.

[4] T. Thein, M. M. Myo, S. Parvin, and A. Gawanmeh, "Reinforcement learning based methodology for energy-efficient resource allocation in cloud data centers," J. King Saud Univ. - Comput. Inf. Sci., 2018.

[5] W. Iqbal, M. N. Dailey, D. Carrera, and P. Janecek, "Adaptive resource provisioning for read intensive multi-tier applications in the cloud," Futur. Gener. Comput. Syst., vol. 27, no. 6, pp. 871–879, 2011.

[6] A. Benoit, L. Marchal, J. F. Pineau, Y. Robert, and F. Vivien, "Scheduling concurrent bag-of-tasks applications on heterogeneous platforms," IEEE Trans. Comput., vol. 59, no. 2, pp. 202–217, 2010.

[7] P. S. H. Darius and E. G. M. Kanaga, "Bag-of-Tasks Intelligent Scheduling Agent ( BISA ) in Cloud Computing," Adv. Comput. Commun. Paradig., pp. 239–246, 2018.

[8] O. H. Ibarra and C. E. Kim, "Heuristic Algorithms for Scheduling Independent Tasks on Nonidentical Processors," J. ACM, vol. 24, no. 2, pp. 280–289, 1977.

[9] I. A. Moschakis and H. D. Karatza, "A meta-heuristic optimization approach to the scheduling of bag-of-tasks applications on heterogeneous clouds with multi-level arrivals and critical jobs," Simul. Model. Pract. Theory, vol. 57, pp. 1–25, 2015.

[10] R. N. Calheiros and R. Buyya, "Energy-efficient scheduling of urgent bag-of-tasks applications in clouds through DVFS," Proc. Int. Conf. Cloud Comput. Technol. Sci. CloudCom, vol. 2015-Febru, no. February, pp. 342–349, 2015.

[11] A. Iosup, O. Ozan Sonmez, S. Anoep, and D. Epema, "The performance of bags-of-tasks in large-scale distributed systems," in Telecommunications Policy - TELECOMMUN POLICY, 2008, pp. 97–108.

[12] M. Maheswaran, A. Shoukat, H. J. Siegel, D. Hensgen, and R. F. Freund, "Dynamic Matching and Scheduling of a Class of Independent Tasks Onto Heterogeneous Computing Systems," in Proceedings of the Eighth Heterogeneous Computing Workshop, IEEE Computer Society, 1999, p. 30.

[13] W. Cirne, D. Paranhos, F. Brasileiro, L. Fabrício, and W. Góes, "On the Efficacy, Efficiency and Emergent Behavior of Task Replication," Large Distrib. Syst. Parallel Comput., vol. 33, no. 3, pp. 213–234, 2007.

[14] J. O. Gutierrez-Garcia and K. M. Sim, "A family of heuristics for agent-based elastic Cloud bag-of-tasks concurrent scheduling," Futur. Gener. Comput. Syst., vol. 29, no. 7, pp. 1682–1699, 2013.

[15] Z. Zheng, R. Wang, H. Zhong, and X. Zhang, "An approach for cloud resource scheduling based on parallel genetic algorithm," ICCRD2011 - 2011 3rd Int. Conf. Comput. Res. Dev., vol. 2, pp. 444–447, 2011.

[16] H. Zhong, K. Tao, and X. Zhang, "An approach to optimized resource scheduling algorithm for open-source cloud systems," Proc. - 5th Annu. ChinaGrid Conf. ChinaGrid 2010, pp. 124–129, 2010.

[17] H. C. Lau, W. C. Wan, S. Halim, and K. Toh, "A software framework for fast prototyping of meta-heuristics hybridization," Int. Trans. Oper. Res., vol. 14, no. 2, pp. 123–141, 2007.

[18] S. Russell and P. Norvig, Artificial Intelligence: A Modern Approach, 3rd ed. Upper Saddle River, NJ, USA: Prentice Hall Press, 2009.

[19] R. S. Sutton and A. G. Barto, Introduction to Reinforcement Learning, 1st ed. Cambridge, MA, USA: MIT Press, 1998.

[20] R. N. Calheiros, R. Ranjan, A. Beloglazov, C. A. F. De Rose, and R. Buyya, "CloudSim: A Toolkit for Modeling and Simulation of Cloud Computing Environments and Evaluation of Resource Provisioning Algorithms," Softw. Pr. Exper., vol. 41, no. 1, pp. 23–50, Jan. 2011.

# Four-Class Motor Imagery EEG Signal Classification using PCA, Wavelet and Two-Stage Neural Network

Md. Asadur Rahman[1]

Department of Biomedical Engineering
Khulna University of Engineering & Technology (KUET)
Khulna-9203, Bangladesh

Farzana Khanam[2]

Department of Biomedical Engineering
Jashore University of Science and Technology (JUST)
Jashore-7408, Bangladesh

Md. Kazem Hossain[3], Mohammad Khurshed Alam[4], Mohiuddin Ahmad[5]

Department of Electrical and Electronic Engineering, Khulna University of Engineering & Technology (KUET)
Khulna-9203, Bangladesh

*Abstract*—Electroencephalogram (EEG) is the most significant signal for brain-computer interfaces (BCI). Nowadays, motor imagery (MI) movement based BCI is highly accepted method for. This paper proposes a novel method based on the combined utilization of principal component analysis (PCA), wavelet packet transformation (WPT), and two-stage machine learning algorithm to classify four-class MI EEG signal. This work includes four-class MI events by an imaginary lifting of the left hand, right hand, left foot, and Right Foot. The main challenge of this work is to discriminate the similar lobe EEG signal pattern such as left foot VS left hand. Another critical problem is to identify the MI movements of two different feet because their activation level is very low and show an almost similar pattern. This work firstly uses the PCA to reduce the signal dimensions of the left and right lobe of the brain. Then, WPT is used to extract the feature from the different class EEG signal. Finally, the artificial neural network is trained into two stages – 1st stage identifies the lobe from the signal pattern and the 2nd stage identifies whether the signal is of MI hand or MI foot movement. The proposed method is applied to the 4-class MI movement related EEG signals of 15 participants and found excellent classification accuracy (>74% on average). The outcomes of the proposed method prove its effectiveness in practical BCI implementation.

*Keywords*—*Brain-computer interface; electroencephalogram; motor imagery; principal component analysis; wavelet packet transformation; artificial neural network; classification*

## I. INTRODUCTION

Brain-computer interface (BCI) creates a communication system between computer and brain functionality to control the other devices. BCI can be used to control devices like a wheelchair, room light, fan, etc. that assists physically challenged people. There are many modalities such as Electrocorticogram, Electroencephalogram (EEG), Functional near-infrared Spectroscopy, Magneto-encephalogram, Functional Magnetic Resonance Imaging, etc. to read the brain signal. Among these modalities, EEG is the most familiar and cheapest technique to record brain signal [1].

EEG is a non-invasive way to record neuro-electric signal which is often called an EEG signal. Since a number of brain stimuli have been proposed for implementing the BCI system. Among all other stimuli, motor imagery (MI) movement is the highest choice for the researchers [2]. MI has a special benefit because it needs no additional setup like visual stimuli [3].

An MI movement is a process where a candidate imagines the real movement execution and corresponding neuro-electric activities are recorded by the EEG modality. There are different types of EEG-based BCI like as simple and compound limb motor imagery [4], continuous arm movement from EEG signals [5] and individual finger movements from one hand using human EEG signals [6], etc. In the present works of literature of MI, most of the research works [7-12] are related to two-class or three-class such as left hand vs. right hand and left hand, right hand, and foot, respectively. Here both feet are considered as a single class.

Multiclass MI movement *i.e.*, more than 3 class classification task is always a challenging issue when the discrimination is necessary between two-foot movement. Two hands and one-foot MI movement classification provide acceptable classification accuracy though, four-class MI movement *i.e.*, imagery left hand (iLH), right hand (iRH), left foot (iLF), and right foot (iRF) classification is still not an acceptable range. To meet this challenge, a very handful research works [13-15] have been proposed those considered the fourth class as tongue movement, not an approach to discriminate between two imagery feet-based EEG signal along with the other two imagery hand. The approach of excluding the two-foot is due to having an almost similar pattern and neuro-activation of imagery feet movements. The precise difference between the patterns of two-foot movements is the main challenge to attain the acceptable classification accuracy regarding 4-class MI movement classification. Therefore, very wise feature selection and classifier models are two important tasks to meet this 4-class EEG signal classification regarding MI movements.

A number of feature extraction methods from EEG signal have been proposed in recent years. Within these feature extraction techniques, there are autoregressive (AR) methods [16, 17], wavelet transforms (WT) method [17-21], and phase-space reconstruction approach [10], CSP based methods [13, 14, 22], empirical mode decomposition [23-25], etc. For a wide range of pattern recognition, wavelet packet transformation (WPT) provides excellent time-frequency features. The WPT coefficients during EEG signal decomposition are widely used in EEG signal classification. But, for multiple class, the WPT based feature extraction has two important limitations: *i*) Structuring the features and *ii*) Selection of the bases [26]. The features are structured by WPT coefficients those are considered to yield the significant pattern of the different classes EEG signal. Besides the feature structuring, the proper base selection is the other step by which the structured features can show the highest discriminative characteristics among the classes. As a result, the WPT based features of the four classes are not enough criteria for satisfying classification accuracy. Therefore, there is a scope to add some innovative approach to meet the challenge.

In this paper, we have proposed a novel method for the aforementioned multiple class EEG signal classification utilizing a principal component analysis (PCA), WPT, and two-stage artificial neural network (ANN). The proposed method utilizes the concept of PCA to reduce the signal dimension. The reduced dimension EEG signals are fed to the WPT algorithm to extract its features. After that, we have modeled two ANN sequentially for identification of the lobe-origin and limb-origin of the signal. The proposed ANN-based two-stage model finds the signal's origin at first, *i.e.*, is the signal comes from either left lobe or right lobe and secondly, it decides whether the signal is of either lower limb or upper limb, i.e., foot or hand. In this work, four-class motor imagery EEG data were collected in our laboratory from 15 participants. The data were preprocessed and separated according to the tasks. Then the dimensionality reduction and feature extraction were conducted by PCA and WPT, respectively. Finally, the two-stage ANN model was trained with the training data set. The proposed method was finally applied on the testing data set and we found that the proposed algorithm can improve classification performance (on average =11.25%) than the WPT-single ANN model.

The rest of the article is organized as: The materials used in this research work are described in Section II. The proposed methodology is elaborately explained in Section III. The results are presented with related discussions in Section IV. Finally, the whole work has been concluded briefly in Section V.

## II. Materials

### A. Participants

Fifteen healthy adult male subjects (age=23±2.5 years) participated in this experiment. The participants did not have any psychological disorder and they were visually corrected person. In addition, all participants were not in any medication last one month. The participants claimed themselves having no pain in muscle and no mental disorder. All the subjects were right-handed based on Edinburg Handedness Inventory to avoid the variation of EEG signal pattern in left and right lobes during data collection. The participants were verbally informed about the protocol and they also practiced the protocol several timed prior to proceed the actual data acquisition session. Their written consent was taken from all the participants prior to the EEG data collection. Since it was a volunteer-based contract, no participant was paid for the data collection.

### B. Data Acquisition Protocol

The subject was asked to sit on a chair before a computer screen at a convenient distance. Then he was told to take rest for a while to reduce the physiological artifacts which can be accumulated with the electrical activities of the brain. There were four sequential tasks that every participant had to do. These tasks were: iLH, iRH, iLF, and iRF movements. The experimental protocol is sketched in Fig. 1. The protocol was started with 5s rest. It was followed by 10s task and 20s rest. The tasks were performed, sequentially as given in Fig. 1 by the participants. One participant performed 5 trials per each session and 4 sessions per day. Each participant participated in the data acquisition in two different data in a single week. Eventually, forty trials of each task were acquired form each participant. The protocol does not violate the Declaration of Helsinki [27].

### C. Data Acquisition and Data Description

The EEG data were acquired utilizing the 9-channel B-Alert X10 wireless device in the Neuroimaging Laboratory, Department of Biomedical Engineering, Khulna University of Engineering & Technology (KUET). The data logging in computer memory was conducted by Acqknowledge 4.4 software [28]. The acquisition sampling rate of the data was 256 Hz. The collected 9 channel EEG data covers the following position of the scalp F3, Fz, F4, C3, Cz, C4, P3, Pz, and P4 according to the international 10/20 system. Therefore, this device covers the frontal, central, and parietal lobe of the brain during EEG data acquisition. The data acquisition procedure of this research work with the 9-channel B-Alert X10 wireless device is graphically presented in Fig. 2.

| Rest 20 Sec | iLH 10 Sec | Rest 20 Sec | iRH 10 Sec | Rest 20 Sec | iLF 10 Sec | Rest 20 Sec | iRF 10 sec |
| --- | --- | --- | --- | --- | --- | --- | --- |

Fig. 1.    Data Acquisition Protocol with a Time Schedule of the Tasks.

Fig. 2. MI EEG Data Acquisition Procedure Utilizing the B-Alert X-10 Wireless Device with a Computer-Assisted Command.

## III. METHODS

### A. Proposed Method

Our proposed method is an offline EEG signal classification method. After the raw EEG signal acquisition, several processing steps were conducted to prepare the signal for machine understanding format. At first, the raw EEG signals were filtered with a 50 Hz notch filter to remove the power line noise from the signal. A $2^{nd}$ order band stops IIR filter was designed considering its cut off band from 49 to 51 Hz and utilized as a notch filter for the EEG signal. Secondly, the signals were filtered with $3^{rd}$ order bandpass IIR filter with a passband from 4 to 100 Hz. The eye blink effect was removed from the signal utilizing the method described in [29, 30] with their Matlab based automatic EEG artifact rejection toolbox, EAWICA.

After filtering the data, the EEG signals were separated according to the tasks. These filtered data were used as input of our proposed method to model the 4-class MI movement prediction. The proposed method to construct the ANN-based predictive model can be presented by the following block diagram given in Fig. 3.

A complex with several steps is shown in Fig. 3 where the input EEG signal is considered as preprocessed and separated for further analysis. Then the dimension on the signals is reduced applying the PCA. Since the EEG signals were acquired nine channel device, the left lobe channels F3, C3, and P3 are reduced along with the neutral positions Fz, Cz, and Pz to a single channel with the combination of these six channels with their correlation weight utilizing the PCA concept. On the other hand, the right lobe channels F4, C4, and P4 are also reduced in the same way with PCA. Therefore, applying the PCA the 9-channel EEG signal was reduced to two channel signals. During PCA, we considered only the highest PCA ($1^{st}$ principal components) factors.

According to the proposed method, WPT was used to extract the features from these two reduced signals of every task, separately. After that the features of 4-class have been broadly divided into two classes based on the lobe of operation *i.e.*, *i*) iLH & iLF and *ii*) iRH & iRF. These wider two classes of features were used to train the ANN to predict the input signal whether it is of left or right. Additionally, another two classes were formed by the features of iLH+iRH and iLF+iRF and trained the ANN so that it can predict whether the signal is off hand or foot.

Therefore, two separate ANN were utilized to make a hybridized two-stage ANN predictive model that would be able to predict a scratch signal of any of the four classes. For any MI EEG signal, this model will predict the left or right limb signal, firstly and then in the second stage, the model will find it as a signal of hand or foot. The classification procedure of the two-stage ANN is given in Fig. 4. In this figure, the block preprocessing includes the filtering and dimensionality reduction by PCA as explained before.



Fig. 3. The Block Diagram of the Proposed Method for Training and Constructing the Two-Stage ANN-based Predictive Model Along with PCA and WPT based Feature Extraction.

Fig. 4.    The Proposed Two-Stage ANN-based Predictive Model and its Classification Technique from a Raw Signal to Decision.

## B. Principal Component Analysis

A matrix $\Delta$ consists of data of *n* dimension. Now, a matrix *H* can be designed which characterizes the eigenvectors arranged as the eigenvalues of the covariance matrix of $\Delta$. Then we get the PCA of the data $\Delta$ in the form of *P* as,

$$P = H^T \Delta \tag{1}$$

The eigenvectors can also be named as the principal components. To project the data, if first rows (*R*) of *P* are selected, the data becomes of *R* dimensional from *d* dimensions. This transformation is performed by singular value decomposition (SVD). Matrix decomposition can describe the procedure to perform PCA by SVD. Suppose, the matrix, $\Delta$ can be decomposed using SVD as

$$\Delta = \Psi \kappa \theta^T \tag{2}$$

Here, $\Psi$ is an *n×m* matrix with orthonormal columns ($\Psi^T \Psi = I$); $\theta$ is an *m×m* orthonormal matrix ($\theta^T \theta = I$), and $\kappa$ is an *m×m* diagonal matrix with positive or zero elements which is recognized as a singular value. Besides, the covariance matrix can be calculated, C of $\Delta$ as,

$$C = \frac{1}{N} \Delta \Delta^T = \frac{1}{N} \Psi \kappa^2 \Psi^T \tag{3}$$

As the singular values are organized in descending order and if *n < m*, the first n columns in $\Psi$ corresponds to the sorted eigenvalues of matrix C and if *m≥n*, the first m corresponds to the sorted non-zero eigenvalues of C. Therefore, eventually the transformed data can be written as,

$$P = H^T \Delta = H^T H \kappa \theta^T \tag{4}$$

## C. Wavelet Packet Transformation for Feature Extraction

The WPT is a concept that is different from conventional wavelet transformation (WT). WPT decomposes both the approximate coefficients and the detailed coefficients. The WPT may be considered as a subspace tree. We can present the original signal as $\Pi_{0,0}$ which reflects the root mood of the tree in the original signal space. Generally, the notation j and k in $\Pi_{j,k}$ denotes the scale and sub-band space. The WPT decomposes an original signal $\Pi_{j,k}$ into two different subspaces: an approximation space $\Pi_{j,k} \to \Pi_{j+1,2k}$ and a detailed space, $\Pi_{j,k} \to \Pi_{j+1,2k+1}$. This space decomposition utilizes the concept of dividing the orthogonal basis function $\{\phi_j(t - 2^j k)\}_{k \in Z}$ of the original signal space into two new orthogonal bases, i) $\{\phi_{j+1}(t - 2^{j+1}k)\}_{k \in Z}$ of approximate space $\Pi_{j+1,2k}$ and ii) $\{\psi_{j+1}(t - 2^{j+1}k)\}_{k \in Z}$ of detailed space $\Pi_{j+1,2k+1}$. Here $\Pi_{j,k}(t)$ and $\Psi_{j,k}$ represents the scaling and wavelet functions, respectively. These functions are equated as [26]:

$$\phi_{j,k}(t) = \frac{1}{\sqrt{12^j 1}} \phi(\frac{t - 2^j k}{2^j}) \tag{5}$$

$$\psi_{j,k}(t) = \frac{1}{\sqrt{12^j 1}} \psi(\frac{t - 2^j k}{2^j}) \tag{6}$$

Fig. 5. Graphical Representation of Wavelet Packets Decomposition Method that Decomposes $\Pi_{0,0}$ into Tree-Structured Subspaces.

Here $2^j$ is the scaling parameter that measures the scaling or compression degree of the original signal. In addition, $2^j k$ is the location parameter or translation parameter that indicates the time location of the wavelet. The aforesaid process can be repeated $J$ times, where $J$ must be less than $\log_2 N$. Here, $N$ is the total number of samples in the original signal. This process of WPT $J \times N$ founds coefficients. Therefore, at any level of transformation $j$ $[j = 1,2,...,J]$, the tree has $N/(2j)$ coefficient blocks. This iterative process in a WPT can be treated as a tree-like structure, where the tree nodes represent the subspaces of different frequency localization characteristics. The corresponding decomposition procedure can be presented as Fig. 5 [26, 31].

### D. Artificial Neural Network

ANN replicates the functional concept of the human brain. The multilayer feedforward ANN has three basic layers: an input layer, an output layer, and a hidden layer. A typical model of a feed-forward network with its prominent layers is given in Fig. 6. The four outputs are chosen in the output layer because total data are to be classified by this work into four classes.

In a supervised neural network, the input layer is formed based on the size of the features and the output layer is chosen upon its number of classes to be classified. The significant layer is a hidden layer which is connected to the input and output with single or multiple layers while in case of multiple layers are often referred to multilayer neural network. The connected manners are actually a mathematical function with a predefined function which is also known as a neuron. Generally, a neuron at $j$ label receives an input $p_j(t)$ from the previous neuron at a discrete time, $t$. Suppose, the activation function of the neuron at this level is $a_j(t)$ where a threshold, $\theta_j$ is also chosen. Therefore, the activation function, $\zeta$ computes the next level activation, $a_j(t+1)$ from the current information as [32],

$$a_j(t+1) = \zeta\big(a_j(t), p_j(t), \theta_j\big) \tag{7}$$

Inside an ANN, the output of a neuron $i$ performs as an input to a neuron $j$ and each connection is assigned with a weight $w_{ij}$ with a bias term $w_{0j}$ (sometimes). A propagation function calculates the input of the neuron $j$ from the output, $o_i(t)$ of the neuron $i$ with the assigned bias value as [33],

$$p_j(t) = \sum_i o_i(t) w_{ij} + w_{0j} \tag{8}$$

In a supervised learning method of an ANN, a set of given pairs of input features, $x$ and output, $y$ ( $x \in X$ , $y \in Y$ ) are fed to find a function, $\Lambda : X \to Y$ subject to the assumed class of functions. By this predictive function, it is expected that it could be able to infer the applied data mapping. A cost function is used to relate the expected and actual mismatch of the inferring data map that has prior knowledge about the domain of interest. A commonly used cost function for the ANN is a mean-square error that aims to minimize the mean square error between $\Lambda(x)$ and the target $y$. In a multilayer perceptron, ANN network utilizes the gradient descent algorithm to optimize the cost function by the backpropagation algorithm.



Fig. 6. A Typical Model of an ANN with the Structure of its basic Layers.

## IV. RESULTS AND DISCUSSIONS

The results regarding the proposed analytical methods were demonstrated by Matlab 2018a [34]. Utilizing the Acqknowledge software the raw EEG data were converted to .mat file to make it compatible with Matlab. The filtering, feature extraction, classification, etc. processing was conducted utilizing the different toolboxes of Matlab 2018 in an offline fashion. The raw signal was filtered with several steps as notch filtering, bandpass filtering, and eye blink removal. The following procedures were applied and the resulting effects on the EEG signal are presented in Fig. 7. Then, the dimensionality reduction of the six-channel EEG signal was performed and one-dimensional EEG signal is prepared from a lobe using PCA. PCA converts the signal taking the maximum variations in the signal and avoiding the similarity among the signal which helps to get maximum frequency contents without taking the curse of dimensionality. The resulting signals are graphically presented in Fig. 8.



Fig. 7. Preprocessing Steps.



Fig. 8. The Proposed Method of Dimensionality Reduction of the Channels by PCA.

TABLE I.        TOPOPLOTS OF THE FOUR PARTICIPANTS REGARDING DIFFERENT IMAGERY STIMULI



From the strength of the signal power regarding four imagery movements, we found various kinds of pattern. The EEG signals of four randomly chosen participants were analyzed to observe the neural activation pattern due to the different applied stimuli. The neural activations of the aforesaid conditions were prepared based topographic plot or topoplot on demo human scalp. Utilizing the open-source Matlab based function of topoplot (which is solely designed for the 9 channel B-alert EEG data and available in [35]); the activation based topolots of aforesaid stimuli are given in Table I. The resulting topoplots demonstrate that the imagery hands movements (both iLH and iRH) show consistent pattern in their neural activations. On the other hand, there are some discrepancies in the neural activations for imagery feet movements. This is the cause we utilized the two stages ANN training method to recognize this pattern with discrepancies of the feet movements.

In the training session of the ANN, the wavelet tree coefficients were used as features. The training and testing features were separated as a 5-fold cross-validation technique. The classification accuracies were calculated utilizing the proposed method as a subject dependent approach. It is found during the training process that the first stage training performance is better than the second stage training. It may occur due to the similarity of the EEG based neural activation of the same lobe either from hand or foot. The training sessions of the features of subject 1 are given in Fig. 9. In this figure, the stage I training and validation performances are given in Fig. 9(a) and the stage II training and validation performances are given in Fig. 9(b). It is clear from the figures that, the performances in case of stage II are inferior to that of the stage I. It has a significant impact on the classification accuracy. Such training and testing were conducted for an individual subject to present the classification accuracy.

Fig. 10. The Training and Testing Ratio of the 5-Fold Cross-Validation Technique.

Since there are 40 trials of each imagery task, according to the 5-fold cross-validation technique 32 trials were used to train and validation of the network and the rest 8 trials of each task were used to test the accuracy of the trained network. The selection of the training and testing trials from the 40 trials were performed 5 different sets as the presentation in Fig. 10. The final classification accuracy was estimated from the average value of the five testing results. Using this consideration, the classification accuracy of the subject one was calculated and the classification performance of the proposed method is given by a confusion matrix in Fig. 11. This result is of the 1st iteration of the 5-fold cross-validation technique where we found that the classification accuracy is 81.3%. Similar four more iterations were performed to get the average classification accuracy. We found the classification accuracy for rest four trials as 68.75%, 78.2%, 71.9%, 65.7%, respectively. Therefore, the average four-class classification accuracy for participants is 73.17%.



Fig. 9. Training and Testing Performances of the Proposed Network in its Stage I (a) and Stage II (b) for Participants 1.



Fig. 11. Confusion Matrix of Classification Accuracy of 1st Iteration Among the 5 Iterations.

TABLE II.    CLASSIFICATION ACCURACIES FOUND UTILIZING THE SINGLE-STAGE AND TWO-STAGE ANN CLASSIFIERS

| Sub. Num. | Classification Accuracies (%) | |
|---|---|---|
| | *One-stage ANN* | *Two-stage ANN* |
| 1 | 62.24 | 73.17 |
| 2 | 63.50 | 74.65 |
| 3 | 52.5 | 68.25 |
| 4 | 60.25 | 81.46 |
| 5 | 58.24 | 75 |
| 6 | 45.8 | 65.5 |
| 7 | 50.8 | 71.25 |
| 8 | 60.45 | 75.50 |
| 9 | 62.50 | 78.25 |
| 10 | 58.90 | 72 |
| 11 | 58.90 | 72 |
| 12 | 61.45 | 82.17 |
| 13 | 62.45 | 86.25 |
| 14 | 60.50 | 73.40 |
| 15 | 55.5 | 70.25 |
| Average ±std | 58.26±5.04 | 74.60±5.50 |

Furthermore, the similar feature extraction method was applied for classification utilizing stage ANN classifier. The results due to the conventional one-stage classifier and the proposed two-stage classifiers for 15 participants are given in Table II. Here the average classification accuracies from 5-fold cross-validations are considered. The outcomes suggest that the two-stage classification accuracies are better than that of the conventional one-stage classification method with the same features. In average, the two-stage classification technique provides 74.60% accuracies whereas, the single-stage classifiers give use 58.26 % classification accuracies regarding the four-class motor imagery EEG signal.

## V.    CONCLUSIONS

By this proposed work, it has been found that not only the innovative feature extraction is mandatory but also the classifier setup with an appropriate approach is a concerning issue. The feature extraction method utilizing the PCA based wavelet packet transformation is although an excellent approach to find the properties of the EEG signal; it proves failure to classify the four-class motor imagery signals with satisfactory accuracy. On the other hand, the proposed two-stage classifier improves the classification accuracy at 16.34% on average which is a remarkable outcome. This approach can be applied in any higher-class classifier than two-class. Therefore, this proposal hopefully outtakes the conventional approach in practical BCI implementation.

One potential drawback of the proposed multi-stage training-based classifier is its time requirement for the training stage. If the classification accuracy becomes the first priority, the said limitation can be avoided in the implementation.

REFERENCES

[1] L. F. Nicolas-Alonso and J. Gomez-Gil, "Brain computer interfaces, a review," Sensors, vol. 12, pp. 1211-1279, 2012.

[2] K. Hong, M. J. Khan, and M. J. Hong, "Feature extraction and classification methods for hybrid fNIRS-EEG brain-computer interfaces," Frontiers in Human Neuroscience, vol. 12, no. 246, 2018.

[3] K. Kitahara, Y. Hayashi, S. Yano, and T. Kondo, "Target-directed motor imagery of the lower limb enhance event-related desynchronization," PLOS One, vol. 12, no. 9, pp. 1-15, 2017.

[4] Y. Weibo, S. Qiu, H. Qi, L. Zhang, B. Wan, and D. Ming. "EEG feature comparison and classification of simple and compound limb motor imagery." J. Neuroeng. Rehabil 10 (2013).

[5] W. J. Seok, K. R. Muller, and S. W. Lee. "Classifying directions in continuous arm movement from EEG signals." 3rd International Winter Conference on BrainComputer Interface (BCI), 2015, pp. 1-2.

[6] K. Liao, R. Xiao, J. Gonzalez, and L. Ding. "Decoding individual finger movements from one hand using human EEG signals." PLoS One, vol. 9, no. 1, 2014.

[7] M. Li, W. Zhu, H. Liu, and J. Yang, "Adaptive feature extraction of motor imagery EEG with optimal wavelet packets and SE-isomap," Applied Sciences, vol. 7, no. 390, pp. 1-18, 2017.

[8] Y. Ma, X. Ding, Q. She, Z. Luo, T. Potter, and Y. Zhang, "Classification of motor imagery EEG signals with support vector machines and particle swarm optimization," vol. 2016, Article ID. 4941235, pp. 1-8, Computational and Mathematical Methods in Medicine, 2016.

[9] Z. Tang, C. Li, J. WU, P. LIU, and S. Cheng, "Classification of EEG-based single-trial motor imagery tasks using a B-CSP method for BCI," Frontiers of Information Technology & Electronic Engineering, 2018.

[10] R. Djemal, A. G. Bazyed, K. Belwafi, S. Gannouni, and W. Kaaniche, "Three-class EEG-based motor imagery classification using phase-space reconstruction technique," Brain Sciences, vol. 6, no. 36, pp. 1-19, 2016.

[11] L. Cao, B. Xia, O. Maysam, J. Li, H. Xie, and N. Birbaumer, "A synchronous motor imagery based neural physiological paradigm for brain-computer interface speller," Frontiers in Human Neuroscience, vol. 11, no. 274, pp. 1-9, 2017.

[12] J. Petersen, H. K. Iversen, and S. Puthusserypady, "Motor imagery based brain-computer interface paradigm for upper limb stroke rehabilitation," International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC), Honolulu, HI, 2018, pp. 1960-1963.

[13] L. Wang and X. Wu, "Classification of four-class motor imagery EEG data using spatial filtering," International Conference on Bioinformatics and Biomedical Engineering, Shanghai, 2008, pp. 2153-2156.

[14] A. Mahmood, R. Zainab, R. B. Ahmad, M. Saeed, and A. M. Kamboh, "Classification of multi-class motor imagery EEG using four band common spatial pattern," Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC), Seogwipo, 2017, pp. 1034-1037.

[15] S. Ge, R. Wang, and D. Yu, "Classification of four-class motor imagery employing single-channel electroencephalography," PLOS One, vol. 9, no. 6, pp. 1-7, 2014.

[16] Y. Zhang, X. Ji, and Y. Zhang, "Classification of EEG signals based on AR model and approximate entropy," International Joint Conference on Neural Networks (IJCNN), Killarney, 2015, pp. 1-6.

[17] Y. Zhang, B. Liu, X. Ji, D. Huang, "Classification of EEG signals based on autoregressive model and wavelet packet decomposition," Neural Processing Letter, vol 45, no. 2, pp. 365-378, 2017.

[18] H. U. Amin, A. S. Malik, R. F. Ahmad, N. Badruddin, N. Kamel, M. Hussain, and W-T. Chooi, "Feature extraction and classification for EEG signals using wavelet transform and machine learning techniques," Australasian Physical & Engineering Sciences in Medicine, vol. 38, no. 1, pp. 139-149, 2015.

[19] F. Sherwani, S. Shanta, B. S. K. K. Ibrahim, and M. S. Huq, "Wavelet-based feature extraction for classification of motor imagery signals," IEEE EMBS Conference on Biomedical Engineering and Sciences (IECBES), Kuala Lumpur, 2016, pp. 360-364.

[20] A. Sharmila and P. Mahalakshmi, "Wavelet-based feature extraction for classification of epileptic seizure EEG signal," Journal of Medical Engineering & Technology, vol. 41, no. 8, pp. 670-680, 2017.

[21] M. A. Rahman, M. M. O. Rashid, F. Khanam, M. K. Alam, and M. Ahmad, "EEG based brain alertness monitoring by statistical and artificial neural network approach," International Journal of Advanced Computer Science and Applications, vol. 10, no. 1, 2019.

[22] P. K. Saha, M. A. Rahman, and M. N. Mollah, "Frequency domain approach in CSP based feature extraction for EEG signal classification," International Conference on Electrical, Computer and Communication Engineering (ECCE), 7-9 February 2019, Cox's Bazar, Bangladesh.

[23] P. A. Munoz-Gutierrez, E. Giraldo, M. Bueno-Lopez, and M. Molonas, "Localization of Active Brain Sources From EEG Signals Using Empirical Mode Decomposition: A Comparative Study," vol. 12, no. 55, pp. 1-14, 2018.

[24] J. Sokhal, B. Garg, S. Aggarwal, and R. Jain, "Classification of EEG signals using empirical mode decomposition and lifting wavelet transforms," International Conference on Computing, Communication, and Automation (ICCCA), Greater Noida, 2017, pp. 1197-1202.

[25] T. Nazneen, M. A. Rahman, and M. N. Mollah, "Towards the effective intrinsic mode functions for motor imagery EEG signal classification," International Conference on Electrical, Computer and Communication Engineering (ECCE), 7-9 February 2019, Cox's Bazar, Bangladesh.

[26] R. N. Khushaba, S. Kodagoda, S. Lal, and G. Dissanayake, "Driver drowsiness classification using fuzzy wavelet-packet-based feature-extraction algorithm," IEEE Transactions on Biomedical Engineering, vol. 58, no. 1, pp. 121-131, 2011.

[27] World medical association declaration of Helsinki-ethical principles for medical research involving human subjects, Adopted by 64th WMA General Assembly, Fortaleza, Brazil, Special Communication: Clinical Review & Education. (2013).

[28] Acqknowledge Software (v 4.4). Available in: https://www.biopac.com/demo/acqknowledge-4-demo/

[29] N. Mammone and F. C. Morabito, "Enhanced automatic wavelet independent component analysis for electroencephalographic artifact removal," Entropy, vol. 16, no. 12, pp. 6553-6572, 2014.

[30] N. Mammone, F. L. Foresta, and F. C. Morabito, "Automatic artifact rejection from multichannel scalp EEG by wavelet ICA," IEEE Sensors Journal, vol. 12, no. 3, pp. 533-542, 2012.

[31] M. K. Wali, M. Murugappan, and B. Ahmmad, "Wavelet packet transform based driver distraction level classification using EEG," Mathematical Problems in Engineering, vol. 2013, Article ID 297587, 10 pages, 2013.

[32] A. Zell, Simulation Neuronaler Netze [Simulation of Neural Networks] (in German) (1st ed.). 1994, Addison-Wesley. ISBN 978-3-89319-554-1.

[33] C. W. Dawsin, "An artificial neural network approach to rainfall-runoff modelling," Hydrological Sciences Journal, vol. 43, no. 1, pp. 47–66, 1998.

[34] Matlab Product (Matlab 2018a). Available in: https://www.mathworks.com/company/newsroom/mathworks-announces-release-2018a-of-the-matlab-and-simulink-product-families.html

[35] M. A. Rahman, "Topoplot for B-Alert X-10 9-Channel EEG Signal," https://www.mathworks.com/matlabcentral/fileexchange/69991-topoplot-for-b-alert-x-10-9-channel-eeg-signal, MATLAB Central File Exchange. Retrieved April 2, 2019.

# A Calibrating Six-Port Compact Circuit using a New Technique Program

Traii Moubarek[1, 2], Mohannad Almanee[2], Ali Gharsallah[3]

Dept. of physics, Lab CSEHF, Faculty of sciences of Tunis, El-manar, Tunisia[1, 3]
Dept. of electronics, Riyadh College of Technology, tvtc, Saudi Arabia[1, 2]

*Abstract*—In this paper, a calibration of six-port reflectometer using a new technique program is presented. It has been shown that a calibration procedure is based on explicit method, the method that capturing the output wave forms of six-port junction and determines the complex relationship between the two waves present at the input from the value of four outputs. The number of calibrating standards and the computation effort required are the most important parameters in selecting a calibration technique. Comparison between the results obtained from the new calibration method program with measurement results show the validity of the method proposed. This calibration technique can be used in general six-port direct digital receiver.

*Keywords—Calibration technique; digital receiver; explicit method; reflectometer; S parameters*

## I. INTRODUCTION

Recently, six-port concept has been successfully applied in modern communication receivers. For this reason, calibration procedure is very important for the measurement accuracy. Many microwave applications require the determination of the reflection and or the transmission properties of a device under-test (DUT) over a specified frequency band. Various calibration procedures for six-port junction have been suggested along the years [1-2]. Among the many programs which have been proposed for the calibration of six-port reflectometer, Hassan [3], introduced, the accuracy calibration of a six-port reflectometer as an alternative to the conventional ANA.

On the other hand, Yanyang and Frigon presented a new type of six-port radio for ultra wideband containing one power detector and one variable reference load [4]. Yanyang pointed out that there exist exact relations between some system parameters and integral quantities of the detected values for a number of reference plane positions [5]. This technique is advantageous for permitting automatic procedures; nevertheless it requires a great amount of time to be accomplished and demands a large number of known standards necessary to the network training. Another technique requires Schottky diodes with homodyne detection, but this procedure would be entirely impractical to analyze a calibration procedure [6].

Eventually, in this paper, a novel calibration method for a compact six-port junction in SHF band is presented, in which a remarkable improvement is achieved. The main objective is to minimize the number of known standards required to increase the dynamic range of the six-port reflectometer and accuracy of the measurements. Following this, the simplicity of the method

provides for automatic calibration of modern applications like wireless receivers. The comparison between theoretical results and those obtained from graphic interface and experimental results is presented to demonstrate the validity of the system developed.

The analysis of the SPR, which is based on an analytical description of the system behavior, is given in Section II. In Section III-A, six-port calibration techniques have been widely described. Finally the comparison between practical results, developed program and those obtained by commercial software is presented in Section III-B to demonstrate the validity of the system developed.

## II. SYSTEM DESIGN

The six-port reflectometer provides an alternative method of implementing the ANA [7]. In common with the existing design, CST Microwave Studio is chosen due to its ability for the simulation and analysis of these items.

### A. The Proposed Six-Port Circuit

The layout and prototype of six-port junction used in calibrating system is shown in Fig. 1, experimental results have been performed by using a HP8722 network analyzer.

It was designed with RT/Duroid 6010LM having a thickness h = 0.508 mm and a relative dielectric constant $\varepsilon r$ = 3.38.

The configuration of a compact six-port junction is designed using elliptic-disc coupler with delay line us shown in Fig. 1(a). by suitable choice of larger spacing between line delay of compact elliptic coupler [8], a flat-coupling response is obtained for the coupled ports, Table I shows the dimension of miniaturized six-port prototype.

TABLE I. PARAMETERS OF THE PROPOSED SIX-PORT

| Parameter values | Values, mm | | |
|---|---|---|---|
| | *Initial values* | *Optimum values* | *values in [8]* |
| Semi-major (Sm) elliptical axis | 21.5 | 21.44 | 25.71 |
| *W* (width of the impedance steps) | 2.25 | 2.23 | 2.31 |
| Ds (length of the impedance) | 4.95 | 4.91 | 5.2 |
| Slot s (widh of delay line in compact elliptic coupler) | 1.8 | 1.85 | - |
| D (dimension of six-port junction) | 73.1 | 72.82 | 87.14 |

As can be seen from Fig. 1(b), ports 1 and 2 are connected to a local oscillator and received RF signals, respectively. The other four ports (ports 3–6) are output ports and are connected to power detectors.

Fig. 2 compares the simulated and measured scattering parameters of fabricated single layer six-port prototype.



(a)



(b)

Fig. 1. Photograph of the Proposed Six-Port Prototype, (a) Layout Design, (b) Fabricated Prototype.



(a)



(b)

Fig. 2. Scattering Parameters of Single Layer Six-Port, (a) Simulation Results, Measured Results.

To reconfigure the magnitude a characteristic, the six-port junction is adapted from 3 to 6 GHz with a resonance frequency around 4.5 GHz. This figure shows that the measured transmission coefficients of the proposed six-port circuit has some ripples in measurement results compared to the simulated ones, which may be caused by the SMA feed connector. It can be noted that the return loss is better than 15 dB over the entire operating frequency band of interest.

The proposed six-port prototype has the advantages of small volume and low cost compared to published design [8].

### B. Reflectometer System with Proposed Six-Port

The system generates a signal $\rho(n)$ in the digital domain representing the complex ratio between the two input signals of six-port junction above.

Four diode power detectors ($D_3$, $D_4$, $D_5$, and D6) enable the use of ac detection. The power detectors connected at output ports 3, 4, 5 and 6 respectively comprise silicon Schottky diodes (Hughes 47436H- 11) and RC low-pass filter.

Fig. 3 shows the measurement of the output power levels from six-port junction; also convert the AC voltage with high frequency to DC voltage with low frequency which is then passed to the calibration routine. The voltage values measured at the power detectors outputs are given by the following equation [9]:

$$V_i(t) = a_i . V_{LO}^2 + b_i . V_{RF}^2(t) + c_i . V_{RF}(t) . \cos(\Theta(t) - \emptyset_i) \quad (1)$$

Where:

- $i = 3, \ldots 6$

- $\emptyset_i = \varphi_i - \gamma_i$,

- $c_i$: Depends on $a_i$, $b_i$ and $V_{LO}$ which is supposed to be Constant.

The use of four Schottky diode detectors presents a cost-effective solution and offers a good alternative for applications where speed is important, they present linear characteristics at low power levels and high power levels respectively [10].

Four analog-to-digital converters (A ∕D) and a digital signal processor (DSP) connected to pc analyser.

Input signal $v_{rf}(t)$ and $v_{LO}(t)$ with various phases as below:

$$v_{LO}(t) = V_{LO} \cos(wt) \quad (2)$$

$$v_{rf}(t) = Re\left[V_{RF}(t)e^{j(wt+\theta(t))}\right] \quad (3)$$



Fig. 3. Power Detector at Port 3 (i=3).

Let us consider the cases where the incident and reflected signals **a** and **b** respectively has different frequencies, as in the following:

$$a = |a|. e^{j(2\Pi f_1 t + \phi_1)} \qquad (4)$$

$$b = |b|. e^{j(2\Pi f_2 t + \phi_2)} \qquad (5)$$

As illustrated in Fig. 4, the block diagram of six-port reflectometer with power detector, the manufacturer part labels provided in the figure correspond to those components used to implement the method in this work. The out coming RF signals at the reflectometer ports 3, 4, 5 and 6 have the form:

$$V_{rfi}(t) = \sqrt{a_i}. V_{LO}. \cos(wt + \emptyset_i) + \sqrt{b_i}. V_{RF}(t). \cos(wt + \theta(t) + \theta_i) \qquad (6)$$

Where:

- $i = 3,..,6$,

- $\varphi_i$ is the phase of $v_{LO}(t)$ at port i with respect to port 1,

- $\gamma_i$ is the phase of $v_{rf}(t)$ at port i relative to port 2,

- $a_i, b_i$ depend on the circuit characteristics.

The term $V_{RF}(t)\cos(\theta(t) - \emptyset_i)$ is a projection of $v_{rf}(t)$. After a trigonometric calculation, (6) becomes:

$$\hat{U}_i(t) = b_i V_{RF}^2(t) + c_i \cos\emptyset_i U(t) + c_i \sin\emptyset_i V(t) \qquad (7)$$

$$\hat{U}_i(t) = v_i(t) - a_i V_{LO}^2; \quad i=3,..,6 \qquad (8)$$

$$U(t) = V_{RF}(t)\cos\theta(t) \qquad (9)$$

$$V(t) = V_{RF}(t). \sin\theta(t) \qquad (10)$$

Calibration method of the six-port system is based in previous calculation; this calibrating program consists in finding a signal $\rho(n)$ in the complex domain.



Fig. 4. Block Diagram of Six-Port Reflectometer with Power Detector.

## III. CALIBRATING TECHNIQUE PROPOSED

Six-port calibration techniques have been widely reported in the literature [11], [12]. All of the methodologies that begin with the relation between the power measurements and the reflection coefficient are suited only for an intrinsically narrow-band system.

To over this drawback, we study in this paper the six-port reflectometer with power detectors, a measurement port and four sidearm ports to which power detectors are connected. It measures the reflection coefficient of a device under test (DUT) in terms of magnitude and phase using four sidearm power readings and eleven system parameters composed of three real quantities and four complex quantities [13].

Thus, the key problem in the practical use of the six-port reflectometer is how to determine the eleven frequency dependent system parameters, with fewer calibration standards and with less computational effort. For digital processing of out coming six-port reflectometer, it is essential to find a mathematical modeling of output quantities.

Below the model used to express the detected voltages $V_1$ to $V_4$ depending on the module $|\Gamma|$ and phase $\Phi$ of the reflection coefficient $S_{11}$ of the DUT [14]:

$$V_i = a_{i0} + a_{i1}. |\Gamma|^2 + a_{i2}. |\Gamma|. \sin(\emptyset) + a_{i3}. |\Gamma|. \cos(\emptyset) + a_{i4}. |\Gamma|^2. \sin(2\emptyset) + a_{i5}. |\Gamma|^2. \cos(2\emptyset) \qquad (11)$$

i=1…6

By using the real and imaginary parts of $\Gamma$, it is possible to rephrase the previous model by:

$$V_i = b_{i0} + b_{i1}. I + b_{i2}. Q + b_{i3}. I^2 + b_{i4}. Q^2 + b_{i5}. I. Q + \cdots + b_{in} \qquad (12)$$

Where:

$I = Re(\Gamma), Q = Im(\Gamma)$

$b_{ij}$: Calibration constant.

$n$: Index number of minimum standards.

We present below three different levels retainer model:

$$V_i = b_{i0} + b_{i1} + b_{i2}Q + b_{i3}(I^2 + Q^2) \qquad (13)$$

=> First order for $i = 1, ..., 4$

$$V_i = b_{i0} + b_{i1}I + b_{i2}Q + b_{i3}I^2 + b_{i4}Q^2 + b_{i5}IQ \qquad (14)$$

=> Second order for $i = 1, ..., 4$

$$V_i = b_{i0} + b_{i1}. I + b_{i2}. Q + b_{i3}. I^2 + b_{i4}. Q^2 + b_{i5}. IQ + b_{i6}. (I^2 - 3IQ^2) + b_{i7}. (Q^2 - 3I^2Q) \qquad (15)$$

=> Third order for $i = 1, ..., 4$

The $b_{ij}$ terms representing calibration constants depending only on the terms $a_{ij}$. The advantage of using such a development (13), (14) and (15) is to assign the same variation ranges to two quantities of interest (I $\in$ [-1, 1] and Q$\in$ [-1, 1]) that we are trying to determine. Indeed, given the form of expressions, it appears that resolution requires the implementation of a digital resolution in two dimensions.

These equations are used to express the magnitudes of interest (I and Q) as a function of the measured voltages on each of the detectors ($V_3,..,V_6$) [15]. I and Q are related by the reflectometer calibration constants.

This calibration consists in solving the equations of the model by considering the known loads and measuring the tensions of these standard loads. Inputs correspond to a matrix of the measured voltages $V_M$ while the resulting matrix contains the terms I and Q. The calibration constants are grouped in two known matrices $B_0$ and $B_M$.

Note that the parameters that depend on the order of the model are: The number of standard loads required for calibration {(4 loads: order 1); (6 loads: order 2); (8 loads: order 3)}. The dimension of the calibration matrix {($3\times4$, order 1); ($5\times4$, order 2); ($7\times4$, order 3)}. In order to optimize the accuracy/ computation time, we reduce the number of used standard and minimize the number of test, we will adopt the second model. i denotes the ith outputs.

$$V_i = b_{i0} + b_{i1}.I + b_{i2}.Q + b_{i3}.I^2 + b_{i4}.Q^2 + b_{i5}.IQ \qquad (16)$$

For a load with reflection coefficient Γ, the four equations form a linear system of five variables. The type $(I, Q, I^2, Q^2, IQ)$ with:

$$V_M = B_M.\Gamma_M + B_0 \qquad (17)$$

$$V_M = \begin{bmatrix} V_1 \\ V_2 \\ V_3 \\ V_4 \end{bmatrix} \qquad (18)$$

$$B_M = \begin{bmatrix} b_{11} & b_{12} & b_{13} & b_{14} & b_{15} \\ b_{21} & b_{22} & b_{23} & b_{24} & b_{25} \\ b_{31} & b_{32} & b_{33} & b_{34} & b_{35} \\ b_{41} & b_{42} & b_{43} & b_{44} & b_{45} \end{bmatrix} \qquad (19)$$

$$\Gamma_M = \begin{bmatrix} I \\ Q \\ I^2 \\ Q^2 \\ I.Q \end{bmatrix} \qquad (20)$$

$$B_0 = \begin{bmatrix} b_{11} \\ b_{21} \\ b_{31} \\ b_{41} \end{bmatrix} \qquad (21)$$

$V_M$: Matrix of the detected voltages;

$\Gamma_M$: Matrix containing the appropriate quantities I and Q;

$B_M$: Matrix of the calibration coefficients of the I and Q magnitudes;

$B_0$: Matrix of the DC components of the voltages $V_i$,

$i = 1, ..., 4$.

Knowledge of $B_M$ matrices, $B_0$ is associated with the measurement of the four voltages $V_1$ to $V_4$ allows then the determination of I and Q by reversing the previous matrix system:

$$\Gamma_M = B_M^{-1}(V_M - B_0) \qquad (22)$$

The solution of the above matrix equation leads to the calculation of $B_M$ and $B_0$ corresponds to the calibration of the six-port junction.

### A. Algorithm and Interface

The algorithm is based on the mathematical development of explicit calibration above. The inputs in this program are determined in the previous paragraph. We can summarize the different steps in the calculation by the following chart:

The S parameters of six-port calibrated can be regarded as being the same at the frequency difference $f_2-f_1$ if this difference is very small. Let us consider the cases where the incidents wave "**a**" and the reflected wave "**b**" are different in frequency as in the equations (3) and (4) above. The equivalent reflection coefficient becomes:

$$\Gamma = \frac{b}{a} \, e^{j[2\pi(f_2-f_1)-(\emptyset_2-\emptyset_1)]} \tag{23}$$

The reflection coefficient $\Gamma$ becomes a time-dependent vector whose amplitude is invariant and the phase is around at a constant angular speed $2\pi\Delta f, where \; \Delta f = f_2 - f_1$

The voltage at a detector port of the six-port is a vector summation of portions of **a** and **b** presented at the port, the output voltage waveforms of the power detectors are:

$$V_{out,i} = \frac{1}{2}(|a_i|^2 + |b_i|^2) + |a_i|.|b_i|.\cos(2\pi\Delta ft + \Delta\emptyset_i) \tag{24}$$

Where $a_i$ and $b_i$ are waves corresponding to **a** and **b** at port i, $i = 3, ... , 6$

The four output voltage waveforms ($v_1$, $v_2$, $v_3$ and $v_4$) of a six-port reflectometer are shown in Fig. 5:

We note that a period of the waveform corresponds to a whole circle of $\Gamma$ rotating in the complex plane. We select the samples at an equal amplitude space in the whole voltage swing in this channel; these samples will present a group of equally spaced terminations on a circle in the $\Gamma$ plane.

To ensure the communication of the user with the program, we have made use of a graphical interface.

The final circuit of six-port reflectometer is designed and simulated at Fig. 6 to validate results introduced in Fig. 5. The Graphical interface based at the above program, offers us the possibility to calculate the reflection coefficient of a device under test. The DUT is connected to the input of six-port reflectometer.



(a)



(b)

Fig. 5.   (a): V1, V4 and (b): V2, V3 Output Wave forms of a Six-Port Reflectometer, Simulation Results.



Fig. 6.   Graphical Interface of Six-Port Recflectometer.

### B. Validation by Measuring Test Loads

To validate calibration technique performance of miniaturized six-port junction, circuit is realized with screen printing technique; this process is used due to its ability to print a flexible substrate. AC voltage is the output of six-port junction and the input of four power detectors circuit to convert it to dc voltage; we need to connect at the output of power detectors LM324N to increase the output voltage.

LM324N is a 14pin IC consisting of four independent operational amplifiers compensated in a single package. Signal at output of op-amps implement microcontroller mikroC PRO for PIC to validate calibration technique.

The final hybrid prototype and photograph experimental setup system shown in Fig. 7.

Results of reflection coefficients performed is comparing with the program above those found by experiment prototype. The load under test is composed of lumped elements in series (Inductance, Resistance and Capacity), by varying the amplitude of either the input signal, samples of Γ′s well distributed over the whole Smith chart can be obtained.

In order to test the validity of the six-port calibration, a calculation of relative and average errors are shown, results are presented in Table II.

Table above lists the calculation errors obtained for different samples. It can be seen that the relative error for the real parts of the reflection coefficients is about 6% except for very low reflection loads. In another hand it is noticed that the relative error for the imaginary parts of the reflection coefficients is about 4.8%.

In most case, the difference of errors between the real and imaginary parts is in order of 0.01. Note that the readings are taken from a real time continuous display of reflection coefficients. Analysis results are confirmed schematically by the calibration technique studied in Smith chart below.

Fig. 8 shows a maximum of agreement between experimental results, simulation and algorithm taking samples of plane calibration in terms of reflection coefficient. It's very clear for large impedances and inductive loads (positive-imaginary part). Note that the error increases inversely with the impedance, this is due to the relative error of the load impedance.

This error is higher for the low impedances than for high impedance of the load. But errors have several origins, particularly for millimeter frequencies, the achievement of a reflectometer generates imperfections that away from the perfect model.

As well as the adaptation errors, isolation and directivity of the source and the load. It will be impossible to get such perfect 50Ω terminated transmission lines or lossless, in addition to the accuracy of the model of the output voltages.



Fig. 7. Displaying Meter Results with Graphical Interface.

TABLE II. IMPEDANCE AND REFLECTION COEFFICIENT LOADS

| R(S11) | I(S11) | R(OHM) | Z(H OU F) | R | Q | DELTA (R) | DELTA(Q) | DR² | DQ² |
|--------|--------|--------|-----------|---|---|-----------|----------|-----|-----|
| 0 | 0 | CA | CA | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | CO | CO | 0,93 | 0,04 | 0,07 | -0,04 | 0,0049 | 0,0016 |
| -1 | 0 | CC | CC | -1 | 0,07 | -0,05 | -0,07 | 0,0025 | 0,0049 |
| 0,5 | 0,5 | 50 | 8E-09 | 0,55 | 0,52 | -0,05 | -0,02 | 0,0025 | 0,0004 |
| 0,4 | 0,7 | 20,588 | 6,6E-09 | 0,48 | 0,72 | -0,08 | -0,02 | 0,0064 | 0,0004 |
| 0,3 | 0,3 | 70,68 | 4,1E-09 | 0,35 | 0,31 | -0,05 | -0,01 | 0,0025 | 0,0001 |
| 0 | 1 | 0 | 4E-09 | -0 | 0,98 | 0,022 | 0,02 | 0,0005 | 0,0004 |
| -0,471 | 0,883 | 0 | 2,4E-09 | -0,3 | 0,98 | -0,131 | -0,097 | 0,0172 | 0,00941 |
| Addition | | | | | | 6,841 | 10,339 | 0,1974 | 0,24743 |
| Relative error | | | | | | 0,065 | 0,04811 | Relative error | 0,065 |
| Average error | | | | | | 0,0202 | 0,02261 | Average error | 0,0202 |

Fig. 8.    Smith Chart of the Reflection Coefficients with Algorithm, Simulation and Experiment Prototype.

## IV.    CONCLUSION

A novel explicit calibration for six-port reflectometer incorporating new four-port correlators has been proposed and shown. A linearization procedure and an AC detection technique are used to improve the measurement accuracy. This solution makes a good compromise between the number of calibration standards, the computational cost and the accuracy.

Regarding calibration considerations, a further improvement in the technique will be used to achieve measurements in the 1–10 GHz frequency range. Therefore the proposed technique program is suitable for wireless communication system and can successfully replace the classical calibration methods.

## ACKNOWLEDGMENT

## REFERENCES

[1] K. Haddadi, Tuami Lasri, "Formulation for Complete and Accurate Calibration of Six-Port Reflectometer ", IEEE Trans, Microwave Theory Tech, vol. 60, no 3, pp. 574- 581, March 2012.

[2] Kamil Staszek, "Six-Port Calibration Utilizing Matched Load and Unknown Calibration Loads," IEEE Trans, Microwave Theory Tech, vol. 66, no 10, pp. 4617-4626, oct 2018.

[3] Abul Hasan, Mohamed Helaoui, "Formulation for Complete and Accurate Calibration of Six-Port Reflectometer," IEEE Trans, Microwave Theory Tech, vol. 60, no 12, pp. 574-581, 2012.

[4] Abdullah O. Aolopade, Mohamed helaoui, "High performance homodyne six-port receiver using memory polynomial calibration," 2014 IEEE 27 th Canadian conference on Electrical and Computer Engineering CCECE 2014.

[5] Kamil Staszek "Six-port Calibrating Utilizing Matched Load and Unknown Calibration Loads," IEEE Transactions on Micrwave Theory and Techniques vol 66, no 10,  pp. 41–44, 2018.

[6] Traii Moubarek, Ali Gharsallah, "A Six-Port Reflectometer Calibration Using Wilkinson Power Divider," American Journal of Engineering and Applied Sciences, vol. 10, pp.38-44, Jun 2016.

[7] K. Haddadi, D. Glay, and T. Lasri, "Homodyne dual six-port network analyzer and associated calibration technique for millimetre wave measurements," Int. Symp. Circuits Syst., Island of Kos, Greece, IEEE, pp. 1–4, May 2006.

[8] Traii Moubarek, Mourad Nedil, Ali Gharsallah, Tayeb A. Denidni, "A New wideband Six-port Junction using single layer Technology,"International Journal of Information Sciences and Computer Engineering, vol. 2, NO. 2,  pp.14-48, Jun 2011.

[9] Fadhel M. Ghannouchi, Abbas Mohammadi, "The six-port Technique with Microwave and Wireless Applications," Aretch House Microwave Library, 2009.

[10] S.Wafi, "Study of a system based on microwave junction six-port with associated calibration technique, " Master Degree, El Manar University, Tunisia, 2011.

[11] C. Potter and A. Bullock, "Non-linearity correction of microwave diode detectors using a repeatable attenuation step," Microwave Journal, vol. 36, pp. 272, 274 and 277–279, May 1993.

[12] Kamil Staszek, Sarah Linz, Fabian Lurz, Sebastian Mann, and Robert Weeigel, Alexander Koelpin, "Improved calibration procedure for six-port based precise displacement measurements," 2016 IEEE Topical Conference on Wireless Sensors and Sensor Networks (WiSNet), 2016.

[13] Jigisha Das, Srimoyi Roy, Rijubrata    Pal, Mrinal Kanti Mandal, "Rive points method of calibration for six-port receivers," 2018 3 rd International Conference on Microwave and Photonics (ICMAP), 2018.

[14] K.Haddadi, C. Loyez, L.Clavier, D.Pomorski, S.Lallemand, "Six-port reflectometer in WR15 metallic waveguide for free-space sensing applications," 2018 IEEE Topical Conference on Wireless Sensors and Sensor Networks (WiSNet), pp. 80-83, March 2018.

[15] Vladimir Bilik, Six-port Measurement  Technique : Principles, Impact, applications, Slovak University of Technology.

# Hybrid Concatenated LDPC Codes with LTE Modulation Schemes

Mohanad Alfiras[1]

Communication and networks
Engineering, Gulf University,
Manama, Kingdom of Bahrain

Wael A. H. Hadi[2], Amjad Ali Jassim[3]

Communication Engineering Department
Engineering college, Al-Technology University
Baghdad, Iraq

*Abstract*—In a communication system, the LDPC code is considered as a good performance error correcting code which reaches near Shannon limit. In this paper a hybrid LDPC code is proposed, the hybrid term here refers to the serial concatenation of parallel LDPC codes group and a single serial LDPC code. The outer two parallel LDPC codes encoder represents outer encoder where the single LDPC encoder represents the inner encoder. This study also emphases on the performance of a hybrid coding system in consideration with three modulation schemes. The modulation schemes include quadrature phase shift keying (QPSK) and two types of quadrature amplitude modulation; 16-QAM and 64-QAM. These modulation schemes are selected due to their importance in modern communication applications, such as long term evolution (LTE); such schemes are the standard modulation schemes used with LTE system. This study investigates different LDPC code rates such as 1/2 and 1/3 and simulates the AWGN communication channel using MATLAB. The simulation results show improvement in bit error rate (BER) when using 1/3 LDPC code rate in the designed system rather than 1/2, but it also increases the system complexity. In the end, all simulation results, the comparison between different cases of LDPC code rates and system performance are summarized.

*Keywords—Coding; modulation; Hybrid; concatenation; low density parity check*

## I. INTRODUCTION

G. David Forny introduced the concept of concatenated codes in 1965 as a method to improve system code performance [1]. Earlier to this, Gallager in his Ph.D. dissertation at M.I.T invented LDPC codes in 1960 [2]. This paper focuses on the concept that uses the LDPC codes concatenation as an error correction code to improve communication system performance against errors raised during signal transmission through the noisy channel. The code concatenation improves the performance of the error correction code as a group against transmission errors. The simplest form of code concatenation is observed in serial concatenation between two error correction codes. The codes that are involved in serial concatenation can be of different or same type. For example, serial concatenation between two convolutional encoders [3] or between two different error correction codes, such as, read Solomon code and convolutional encoder [4] and between read Solomon code and LDPC code [5]. The error correction code concatenation also has another form which consists of parallel concatenation between two identical error correction codes [6]. For example, the parallel concatenation between two convolutional encoders

forms a well-known Turbo code. The concatenation between two parallel convolutional encoders offers the significant BER Turbo code performance, but with increased decoder complexity. The decoder uses complicated algorithms to calculate decoding frame estimation such as SOVA and BCJR algorithms [7]. LDPC codes can be used with the same concept of code concatenation. Therefore, researchers focus on using serial concatenation of LDPC codes of the same code type and rate [8]. In order to form a serially concatenated codes, consider concatenation between LDPC codes of different types as compatible pairs [9]. In order to achieve that, the inner LDPC encoder takes code rates which are fitted in data rate with outer LDPC encoder. This strategy takes benefit of increasing only the inner LDPC encoder size and reduce the system complexity as compared when using two large size LDPC codes. Also, parallel concatenation is applicable for LDPC codes using two identical LDPC codes with a simple modification in the receiver to avoid increasing system decoder complexity [10]. This modification includes taking the sum of the two LDPC decoders that are based on bit flipping algorithm [11]. In many communication applications such as deep space communication systems, there is a need for accurate BER performance against raised communication errors rather than the system complexity [12]. Therefore, the system complexity could be acceptable in some communication applications where the actual goal of the system is to achieve a good BER performance [13]. This work focus on using short length irregular LDPC codes connected in Hybrid form, the term Hybrid concatenation used here refers to two types of concatenated codes, parallel concatenation between two identical LDPC encoders and serially concatenated with inner LDPC encoder. The concatenation strategy should be carefully designed to get better system performance and low system complexity as much as possible. So, different code rates and different modulation schemes are investigated by comparing different system designs.

## II. PROPOSED SYSTEM TRANSMITTER

The proposed hybrid system transmitter consists of three main stages. First, represented by identical irregular LDPC codes connected in parallel to construct outer encoder. Second, the produced code words from two LDPC encoder are multiplexed to prepare input for the next stage of inner encoder LDPC with length equal to twice individual outer LDPC code. The inner code input frame length will be two times the output of a single LDPC code used to construct the outer parallel

group. Third, the modulation scheme such as QPSK, 16-QAM or 64-QAM considered in modulation stage [14, 15] as LTE standard modulation schemes. The system transmitter is shown in Fig. 1.

In order to process the input data, the Hybrid system transmitter starts by the LDPC1 and LDPC2 encoders which are designed to be identical. Each code of data rate equal to $1/n_1$. There is an interleaver between these parallel LDPC1 and LDPC2 encoders which is denoted by $\pi_1$ in Fig. 1. The interleaver rearrange the order of input data frame to construct a new code word which differs from the code word of LDPC1 encoder. The interleaver could be random or another type. Here, the main usage of interleaver is to overcome the effects of burst errors, if found in the received codeword and change it to suppurated errors that could be handled and corrected by LDPC decoder. Each LDPC1 and LDPC2 encoder of code rate

$$k/n_1 = k/n_2 = k/n \qquad (1)$$

Where 'k' represents the length of the input data frame, and the 'n' represents the length of the produced code word. Therefore, after the parallel encoding process and multiplexer, it produces a codeword of length

$$1/(2 \times n_1) = 1/(2 \times n_2) = 1/(2 \times n) \qquad (2)$$

Inner LDPC encoder is designed to be of different length, not the same as LDPC codes which are used with an outer parallel group. The length of the input data which inner LDPC takes is the same as the length of the produced codeword by the outer parallel group. Then the inner LDPC code rate $1/n_3$, in terms of input is

$$R = (2 \times n_1)/(2 \times n_1 \times n_3) \qquad (3)$$

The symbol $\pi_2$ in system stands for random interleaver used to enhance system performance against burst errors. In general, for LDPC codes, the larger length LDPC encoder gives the best performance as compared with shorter ones. However, larger LDPC code length means more complex LDPC decoder in system receiver. The inner LDPC3 will be an effective inner encoder, without using a larger encoder.



Fig. 1. Hybrid LDPC Code System Transmitter.

### III. PROPOSED SYSTEM RECEIVER

The hybrid system receiver starts with the demodulation process which is the same as the modulation scheme used in the system transmitter. The decoding process starts in reciprocal order of the transmitter encoder process. It starts by decoding LDPC3 which represent inner encoder of the transmitter. This decoder uses the LLR (Log Likelihood Ratio algorithm) [16]. The LDPC decoder gives estimation output which is denoted by $\hat{E}_3$. The decoding algorithm LLR is described as follow [16]:

The input to the LDPC decoder is the log-likelihood ratio (LLR), $L(c_i)$, which is defined by

$$L(c_i) = log\left(\frac{\Pr(c_i=0|\, channel\ output\ for\ ci)}{\Pr(c_i=1|channel\ output\ for\ ci)}\right) \qquad (4)$$

Where $c_i$ is an ith bit of the transmitted codeword $c$. There are three key variables in the algorithm: $L(r_{ji})$, $L(q_{ij})$, and $L(Q_i)$. $L(q_{ij})$ is initialized as $L(q_{ij}) = L(c_i)$. For each iteration, update $L(r_{ji})$, $L(q_{ij})$, and $L(Q_i)$ using the following set of equations [16]:

$$L(r_{ji}) = 2\ \text{atanh}(\prod_{i' \in V_{j\backslash i}} \tanh(\tfrac{1}{2}\ L(q_{i'j})))$$

$$L(q_{ij}) = L(c_i) + \sum_{j' \in c_i\backslash j} L(r_{j'i})$$

$$L(Q_i) = L(c_i) + \sum_{j' \in c_i} L(r_{j'i}) \qquad (5)$$

$\hat{E}3$ represents the estimated output from LDPC$_3$, which is passed to de-multiplexer to redirect $\hat{E}_3$ into two groups which are inputs for LDPC decoder one and LDPC decoder two, respectively. The LDPC decoders of one and two use the same described decoding algorithm LLR. Such a process at the end produce another two estimations denoted by $\hat{E}_1$ and $\hat{E}_2$ refers to each decoder of the parallel group. The two estimations then summed before the decision. The decision represents the final receiver stage to produce received data. Fig. 2 shows the proposed system receiver.

Where the symbol $\pi^{-1}$ refers to random de-interleaver.



Fig. 2. Proposed Hybrid System Receiver.

### IV. SIMULATION PARAMETERS

The simulation includes the generation of two sets of irregular LDPC codes. First of rate 1/2 and second of rate 1/3. Table I and Table II summarize the description where the code is described by $C_b$ (N, K), N codeword length, and K input data length.

The simulation takes in consideration of multiple hybrid system designs; it discusses the increase in the length of generated irregular LDPC codes and also designs the hybrid system with two LDPC code rates, first represented by rate 1/2 and the second of rate 1/3. It gives us more sense about increasing LDPC code rate and length and its effect in system BER performance.

TABLE I. LDPC CODES RATE 1/2

| Design | Outer encoder parallel group | | Inner encoder |
|--------|--------|--------|--------|
| | LDPC1 | LDPC2 | LDPC3 |
| Case 1 | $C_b$ (48, 24) | $C_b$ (48, 24) | $C_b$ (192, 96) |
| Case 2 | $C_b$ (96, 48) | $C_b$ (96, 48) | $C_b$ (384, 192) |
| Case 3 | $C_b$ (144, 72) | $C_b$ (144, 72) | $C_b$ (576, 288) |
| Case 4 | $C_b$ (192, 96) | $C_b$ (192, 96) | $C_b$ (768, 384) |
| Case 5 | $C_b$ (240, 120) | $C_b$ (240, 120) | $C_b$ (960, 480) |

TABLE II. LDPC CODES RATE 1/3

| Design | Outer encoder parallel group | | Inner encoder |
|--------|--------|--------|--------|
| | LDPC1 | LDPC2 | LDPC3 |
| Case 1 | $C_b$ (72, 24) | $C_b$ (72, 24) | $C_b$ (432, 144) |
| Case 2 | $C_b$ (144, 48) | $C_b$ (144, 48) | $C_b$ (864, 288) |
| Case 3 | $C_b$ (216, 72) | $C_b$ (216, 72) | $C_b$ (1296, 432) |
| Case 4 | $C_b$ (288, 96) | $C_b$ (288, 96) | $C_b$ (1728, 576) |
| Case 5 | $C_b$ (360, 120) | $C_b$ (360, 120) | $C_b$ (2160, 720) |

## V. SIMULATION RESULTS

The proposed system simulation includes the cases and their corresponding generated irregular LDPC codes. It splits the system performance as BER vs. SNR into two groups depending on the LDPC code rate. In each case, a modulation scheme is selected from three types of QPSK which are used for low data rate while its symbol consists of two bits; 64-QAM for high data rate with good quality SNR where each symbol consisted of 6 bits and 16-QAM with 4 bits per symbol. These modulation schemes are used as a standard with LTE application [4]. Fig. 3 to 8 shows simulation results, respectively.

The simulation result values (BER) compares different system parameters listed in Table III and Table IV, respectively.



Fig. 3. LDPC Code Rate 1/2, QPSK.



Fig. 4. LDPC Code Rate 1/2, 16-QAM.



Fig. 5. LDPC Code Rate 1/2, 64-QAM.



Fig. 6. LDPC Code Rate 1/3, QPSK.

Fig. 7.   LDPC Code Rate 1/3, 16-QAM.



Fig. 8.   LDPC Code Rate 1/3, 64-QAM.

TABLE III.    SHOWS SIMULATION RESULTS IN COMPARISON TO HYBRID SYSTEM LDPC RATE 1/2

| Modulation type | Outer Parallel Two LDPC Codes | Inner LDPC code | SNR | BER |
|---|---|---|---|---|
| QPSK | Cb(48, 24) | Cb(192, 96) | 4 | $1.1 \times 10^{-5}$ |
| | Cb(96, 48) | Cb(384, 192) | 2.5 | $2.2999 \times 10^{-5}$ |
| | Cb(144, 72) | Cb(576, 288) | 2.5 | $9.9999 \times 10^{-7}$ |
| | Cb(192, 96) | Cb(768, 384) | 2 | $1.9999 \times 10^{-6}$ |
| | Cb(240, 120) | Cb(960, 480) | 1.5 | $1.2999 \times 10^{-5}$ |
| 16-QAM | Cb(48, 24) | Cb(192, 96) | 10 | $2 \times 10^{-5}$ |
| | Cb(96, 48) | Cb(384, 192) | 9 | $1.9999 \times 10^{-6}$ |
| | Cb(144, 72) | Cb(576, 288) | 7 | 0.000149 |
| | Cb(192, 96) | Cb(768, 384) | 7 | $3.1999 \times 10^{-5}$ |
| | Cb(240, 120) | Cb(960, 480) | 7 | $2.9998 \times 10^{-6}$ |
| 64-QAM | Cb(48, 24) | Cb(192, 96) | 15 | $2.3 \times 10^{-5}$ |
| | Cb(96, 48) | Cb(384, 192) | 13 | $3.0999 \times 10^{-5}$ |
| | Cb(144, 72) | Cb(576, 288) | 12 | $1.2 \times 10^{-5}$ |
| | Cb(192, 96) | Cb(768, 384) | 12 | $9.9997 \times 10^{-7}$ |
| | Cb(240, 120) | Cb(960, 480) | 11 | 0.00014399 |

TABLE IV.    SHOWS SIMULATION RESULTS IN COMPARISON TO HYBRID SYSTEM LDPC RATE 1/3

| Modulation Type | Outer Parallel Two LDPC Codes | Inner LDPC Codes | SNR | BER |
|---|---|---|---|---|
| QPSK | Cb(72, 24) | Cb(432, 144) | 1 | $6 \times 10^{-6}$ |
| | Cb(144, 48) | Cb(864, 288) | 0 | $1.1 \times 10^{-5}$ |
| | Cb(216, 72) | Cb(1296, 432) | -0.4 | $6 \times 10^{-6}$ |
| | Cb(288, 96) | Cb(1728, 576) | -0.6 | $6.9998 \times 10^{-6}$ |
| | Cb(360, 120) | Cb(2160, 720) | -0.8 | $5.9995 \times 10^{-6}$ |
| 16-QAM | Cb(72, 24) | Cb(432, 144) | 6 | $2.1 \times 10^{-5}$ |
| | Cb(144, 48) | Cb(864, 288) | 5 | $1.4 \times 10^{-5}$ |
| | Cb(216, 72) | Cb(1296, 432) | 4 | $5.2 \times 10^{-5}$ |
| | Cb(288, 96) | Cb(1728, 576) | 4 | $8.9997 \times 10^{-6}$ |
| | Cb(360, 120) | Cb(2160, 720) | 4 | $9.9992 \times 10^{-7}$ |
| 64-QAM | Cb(72, 24) | Cb(432, 144) | 10 | $6 \times 10^{-6}$ |
| | Cb(144, 48) | Cb(864, 288) | 9 | $9.9997 \times 10^{-7}$ |
| | Cb(216, 72) | Cb(1296, 432) | 8 | $1.1 \times 10^{-5}$ |
| | Cb(288, 96) | Cb(1728, 576) | 7 | $5.6998 \times 10^{-5}$ |
| | Cb(360, 120) | Cb(2160, 720) | 7 | $1.1999 \times 10^{-5}$ |

## VI. Conclusions

The proposed hybrid system consists of serial concatenation between parallel and serial LDPC codes. The work discusses the different generated irregular LDPC codes. From the simulation, there are two main results; the first result represents the increasing length of LDPC code which shows enhancement in system BER performance as shown in Fig. 3 to 8. The second result is obtained by using LDPC codes of code rate 1/3 instead of rate 1/2 which shows more improvement in system BER performance. This improvement is realized with increased system complexity. However, the designed system is a compromise between cost and performance. Hence, 16-QAM at 7 dB reaches $10^{-5}$ BER and QPSK -2 dB which shows negative performance. Where 16-QAM reaches $10^{-6}$ at 4 dB SNR value. The system looks complicated, but it should be noted that the design uses a short length of irregular LDPC codes as maximum $C_b$ (2160, 720) for rate 1/3 LDPC as an inner encoder. The choice between different hybrid systems introduced in this work comes with two considerations, i.e., the performance and system complexity. The proposed hybrid system could be achieved in a practical application using FPGA. Since the LDPC codes show flexibility in the implementation using such technology. The LDPC decoder algorithms provide simple decoding estimation such as Bit Flipping decoding algorithm. Such consideration can reduce the hybrid system complexity.

### References

[1] G David Forny JR, "Concatenated Codes", Massachusetts Institute of technology research laboratory of electronics Cambridge, technical report 440, December 1965.

[2] Robert G. Gallager, "Low-Density Parity-Check Codes", Cambridge Mass, July 1963.

[3] Deepak Mishra, T.V.S Ram etc, "Concatenated Convolutional Codes for Deep Space Mission", International Journal of Information and Communication Technology Research, Volume 2 No. 6, June 2012.

[4] Kattaswamy Mergu, "Performance Analysis of Reed-Solomon Codes Concatenated with Convolutional Codes over AWGN Channel", Wolaita Sodo University, January 2016.

[5] Z. Shi, C. Fu, and S. Li, "Serial Concatenation and Joint Iterative Decoding of LDPC Codes and Reed-Solomon Codes", Open Research Fund of National Mobile Communication Research Laboratory, Southeast University. 2005.

[6] Satoshi Tajima, Takumi Takahashi, etc. "Iterative Decoding Based on Concatenated Belief Propagation for CRC-Aided Polar Codes", Proceedings, APSIPA Annual Summit and Conference 2018.

[7] M.Srinivasa Rao, G.Vijaya Kumar, and P.Rajesh Kumar, "Optimized BER Performance of Asymmetric Turbo Codes over AWGN Channel", International Journal of Computer Applications(0975 – 8887).Vol. 81 – No.12, November 2013.

[8] Latifa Mostari and Abdelmalik Taleb-Ahmed. "Non-Binary Serial Turbo LDPC Codes Combined with High Order Constellations", Pertanika J. Sci. & Technol. 27 (1): 33 – 47, 2019.

[9] Amjad Ali Jassim, Wael A. Hadi, and Muhanned Alfiras, "Serially Concatenated Low-density Parity Check Codes as Compatible Pairs", International Journal of Engineering & technology, Vol.7, No 4.15 (2018)

[10] Nguyen Tung Hung, Nguyen Van Duan, Do Quoc Trinh, "Parallel and Serial LDPC Decoders for Wifi and Wimax Receivers", Chuyên san Công Singhê thông tin và Truyền thông - Sè (4-2015).

[11] Mohanad Alfiras1,∗, Wael A. H. Hadi2 and Amjad Ali Jassim3, "Parallel Concatenation of LDPC Codes with LTE Modulation Schemes", Applied Mathematics & Information Sciences An International Journal, Appl. Math. Inf. Sci. 12, No. 6, 1165-1176 (2018).

[12] Carlo Condo, "Concatenated Turbo/LDPC Codes for Deep Space Communications: Performance and Implementation", The Fifth International Conference on Advances in Satellite and Space Communications, SPACOMM, January 2013.

[13] Shwetha n, Nagaraj p, J V Narasimham," Design & Implementation of Concatenated Turbo/LDPC Codes for Deep Space Communications", International Journal of Industrial Electronics and Electrical Engineering, ISSN: 2347-6982, Vol.3, Issue-9, pp 6-9, Sept.-2015.

[14] Dr Houman Zarrinkoub, "Understanding LTE with MATLAB from Mathematical Modeling to Simulation and Prototyping", John Wiley & Sons, Ltd, 2014.

[15] Christopher Cox, "An Introduction to LTE LTE, LTE-Advanced, SAE, VoLTE and 4G Mobile Communications", second edition, John Wiley & Sons, Ltd, 2014.

[16] Kun Cheng, etc. "Multi-Code-Rate Correction Technique with IR-QC-LDPC: An application to QKD", 21st IEEE real time conference - colonial Williamsburg, 13 mar 2019.

# Cloud Computing Adoption in Small and Medium-Sized Enterprises (SMEs) of Asia and Africa

## A Cross-Continent Overview of Advantages and Challenges

Babur Hayat Malik[1], Jazba Asad[2], Sabila Kousar[3], Faiza Nawaz[4], Zainab[5]
Farania Hayder[6], Sehresh Bibi[7], Amina Yousaf[8], Ali Raza[9]
Department of Computer Science and Information Technology
University of Lahore, Chenab Campus Gujrat, Pakistan

*Abstract*—**Cloud computing is a rapidly emerging technology over the last few years that has abolished the burden of purchasing heavy hardware and licensed software. Cloud computing has been advantageous to Small and Medium-sized Enterprises (SMEs), but still numerous SMEs have not adopted cloud computing to delve into its appealing benefits. Asia and Africa vary notably regarding their innovative capability. Asia has been competent to advance and sustain world leadership in technological innovations whereas Africa has not developed significantly in these terms. A seldom comparative study has been implemented on the reasons for the innovation gap between these two continents. This article examines and compares the cloud computing adoption from a Geo-regional framework; Asia and Africa. A comparative study is used to organize the findings from China in Asia, and Nigeria in Africa. The article identifies the probable benefits, usage of cloud computing and level of cloud computing adoption amid SMEs in Nigeria and China. The paper explores the margin that subsists amongst the level of cloud computing adoption in SMEs of these two countries and specifies challenges particular to each country intercepting the complete cloud computing adoption and proposes solutions for Nigerian SMEs to beat these challenges. Furthermore, the article contributes proof-supported intrusion for cloud service providers, the government and the capitalism to enhance the cloud computing adoption amid SMEs to eventually determine the enterprises for the probable financial advantage.**

*Keywords*—*Cloud computing; adoption; Asia; Africa; small and medium-sized enterprises; analysis*

## I. INTRODUCTION

Over the last decade, cloud computing has been a major agenda in the computing field. Cloud computing is the on-demand delivery of computer system resources as a service over the network. [1] The features of cloud computing, including scalability, flexibility and pay-per-usage model [2] has the potential to influence the various aspects of social and economic activities globally.

Cloud computing offers enormous benefits to all organizations and enterprises, including SMEs [3]. Small and medium Enterprises (SMEs -are the enterprises in which amount of personnel are less than certain limits, and they are control the data of high sensitivity. Some cases of sensitive data which is controlled by SMEs are: data of intelligence agencies and government federal, financial data of companies,

purchase contracts, company databases, de-identified research data, bank associated data like bank accounts, pin, passwords, balances and dealings, trade secrets, email accounts, drug formulas, accounting records and source codes [4].

The adoption of cloud computing is growing rapidly as it allows enterprises to concentrate on their essential business events, and, thus, efficiency is improved [5]. An adequately adopted cloud provides a plenty of benefits to the enterprises such as unlimited computing power, easy access of data and applications, lower IT expenditure, and build up competitive advantage. Recently, SMEs has shown a great concern in including cloud computing to their overall Information technology (IT) strategies. A recent report by Mckinsey [6] on the adoption rate of cloud service by SMEs informed that, 70% of SMEs have formerly bought at least two cloud service, and 40% have bought six or more cloud services. Nevertheless the touted benefits of cloud computing, its adoption and implementation in SMEs is faced with many challenges including national and international regulations, shortage of industry-specific conformism to principles, security and privacy threats [7][8][9]. Due to these challenges, some enterprises are quiet anxious around the threats of shifting business-critical applications to the cloud.

### A. Contribution and Paper Organization

Asia and Africa vary notably regarding their innovative capability. Asia has been competent to advance and sustain world leadership in innovation and technology, whereas Africa has not developed significantly in these terms. A very little comparative study is implemented on the reasons for innovation gap in these two continents. Motivated by this issue, this study sets out to examine and compare the cloud computing adoption from two regions; Asia and Africa. A comparative analysis is used to organize the findings from China in Asia, and Nigeria in Africa. Outcomes from this examination show that, in Nigeria security, privacy and trust, good internet connection, and level of awareness can pose as interferences to complete adoption of cloud computing. While, in China, these issues are not seen as hindering cloud computing adoption because of cheaper access to computing properties and reliable services in the cloud compared to locally hosted services.

The definite purposes of the paper were:

- To review systematically the extant literature regarding adoption of cloud computing amid SMEs in Nigeria and China.

- To recognize the probable benefits and usage of cloud computing amid SMEs in China and Nigeria.

- To contrast the level of cloud computing adoption amid SMEs in Nigeria and China.

- To determine main problems critical to complete adoption of cloud computing in the two states.

- To propose solution for Nigerian and Chinese SMEs to overcome these challenges.

The rest of the study is structured as follows: Section II presents the technical background of cloud computing and cloud computing adoption, by defining key concepts of cloud computing, discuss service and deployment models, and cloud computing adoption models. Section III presents the methodological approach for the literature search procedure. Section IV presents the explanation about SMEs in China and Nigeria, and also discusses about their impact to national GDP in the two countries. Section V presents the level of acceptance of cloud computing among SMEs in these two states. Section VI discusses the boundary between the adoption and practice of cloud computing in China and Nigeria. Section VII analyzes the potential benefits of cloud computing in Chinese and Nigerian SMEs. Section VIII examines some of the main competitions that impede the total adoption of cloud computing by SMEs in both states and recommends a series of recommendations to overcome these obstacles. Section IX covers the end of this review by highlighting the implications of the results.

## II. TECHNICAL BACKGROUND

### A. Cloud Computing Overview

Cloud computing is "an old idea whose time has (finally) come" [10]. As cloud computing has progressive nature, which is why it is hard to limit it to standard definition [11-13]. In the phrase Cloud Computing, the word "Cloud" is assumed to be evolved from (at least partly) the use of cloud symbol drawn in diagrams or flow charts as a metaphor, depicting some large networked environment or Internet [14].



Fig. 1.    Cloud Computing Environment.

The US-based National Institute for Standards and Technology (NIST) gave guidance on defining cloud computing. The NIST explains cloud computing as "a model for enabling convenient, on-demand network access to a shared pool of configurable computing resources (e.g., networks, servers, storage, applications, and services) that can be rapidly provisioned and released with minimal management effort or service provider interaction" [15]. Conferring to this description, cloud computing is beneficial to both the cloud service provider and cloud service user as shown in Fig. 1. In case of provider, the adaptability of cloud computing resources to scale according to service needs devoid of recompensing for this big scale is unique in the history of Information Technology (IT). In case of user, for retrieving responsive web-based applications, processing power on server-side require minimum system specifications for electronic devices [16].

The cloud computing has five vital features: [17]

- Rapid elasticity: Refers to near-instantaneous provisioning of capabilities of the application delivery infrastructure to expand and contract spontaneously according to needs.

- Broad network access: Potential and control of a cloud computing services can be retrieved via internet or other networks by means of standard protocols or thin/thick customer platforms e.g., work stations, laptops, mobile phones, and PDAs.

- Measured Service: Customers' use of the resources is optimized, reported, monitored and charged with some metering competencies, as an allocation to both the consumer and provider.

- On-demand self-service: Clients can separately use computing competences as desired, deprived of the requirement of human collaboration with the service's provider in the cloud.

- Resource pooling: Simulated and physical properties are dynamically assigned and reassigned to assist frequent consumers by a multi-tenant software conferring to consumers' requirements.

*1) Types of cloud Service models:* Types of cloud service models are shown in Fig. 2.

*a) Infrastructure-as-a-Service (IaaS):* In this form of cloud computing, the cloud service provider delivers the users with a "pay-per-use" cloud computing set-up over the network to set up and run arbitrary software [18].

*b) Platform-as-a-service (PaaS):* PaaS allows customers to generate web service-based programs quickly deprived of expense and complications of purchasing and handling the fundamental computing resources [19].

*c) Software-as-a-Service (SaaS):* In this cloud service, the consumers are provisioned to commercially accessible applications as a service [19].

*d) Data-as-a-Service (DaaS):* In this cloud service, comprises structures data and unstructed content (i.e. Content-as-a-Service). Data or information that is provided from the

cloud is in the form of raw data sets or used up by an analytics interface [20].

*e) Business Process-as-a-Service (BPaaS):* Cloud provided organization facilities that are adjusted to organization forms and related quantifiable business organization [20].

*2) Types of cloud deployment models:* The types of cloud deployment models are shown in Fig. 3.

*a) Public Cloud:* The framework of public cloud exist outside of the companies' own firewall and is provisioned to a big business set or the common community for open use [21].

*b) Private Cloud:* The framework of private cloud exists inside the company's own firewall and is provisioned to multiple employees of single organization (e.g., business units) for exclusive use [22].

*c) Community Cloud:* The framework of community cloud is available to customers of particular network for selective use, that have common interests like necessities, safety, strategy, mission, and agreement considerations [23].

*d) Hybrid Cloud:* The framework of hybrid cloud is a mixture of more than two clouds (mostly public and private), that stay exclusive units but remain composed by copyrighted technology [23].



Fig. 2. Cloud Computing Service Models.



Fig. 3. Cloud Computing Deployment Models.

### B. Cloud Computing Adoption: overview

Cloud computing adoption raises to the approval and contract to utilize cloud-support facilities as a novel method of installing applications. Company's effectiveness can be enhanced by deploying innovative services. Fig. 4 shows the cloud computing adoption models. Types of cloud computing adoption models are discussed below:

- Technology-Organization-Environment (TOE) Framework: The technological organization environment (TOE) is a model at the organizational level and a multi-point charter. The technological, organizational and environmental features of an organizational framework affect the method of adoption of technological invention [24].

- Theory of Reasoned Action (TRA): In this framework, all human actions is prophesied and clarified over three main cognitive modules, namely attitudes (bad luck or favor of behavior), social standards (social impact) and intents [25].

- Theory of Planned Behavior (TPB): In this framework, the perceived behavior control (PBC) by means of a novel inconstant to spread the TRA model is extended. In essence, PCB is resolute by the accessibility of assets, chances and skills, as well as by the assessed status of these assets, chances and capabilities to reach consequences [25].

- Theory of Interpersonal Behavior (TIB): This typical framework mostly illuminates the difficulty of human actions, which is unfair by societal and expressive aspects. Therefore, this model not merely offerings all the features of TRA and TPB, but moreover contributes to practices, facilitates circumstances and effects to recover predicted power [25].

- Technology Acceptance Model (TAM): This model originates from the TRA model. TAM clarifies the user's motivation over three features: perceived uses, ease of use and attitude to use.

- Diffusion of Innovations Theory (DOI): The DOI model observes a variety of inventions by means of announcing four features (time, communication channels, innovation, or social system) that effect the diffusion of a new notion [24].

- Perceived Characteristics of Innovating Theory (PCIT): This model extends to DOI theory by recognizing three further characteristics, namely: image, volunteerism and behavior [25].

- Unified Theory of Acceptance and Use of Technology (UTAUT): This analyzed four precursors for the adoption of information systems. Important events are the anticipation of effort, the expectation of results, social impact and conditions [25].

- Compatibility UTAUT (C-UTAUT): The objective of this model is to increase a improved considerate of how cognitive occurrences are done in the UTAUT model by identifying and testing new boundary conditions [25].

Fig. 4.  Cloud Computing Adoption Models.

## III.  METHODOLOGY

An organized comparative study was conducted on the comparison of cloud computing adoption in SMEs of China and Nigeria.

The following keywords, comprising truncation signs (symbolized by *), and Boolean operators (e.g., OR, AND), were chosen for this paper: (Cloud Computing) OR (Cloud Computing Adoption) OR (Comparison) AND (Analysis) OR (Comparative study) AND (Asia) OR (Africa) OR (China) OR (Nigeria) AND (SMEs) OR (Enterprises).

The following electronic databases were retrieved for the literature search, chosen due to their content being related to the discipline: Google Scholar and browsing over academic databases including IEEEXplore, Springer Link, ScienceDirect, and ACM Digital Library.

## IV.  INFLUENCE OF SMEs TO NATIONAL GROSS DOMESTIC PRODUCT (GDP) IN CHINA AND NIGERIA

Small and Medium-sized Enterprise (SMEs) have developed much significance in the worldwide budget, which cannot only be dignified by the growing quantity of SMEs signifying the 90% of total enterprises globally, but also their momentous part in producing occupation chances [27], decreasing poverty, making profits and helping state's prosperity.

The definition of SMEs is different for each country. In China, there are at least 300 employees in a medium-sized company with total assets and annual revenues that do not exceed 40 million RMB and 30 million as a result. Whatsoever fewer than that is categorized by means of a small-sized enterprise [28]. In Nigeria, the number of employees in SMEs ranges from 11 to 200, with a base of 5 million N and no more than 500 million (without houses and land).

In china, the quantity of SMEs are abundant comprising above 99% of the aggregate amount of businesses, amongst

which are small organizations [26]. Conferring to [29], Chinese SMEs offers 70% GDP, 62.3% of gross trade value, 74.7% of industrial added productivity rate and it also gives 80% of city employment. Consequently, it could be said that SMEs are vital for driving financial development in China. Fig. 5 shows the GDP of Chinese SMEs.

In Nigeria, it is expected that SMEs generate about 60% of the country's GDP and signify several commercial areas i.e. 50% as distributive trade, 10% in service, 10% in the manufacturing and 30% in agriculture as shown in Fig. 6 [30].



Fig. 5.  GDP of Chinese SMEs.



Fig. 6.  GDP of Nigerian SMEs.

China is taking advantage of cloud by delivering economical internet facilities to SMEs as a constituent of a determined virtualization platform [31].

## V.  ADOPTION OF CLOUD COMPUTING BY SMEs IN CHINA AND NIGERIA

It can be said that the adoption of cloud computing by SMEs in China has evolved to specific point (yet not completely settled). According to [32], the cloud services market for initiatives in the continent is estimated at 33.8 billion in 2016, a high-level adoption index for cloud computing by Chinese SMEs. In China, credible companies are being developed to use cloud computing, data transfer, computer facilities, software, full-service and commercial services.

TABLE I.    CONTRAST OF LEVEL OF CLOUD COMPUTING ADOPTION AMID CHINESE SMES AND NIGERIAN SMES

|  | SMEs of China | SMEs of Nigeria | Related Studies |
|---|---|---|---|
| Service Model | Platform as a Service , Infrastructure as a Service , and Software as a Service | Mostly Infrastructure as a Service , and Software as a Service | [32][34][35] |
| Deployment Model | Public cloud service , Private cloud service | Public cloud service , Private cloud service | [35] |
| Types of Industries | Health, Education, Computer Administrations and Software, Wholesale and Retail, Manufacturing, Government | Telecommunication, Computer Administration and Software | [32][36] |
| Barriers in Adoption | High Cost, Uncertainty regarding Security and Steadiness | Scarce Infrastructure, Security and Privacy, Good Internet Connection, Low Level of Awareness | [32][34][35] |
| Objective of usage | Email Services, Portal Website, Media services, Enterprise storage | Web Portal, Email Service, Enterprise storage | [35] |

With domestic maintenance for small organizations and the Internet setting, Chinese restaurants are booming. Conferring to [32], the Chinese network industry has reached 1,880 million CNY. Today, 36% of SMEs have a commercial website, of which 76% are SMEs and 94% of Chinese SMEs have already been tested for online management applications. In Africa, there is no new cloud computing, they are considered weak. According to [33], 24% of Nigerian SMEs (in 2015) used the facilities in their business and remained as a service (IaaS), monitored by software (SaaS).

In a study [33], it was noted that 100% of SMEs applied in the research use one or another system of ICT in their association, but that 41% of them are connected to the Internet. Apparently, 29% of SMEs use cloud applications in their company, and 24% of them presently use cloud facilities in their association. According to the previous research presented in [32] and [33], Chinese and Nigeria SMEs demonstrate the matches and variances amongst the adoption of cloud computing by SMEs in China and Nigeria. A series of comparison level of adoption are discussed in Table I.

Table I shows some of the matches and variances in the adoption of cloud computing by SMEs in China and Nigeria. The motive for this difference can be explained by the mass and financial strength of China, as the major economy globally. Additionally, the reason may be China's success in research and development and its contribution in science. All this gives China a lead in the domain of cloud computing. Although there are variances amid the use of cloud computing in Chinese and Nigerian SMEs, the document addresses a few problems that avert the complete implementation of cloud computing amid SMEs in both states.

## VI.    STRUCTURAL GAP BETWEEN CLOUD COMPUTING ADOPTION AND USAGE IN CHINA AND NIGERIA

The Chinese government led the Chinese IT leadership over tactical venture and support [37]. Government reflects cloud computing as a tactical precedence and contain it in the 12th Five-Year Plan of the country [38] [37]. Conferring to the Chinese administration and private companies, there is a very high investment in China, and China has a special capacity in this field. Cloud computing in China is 156 billion dollars in the coming years. An internet corporation in China named as "Tecent", declared in 2016 that in the coming six years, it proposes to devote 1.6 billion dollars in the cloud computing areas. The Alibaba group similarly declared a US $ 1 billion stock in Aliyun [38]. Consequently, we can reduce the level of acceptance of computer science in China.

Instead, the adoption of cloud computing technology is increasing in unindustrialized states. The road cloud has unlocked to deliver better-quality chances for developing countries [39], mainly in Nigeria. It cannot be said that the adoption of cloud computing in Nigeria has enlarged pointedly, but the IT market potential in Nigeria exceeds $ 100 million per year. According to [40], the subsequent main development is high-level cloud computing in Nigeria between IT professionals, government administrations and activities. But there are quiet obstacles to occupied acceptance in Nigeria, including: Data ownership and security, cloud information.

## VII. ADVANTAGES FOR CLOUD ADOPTION IN CHINA AND NIGERIA

Technical advance in cloud computing and its adoption has given much capability to SMEs in supporting administrations to accomplish their goals, increase economic benefit and delivering improved facilities to customers. SMEs that are unaware of the advantages of accessibilities of cloud computing are expected to be defeat in this greatly economical market. Table II indicates some of the most important advantages of cloud computing adoption.

TABLE II.    ADVANTAGES OF CLOUD COMPUTING ADOPTION AMONG CHINESE SMES AND NIGERIAN SMES

| Advantages | Description |
|---|---|
| Storage | Cloud-supporting systems have no environmental restrictions on records storing. The consumer can easily rise storing strategy deprived of spending much in computer equipment. |
| Expenses | Through recession of measure, cloud technology suppliers can allocate business-class quality to much little expenses SMEs, supporting small organizations to be extra dynamic compared to greater challengers and giving license fees for small businesses. This is also reduced. |
| Protected and Secure Data | Utmost of the data directed to the cloud is encoded, so remote sensors cannot be distinguished, and septic files are not spontaneously archived, somewhere they are, which avoids widespread network infections. |
| Disaster-proof | Cloud computing solutions support small businesses reduce the threat of data loss. Since records is presented and kept ubiquitously, society data can be delivered as rapidly as promising. |
| Reliability and accessibility | As the data is stored off-site and supported on other sites, all IT assets can be rapidly reestablished in case of server failure. |
| Focus on main capability | SMEs no longer have to concern about contracting an IT specialist, as cloud service providers admit this and can emphasis on the main capabilities of the association. |
| Adaptability and Remoteness | The cloud computing provision allows a reliable individual to entrance the essential assets over a protected Internet association and a well-suited device without working in the office. |
| Usage of updated applications | SMEs have entrance to the newest sort of the software (necessary for their work) as the service provider installs updates and manages the software licenses. |

## VIII. CHALLENGES FOR FULL ADOPTION IN CHINA AND NIGERIA

In China, despite of the fact that cloud computing has been considered much crucial and massively established, still many Chinese SMEs has not completely adopted it. The challenges hindering the complete cloud computing adoption amongst Chinese SMEs comprises costs (still expensive), and doubts and fears on the security and steadiness (as cloud computing is still at initial phase) of cloud computing [32].

The main challenges averting the complete adoption of cloud computing in Nigerian SMEs along with the solutions to beat these problems for taking lead of cloud computing are discussed in Table III.

TABLE III.    CHALLENGES AND THEIR SOLUTIONS FOR CLOUD COMPUTING ADOPTION AMONG NIGERIAN SMES

| | Privacy and Security | |
|---|---|---|
| | *Description of challenge* | *Solution of challenge* |
| Challenge 1 | The first and foremost challenge to Nigerian SMEs averting the complete cloud computing adoption is security and privacy [34]. Because according to many SMEs, for having assurance in adopting this novel technology, the security of data and privacy of SMEs needs to follow some strict standards [33]. | The problem of security and privacy of data deposited in a cloud of SMEs can be resolved by availability of cloud services to consumers on free trial for a specific duration. |
| | Good Internet Connection | |
| | *Description of challenge* | *Solution of challenge* |
| Challenge 2 | The second problem for Nigerian SMEs that is responsible for an extensive hindrance to complete cloud computing adoption is the accessibility of good internet connection. Because in utmost regions of Nigeria, the internet speed is gradual causing delay in delivery of data from/to the cloud [33] [41]. | To escalate the bandwidth and internet speed for downloading from the cloud or uploading to the cloud, the interested company needs to offer the enhanced internet connection, which will eventually reassure the clients for cloud computing adoption. |
| | Level of Awareness | |
| | *Description of challenge* | *Solution of challenge* |
| Challenge 3 | The third challenge is the very little knowledge about the advantages of cloud computing adoption in SMEs, which is why organization"s management team is anxious of adopting and approving cloud computing usage in SMEs [32][33][35]. | To eliminate the uncertainty about cloud computing adoption [42] [43], it is essential for companies (both public and private) in the country to have comprehensive information about cloud computing, its benefits and drawbacks. |

## IX. CONCLUSION AND FUTURE WORK

Cloud computing, that has been regarded as rapidly emergent technology, is still unfamiliar to many enterprises particularly small and medium-sized enterprises (SMEs). In china, the level of cloud computing adoption amid SMEs is assumed to be excessive (however not completely adopted). Although, the level of cloud computing adoption in African states mainly in Nigeria is comparatively deficient. A few issues averting the SMEs from complete cloud computing adoption have been debated in this paper. The capitalism and the government need to examine the challenges and their expected solutions to contribute the fundamental framework (great strategies, quick and reasonable internet, consistent power supply etc.) and permitting environment to Nigerian SMEs to be efficient for complete cloud computing adoption.

This study is momentous as it presents the comparative study of cloud computing adoption amid SMEs in the states (China, Nigeria) of two different continents (Asia and Africa), determine the main issues of adoption and propose possibly efficacious solutions for Nigerian SMEs to completely adopt the cloud computing.

REFERENCES

[1] Buyya, Rajkumar, Chee Shin Yeo, Srikumar Venugopal, James Broberg, and Ivona Brandic. "Cloud computing and emerging IT platforms: Vision, hype, and reality for delivering computing as the 5th utility." Future Generation computer systems 25, no. 6 (2009): 599-616.

[2] Dahiru, A. A., Julian M. Bass, and Ian K. Allison. "Cloud Computing: A comparison of adoption issues between Uk and Sub-Saharan Africa Smes." In Eur. Mediterr. Middle East. Conf. Inf. Syst. 2014, Oct. 27th-28th 2014, vol. 2014, no. 2010, pp. 1-12. 2014.

[3] Benton, D. "How cloud computing will influence banking strategies in the future." (2010): 2014. Retrieved 15 March, 2019, from www.accenture.com/banking

[4] Misra, Subhas Chandra, and Arka Mondal. "Identification of a company's suitability for the adoption of cloud computing and modeling its corresponding Return on Investment." Mathematical and Computer Modelling 53, no. 3-4 (2011): 504-521.

[5] Garrison, Gary, Sanghyun Kim, and Robin L. Wakefield. "Success factors for deploying cloud computing." Communications of the ACM 55, no. 9 (2012): 62-68.

[6] Avrane-Chopard, J., Th Bourgault, A. Dubey, and L. Moodley. "Big business in small business: Cloud services for SMBs." RECALL No 25 (2014), Available at: http://www.mckinsey.com/~/media/McKinsey/ dotcom/client_service/ HighTech/PDFs/Big_business_in_small_business _Cloudservices_for_SMBs.

[7] Armbrust, M., Armando Fox, Rean Griffith, Anthony D. Joseph, R. H. Katz, Andy Konwinski, Gunho Lee et al. "A view of Cloud Computing, Communications of the ACM." vol 53 (2010): 5058.

[8] Kern, Thomas, Jeroen Kreijger, and Leslie Willcocks. "Exploring ASP as sourcing strategy: theoretical perspectives, propositions for practice." The Journal of Strategic Information Systems 11, no. 2 (2002): 153-177.

[9] Marston, Sean, Zhi Li, Subhajyoti Bandyopadhyay, Juheng Zhang, and Anand Ghalsasi. "Cloud computing—The business perspective." Decision support systems 51, no. 1 (2011): 176-189.

[10] Fox, Armando, et al. "Above the clouds: A berkeley view of cloud computing." Dept. Electrical Eng. and Comput. Sciences, University of California, Berkeley, Rep. UCB/EECS28.13 (2009): 2009.

[11] Foster, Ian, Yong Zhao, Ioan Raicu, and Shiyong Lu. "Cloud computing and grid computing 360-degree compared." arXiv preprint arXiv:0901.0131 (2008).

[12] Gong, Chunye, Jie Liu, Qiang Zhang, Haitao Chen, and Zhenghu Gong. "The characteristics of cloud computing." In 2010 39th International Conference on Parallel Processing Workshops, pp. 275-279. IEEE, 2010.

[13] Zhang, Shuai, Shufen Zhang, Xuebin Chen, and Xiuzhen Huo. "Cloud computing research and development trend." In 2010 Second international conference on future networks, pp. 93-97. Ieee, 2010.

[14] Sultan, Nabil Ahmed. "Reaching for the "cloud": How SMEs can manage." International journal of information management 31, no. 3 (2011): 272-278.

[15] Mell, Peter, and Tim Grance. "The NIST definition of cloud computing." (2011).

[16] Hogan, Michael, Fang Liu, Annie Sokol, and Jin Tong. "Nist cloud computing standards roadmap." NIST Special Publication 35 (2011): 6-11.

[17] Mujinga, Mathias, and Baldreck Chipangura. "Cloud computing concerns in developing economies." (2011).

[18] Goscinski, A., and Brock, M. (2010). Toward dynamic and attribute based publication, discovery and selection for cloud computing. Future Generation Computer Systems, 26 (7), 947-970

[19] Brohi, Sarfraz Nawaz, and Mervat Adib Bamiah. "Challenges and benefits for adopting the paradigm of cloud computing." International Journal of Advanced Engineering Sciences and Technology 8, no. 2 (2011): 286-290.

[20] Mitchell, Ian. Isherwood, Stephen. The white bok of … cloud adoption. London: Fujitsu Services Ltd, 2011.

[21] Mather, Tim, Subra Kumaraswamy, and Shahed Latif. Cloud security and privacy: an enterprise perspective on risks and compliance. " O'Reilly Media, Inc.", 2009.

[22] Kim, Won, Soo Dong Kim, Eunseok Lee, and Sungyoung Lee.

[23] Dillon, Tharam, Chen Wu, and Elizabeth Chang. "Cloud computing: issues and challenges." In 2010 24th IEEE international conference on advanced information networking and applications, pp. 27-33. Ieee, 2010.

[24] Al-Hujran, Omar, Enas M. Al-Lozi, Mutaz M. Al-Debei, and Mahmoud Maqableh. "Challenges of cloud computing adoption from the TOE framework perspective." International Journal of E-Business Research (IJEBR) 14, no. 3 (2018): 77-94.

[25] Taherdoost, Hamed. "A review of technology acceptance and adoption models and theories." Procedia manufacturing 22 (2018): 960-967.

[26] Yu, Jia, and Jun Ni. "Development strategies for SME e-commerce based on cloud computing." In 2013 Seventh International Conference on Internet Computing for Engineering and Science, pp. 1-8. IEEE, 2013.

[27] Bao, Jinlong, and Xuewen Sun. "A conceptual model of factors affecting e-Commerce adoption by SMEs in China." In 2010 International Conference on Management of e-Commerce and e- Government, pp. 172-175. IEEE, 2010.

[28] Xiangfeng, Liu. "SME development in China: A policy perspective on SME industrial clustering." Asian SMEs and Globalization", ERIA Research Project Report 5 (2007).

[29] He, Yuan. "Sustainable development pattern of small and medium enterprises (SMEs) in China." In 2011 2nd International Conference on Artificial Intelligence, Management Science and Electronic Commerce (AIMSEC), pp. 1593-1596. IEEE, 2011.

[30] https://www.ukessays.com/essays/economics/impact-of-small-and-med iumenterprises-on-development-economic s-essay.php

[31] Kshetri, Nir. "Diffusion and Effects of cloud computing in China: Economic and institutional considerations." In PTC 2013 (Pacific Telecommun. Counc. Conf., pp. 20-21. 2013.

[32] Asia Cloud Computing Association. "SMEs in Asia Pacific: The market for cloud computing 2015." (2015): 1-217.

[33] Ofili, Onyeka Uche. "The use and challenges of cloud computing services adoption among SMEs in Nigeria." European Scientific Journal, ESJ 11, no. 34 (2015).

[34] Otuka, Richard, David Preston, and Elias Pimenidis. "The use and challenges of cloud computing services in SMEs in Nigeria." In proceedings of the European Conference on Information Management and Evaluation, vol. 43, no. 10, pp. 47-55. 2014.

[35] Awosan, R. K. "Factor analysis of the adoption of cloud computing in Nigeria." African Journal of Computing & ICT 7, no. 1 (2014): 33- 42.\

[36] Kshetri, Nir. "Cloud computing in developing economies." Computer 43, no. 10 (2010): 47-55.

[37] Ragland, Leigh Ann, Joseph McReynolds, Matthew Southerland, and James Mulvenon. Red cloud rising: cloud computing in China. Defense Group Incorporated, 2013.

[38] To, Wai-Ming, Linda SL Lai, and Andy WL Chung. "Cloud computing in China: barriers and potential." IT Professional 15, no. 3 (2013): 48-53.

[39] Muhammad, Akilu Rilwan. "Towards cloud adoption in Africa: The case of Nigeria." International Journal of Scientific & Engineering Research 6, no. 1 (2015): 657-664.

[40] Dogo, Eustace Manayi, Abdulazeez Salami, and Salim Salman. "Feasibility analysis of critical factors affecting cloud computing in Nigeria." International Journal of Cloud Computing and Services Science 2, no. 4 (2013): 276.

[41] Abubakar, A. D., Julian M. Bass, and Ian Allison. "Cloud computing: Adoption issues for sub-Saharan African SMEs." The Electronic Journal of Information Systems in Developing Countries 62, no. 1 (2014): 1-17.

[42] Muhammed, Kuliya, Kabir Rumana Isma'il Zaharaddeen, and Abdulkadir M. Turaki. "Cloud computing adoption in Nigeria: Challenges and benefits." International Journal of Scientific and Research Publications 5, no. 7 (2015): 1-7.

[43] Dahunsi, F. M., and T. M. Owoseni. "Cloud computing in Nigeria: The cloud ecosystem perspective." Nigerian Journal of Technology 34, no. 1 (2015): 209-216.

"Adoption issues for cloud computing." In Proceedings of the 7th International Conference on Advances in Mobile Computing and Multimedia, pp. 2-5. ACM, 2009.

# Exploratory Analysis of the Total Variation of Electrons in the Ionosphere before Telluric Events Greater than M7.0 in the World during 2015-2016

Alva Mantari Alicia[1], Zarate Segura Guillermo Wenceslao[2], Sotomayor Beltran Carlos[3]

Brian Meneses-Claudio[4], Roman-Gonzalez Avid[5]

Image Processing Research Laboratory (INTI-Lab), Universidad de Ciencias y Humanidades, Lima, Perú

*Abstract*—This exploratory observational study analyzes the variation of the total amount of vertical electrons (vTEC) in the ionosphere, 17 days before telluric events with grades greater than M7.0 between 2015 and 2016. Thirty telluric events have been analyzed with these characteristics. The data was obtained from 55 satellites and 300 GPS receivers that were downloaded from the Center for Orbit Determination in Europe (CODE). The variations are considered significant only if it is outside the "normal" ranges considered after the statistical analysis performed. The data was downloaded by a program developed in our laboratory. The downloaded data was processed and maps of variations of vTEC generated with a periodicity of 2 hours. The analysis area was considered to be a circular one with a radius of 1000km centered on the epicenter of each earthquake. Variation of vTEC was found during 2015-2016 in 100% of the earthquakes in the range from day 1 to day 17 days before the event, over the circular area of 1000 km radius centered on the epicenter of the earthquake. Of these in 96.55% there are positive variations and a negative one exist in 68.97% of the events. If we observe in the range from day 3 to 17 before the event, a variation was recorded in 100% of the cases, and from day 8 to day 17 before the event in 93.10% of the cases, it is important to emphasize that while the evidence in a period before the event is more likely to find evidence to develop early warning tool for earthquake prevention. This study explores the variation of vTEC as precursor events to each earthquake during 2015-2016; it is a preliminary analysis that shows us the feasibility of analyzing this information as a preamble for an exhaustive association study later. The final objective is to calculate the risk of telluric events which would benefit the population worldwide.

*Keywords*—*Total number of electrons; ionosphere; earthquakes; prevention; risk*

## I. INTRODUCTION

Research related to the prevention of earthquakes is increasing in the last decades, due to the destruction and deaths they generate, this makes it imperative to look for methods that will help us predict them effectively. Earthquakes are movements of the earth's crust, which are a consequence of the sudden release of energy, due to the movement of tectonic plates [1]. The earthquakes are cataloged by different methods, depth or the energy that releases the movement. One of these scales is that of Richter, which quantifies the energy released by seismic movement. These exist in a range of dimensions from micro tremors of less than M2.0 to maximum values of more than *M10.0* [2][3], although there is no reliable report of having reached M10.0, it is presumed that in the prehistoric era was reached these measures, the largest recorded earthquake on this scale is M9.5, in Chile on May 22nd, 1960. The greater intensity, greater destruction and damage to hundreds or thousands of kilometers, for example those in Sichuan (May 12th, 2008, M7.9), Samoa (September 29th, 2009, M8.1) and Haiti (January 12th, 2010, M7.0) were responsible for approximately 380,000 deaths and millions of victims [4]. We can notice that the earthquake greater than *M7.0* is the most dangerous, which there is an occurrence record of up to 18 per year [5] approximately. Our study has focused on earthquakes with these characteristics. Seismic grades greater than *M7.0,* considered the most devastating.

The study of the ionosphere has been of interest among different research groups due to its association with the prediction of telluric events [6]. The ionosphere is the area of the ionized atmosphere at a height greater than 70-80 km, in which free electrons and ions form plasma, all data are taken at 450 km of distance in the ionosphere. Solar radiation is generated by the ionization of the area, which depends on solar activity, location and time of observation [7].

The use of this method analyzes the electron concentration of the ionosphere as an alternative of prediction [8, 9]. It is the Total Quantity of Electrons (TEC) that is analyzed in its relation with the telluric movements. The TEC is defined based on the electronic density and the length of the path traveled by the solar emission. For our present study, we use the vertical number of electrons (vTEC), which is the projection of the TEC on a vertical path calculated as shown in Fig. 1.

To consider that an anomaly can be considered as an ancestor of an earthquake, we must fulfill 4 important requirements: (1) To be related to the deformation or the tectonic mechanism before earthquakes. (2) Be observed in 2 or more events. (3) Possess an amplitude in relation to the distance of the earthquake. (4) Define a risk area. Whose size must be respect to the studied area of the earthquake [10].

The global Ionosphere Maps (GIM) used in our research to analyze the association between the variation of the vTEC and seismic events greater than *M7.0* during the years 2015-2016, were generated by the CODE, which receives data from the satellites and it provides vTEC information in the ionosphere, from a combined analysis of the values received from 55 satellites and 300 GPS receivers [11].

Fig. 1.   Representation of Oblique TEC (sTEC) and Vertical TEC (vTEC).

## II.   POPULATION

### A.  Sample

The events that we are going to study for the association are those manifested during the 2015-2016 year. There were 30 seismic events greater than *M7.0* degrees worldwide. These have been extracted from the database of IRIS Earthquake Browser [12], these have been geolocated on the map of the world using Qgis, as shown in Fig. 2.

These events have been placed in a Table I where it shows the place of the event, the latitude and longitude of the event, the grade of magnitude, and date, all data are shown in Table I.

Circular neighborhoods of 1000 km radius have been generated for each of the events [6], where the variation of vTEC will be analyzed during a predefined period of time of 17 days before the event. We can visualize the neighborhoods for the set of analysis events in the following Fig. 3.

The data of the vTEC variations have been generated for each hour and day of the event; from the data obtained by the code, the mean and the standard deviation of the vTEC variation are calculated from the data by a range of 8 days before and then, the significant variation which is considered over the ranges considered normal after the aforementioned statistic. So, we have maps for the analysis for each day of an event in 2015-2016, there are 17 days before the event and 4 days later, although we will only analyze the data of 17 days previously. The map of vTEC variations in the ionosphere is analyzed in rectangular areas of 2.5 latitude, and 5 longitude. The amount of data handled in this study can be seen in the following Table II.

TABLE I.       ABLE OF SEISMIC EVENTS OF 2015-2016 WITH INTENSITIES GREATER THAN M7.0

| Seismic events 2015-2016 | | | | |
|---|---|---|---|---|
| *Place* | *Latitude* | *Longitude* | *Grade* | *Date* |
| Dorsal Meso-Atlantic North | 52.51 | -32.02 | 7.1 | 13/02/2015 |
| Indonesia | -7.36 | 122.49 | 7.0 | 27/02/2015 |
| Papua Nueva Guinea | -4.78 | 152.58 | 7.5 | 29/03/2015 |
| Nepal | 28.13 | 84.72 | 7.8 | 25/04/2015 |
| Japan | 27.83 | 140.49 | 7.8 | 30/05/2015 |
| Nepal | 27.80 | 86.13 | 7.3 | 12/05/2015 |
| Dorsal Meso-Atlantic South | 35.36 | 17.16 | 7.3 | 17/06/2015 |
| Solomon Islands | -10.46 | 165.1 | 7.0 | 18/07/2015 |
| Indonesia | -2.71 | 138.54 | 7.0 | 28/07/2015 |
| Chile | -31.64 | -71.74 | 8.4 | 16/09/2015 |
| Afghanistan | 36.52 | 70.37 | 7.5 | 26/10/2015 |
| Vanuatu | -14.86 | 167.3 | 7.1 | 21/10/2015 |
| Peru | -10.06 | -71.02 | 7.7 | 24/11/2015 |
| Pakistan | 38.21 | 72.78 | 7.2 | 07/12/2015 |
| Russia | 54.01 | 158.51 | 7.2 | 30/01/2016 |
| Indonesia | -4.95 | 94.33 | 7.8 | 02/03/2016 |
| Japan | 32.79 | 130.75 | 7 | 15/04/2016 |
| Ecuador | 0.38 | -79.92 | 7.8 | 16/04/2016 |
| Vanuatu | -16.04 | 167.38 | 7 | 28/04/2016 |
| Visokoy Island | -56.24 | -26.94 | 7.2 | 28/05/2016 |
| Northern Mariana Islands | 18.54 | 145.54 | 7.7 | 29/07/2016 |
| New Caledonia | -22.48 | 173.11 | 7.2 | 12/08/2016 |
| Argentina | -55.28 | -31.87 | 7.1 | 19/08/2016 |
| Santa Elena, Ascencion and Tristan de Acuña | -0.05 | -17.83 | 7.1 | 29/08/2016 |
| New Zealand | -37.36 | 179.15 | 7 | 01/09/2016 |
| New Zealand | -42.72 | 173.06 | 7.8 | 13/11/2016 |
| El Salvador | 11.96 | -88.84 | 7.2 | 24/11/2016 |
| Solomon Islands | -10.68 | 161.33 | 7.8 | 09/12/2016 |
| Papua New Guinea | -4.5 | 153.52 | 8 | 17/12/2016 |
| Chile | -43.41 | -73.94 | 7.6 | 25/12/2016 |
| Russia | 54.01 | 158.51 | 7.2 | 30/01/2016 |



Fig. 2.   Representation of the 30 Seismic Events Greater than M7.0, Worldwide. Red 2015. Blue 2016.



Fig. 3.   Representation of the 15 Work Areas of 1000 Km Radius in Each Seismic Event Greater than M7.0 during 2015 in Red, and during 2016 in Blue.

| Visualization of Data Managing obtained for the Analysis | | | | |
|---|---|---|---|---|
| *Earthquakes* | *Days* | *Hour* | *Rectangular Area Names* | *Total data input* |
| 30 | 17 | 12 | 5184 | 31'726,080 |

From these data, maps with smoothed neighborhoods were generated with a periodicity of 2 hours (Fig. 4), which are the data we use to make the qualitative comparison in the 17 days prior to each event.

In Fig. 5, we show how the vTEC variations are generated for the 04:00 UT hour, these are calculated from the normal variations of the average normal TECU for the same area and the same time.



Fig. 4.    Extraction of Data for Each Day of Interest, for a Period of 2 Hours.



Fig. 5.    Differences between the Image of the TEC Values and the Calculated Variations in a Window of 8 Days on Average, Only those that have a Variation Outside the Range for the Section and the Hour are Plotted.

The analysis performed has taken the data of the generated images as shown above and the presence of those variations on the circular area with a radius of 1000 km. Centered on the GPS position of the seismic event analyzed. This radius has been calculated from the deduction of deformations and

inclinations on the surface of the Earth as a function of the magnitude of the earthquake being prepared and the distance from the epicenter. These data are analyzed in the present study, in which we look for the association between seismic events during 2015 and 2016, and the appearance of a vTEC variation present up to 17 days before the event.

### III.   METHODOLOGY

The study carried out is an exploratory observational data model, to associate the frequency in which, given a telluric event greater than *M7.0*, we find evidence of a significant variation of vTEC in the ionosphere in a circular area with a radius of one kilometer from the event.

For this reason, the images generated from the CODE data, up to 17 days previously, have been analyzed. For all telluric events of 2015, and we have analyzed frequency, for days, for hours, for positive variation, and negative.

We will show below the work done for one of the events in Peru (Latitude -10.06, Longitude -71.02) on November 24[th], 2015 of magnitude *M7.7*.

Fig. 6 shows the exact location of the seismic event on the world map.

The area of interest to analyze is a circle of radius of 1000 km (Fig. 7).

We searched each map generated every 2 hours in the period of 17 days before the event and 4 days later, the presence of significant variations in the vTEC, such as it is described in the methodology on the area of interest marked in Fig. 7. We show some examples of how these maps have been qualitatively visualized in this initial study (Fig. 8, 9, 10 and 11).



Fig. 6.    Representation of the Seismic Event in Indonesia on March 2[nd], 2016.



Fig. 7.    Representation of the Seismic Event in Peru on March 2[nd], 2016, with the Circular Area of 1000 km. around the Point of the Earthquake.

Fig. 8. Day 15 before the Event on March 2$^{nd}$, 2016 (February 15$^{th}$, 2016).



Fig. 9. Day 14 before the Event on March 2$^{nd}$, 2016 (February 16$^{th}$, 2016).



Fig. 10. Day 6 before the Event on March 2$^{nd}$, 2016 (February 24$^{th}$, 2016).

Fig. 11. Day 4 before the Event on March 2<sup>nd</sup>,2016 (February 26<sup>th</sup>, 2016).

The analysis of the data for this event was carried out in which we can observe the following results. Given the event in Indonesia, a variation was observed in 12 of the 17 days analyzed previously, the presence of a significant vTEC variation, with a total of 66 hours approximately. In 34 of these hours, the variation was positive with a higher load than the expected average, and only in 32 hours was a variation less than normal considered by the statistical calculation. All these variations analyzed on the area of defined influence plotted in the monitoring of variations shown.

## IV. RESULTS AND CONCLUSIONS

The process explained for the analysis was carried out in each of the telluric events of 2015-2016. We analyzed and obtained the following table of results. Table III, where we can observe the event, the number of days, the hours and the type of variation of vTEC that were found.

The objective of this study was to observe and explore the association of a variation of vTEC before the telluric event during 2015-2016 based on qualitative analysis from the images generated by the software from INTI-Lab research team led by Dr. Carlos Sotomayor.

TABLE III.     DATA ANALYSIS EVENT OF 2015-2016

| Register of observed variations of the telluric events 2015-2016 | | | | | |
|---|---|---|---|---|---|
| *Place* | *Year* | *Total days with vTEC variation 17 days before* | *Total hours with vTEC variation 17 days before* | *Positive variation of vTEC* | *Negative variation of vTEC* |
| Dorsal Meso-Atlantic North | 2015 | 3 | 22 | 18 | 4 |
| Indonesia | 2015 | 6 | 22 | 22 | 0 |
| Papua Nueva Guinea | 2015 | 2 | 34 | 16 | 18 |
| Nepal | 2015 | 6 | 38 | 28 | 10 |
| Japan | 2015 | 4 | 24 | 24 | 0 |
| Nepal | 2015 | 4 | 34 | 22 | 12 |
| Dorsal Meso-Atlantic South | 2015 | 3 | 16 | 16 | 0 |
| Solomon Islands | 2015 | 5 | 26 | 20 | 6 |
| Indonesia | 2015 | 3 | 12 | 12 | 0 |
| Chile | 2015 | 3 | 10 | 8 | 2 |
| Afghanistan | 2015 | 2 | 6 | 4 | 2 |
| Peru | 2015 | 7 | 28 | 26 | 2 |
| Pakistan | 2015 | 2 | 6 | 6 | 0 |
| Russia | 2016 | 4 | 28 | 14 | 14 |
| Indonesia | 2016 | 12 | 66 | 34 | 32 |
| Japan | 2016 | 12 | 48 | 40 | 8 |
| Ecuador | 2016 | 12 | 32 | 26 | 6 |
| Vanuatu | 2016 | 8 | 48 | 40 | 8 |
| Visokoy Island | 2016 | 4 | 10 | 10 | 0 |
| Northern Mariana Islands | 2016 | 7 | 22 | 20 | 2 |
| New Caledonia | 2016 | 3 | 12 | 10 | 2 |
| Argentina | 2016 | 1 | 2 | 0 | 2 |
| Santa Elena, Ascencion and Tristan de Acuña | 2016 | 6 | 24 | 8 | 16 |
| New Zealand | 2016 | 3 | 8 | 8 | 0 |
| New Zealand | 2016 | 6 | 42 | 38 | 4 |
| El Salvador | 2016 | 5 | 32 | 30 | 2 |
| Solomon Islands | 2016 | 6 | 22 | 20 | 2 |
| Papua New Guinea | 2016 | 7 | 26 | 24 | 2 |
| Chile | 2016 | 4 | 14 | 14 | 0 |

TABLE IV.    DATA ANALYSIS EVENT OF 2015-2016

| DAYS BEFORE THE EVENT | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *Place* | *0* | *1* | *2* | *3* | *4* | *5* | *6* | *7* | *8* | *9* | *10* | *11* | *12* | *13* | *14* | *15* | *16* | *17* | *T* |
| Dorsal Meso-Atlantic North | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 2 |
| Indonesia | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 6 |
| Papua Nueva Guinea | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 2 |
| Nepal | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 5 |
| Japan | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 4 |
| Nepal | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 4 |
| Dorsal Meso-Atlantic South | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 |
| Solomon Islands | 0 | 1 | 1 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 5 |
| Indonesia | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 3 |
| Chile | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 |
| Afghanistan | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 |
| Peru | 0 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 7 |
| Pakistan | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 |
| Russia | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 4 |
| Indonesia | 1 | 1 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 12 |
| Japan | 1 | 1 | 1 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 1 | 0 | 1 | 1 | 12 |
| Ecuador | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 1 | 12 |
| Vanuatu | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 8 |
| Visokoy Island | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 4 |
| Northern Mariana Islands | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 7 |
| New Caledonia | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 3 |
| Argentina | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 |
| Santa Elena, Ascencion and Tristan de Acuña | 1 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 6 |
| New Zealand | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 3 |
| New Zealand | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 6 |
| El Salvador | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 5 |
| Solomon Islands | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 6 |
| Papua New Guinea | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 7 |
| Chile | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 4 |

In Table IV, we can see that each event had manifestations of significant variation in 100% of the cases and the days in which it was present in a dichotomous manner.

From the table, we see that the minimum manifestation occurred in only 1 of the 17 days before the event and the maximum manifestation occurred in 12 days of the 17 days before the event. It should be noted that in the 96.55% of events previously, there is a positive variation, and that only in 68.97% of events the variation of vTEC is negative. On average in 2015-2016, the variation was maintained for the period of 5.17 days previously; this variation was presented for approximately 24.62 hours. That means the manifestation is not momentary, there is a period of time where variation is in the area of analysis selected for each event during 2015 and 2016. In Fig. 12, we can observe the frequency of events that have a manifestation of the variation of vTEC in the "n" days before the earthquake.



Fig. 12. Frequency of Events Days before the Earthquake.

A vTEC variation was found in a circular area of 1000 km radius at least one day before the seismic event of the 2015-2016 data. It is also important to mention that this manifestation has not occurred once since it has been a behavior that was observed in more than 1 day before the earthquake event. We have observed that a significant variation has been found in 100% of the cases, even in the range from 3 to 17 days before the event; it is 96.55% when range is from day 5 to day 17 previous the event; and finally, if our range is restricted from day 8 to day 17 the variation is presented in 93.10% of cases. These data are necessary, in our case, because we want an algorithmic tool to model a risk and calculate with more prior days for a possible earthquake calculated from the analysis of the vTEC in the ionosphere.

The physical mechanism as to how these events happened before the earthquakes can be explained by the model of ionization of the air due to radiation from the ground [13,14,15]. This model estates that when tectonics plates collide with each other, besides producing earthquakes, radon is produced as a product of these collisions. This gas along with others is then released to the lower atmosphere. The presence of radon in the lower atmosphere will then provoke the ionization of the other gases presented here. However, it has been indicated [15] that other factors like air temperature or relative air humidity can also affect atmospheric ionization. Depending on the level of ionization of the atmosphere positive or negative variation of vTEC can be observed. If the air ionization produces heavy ions, then the air conductivity will decrease causing the increment of electrons in the ionosphere. This in turn will produced the positive variations of vTEC. Such mechanism is very likely to the cause of the positive variations observed before the earthquakes that happened in Indonesia, Nepal, Japan, Chile, Peru, Ecuador, New Zeland and El Salvador. On the other hand, if the level of air ionization is low, light ions will be produced. This type of ions will increase the air conductivity and therefore a decrease on is vTEC expected. This decrease of vTEC translates in to negative variations of vTEC. We suggest that the negative variations observed before the earthquakes that occurred in Papua Nueva Guinea, Nepal, Russia, Indonesia and Santa Elena, Ascencion and Tristan de Acuña, are very likely the product of the aforementioned mechanism. In order to verify this theory of the collision of tectonics plates been the cause of the variations of vTEC, studies of radon in the locations where the earthquakes happen would be beneficial [16].

It is important to mention that this study is preliminary because it has served to explore the association between telluric events and the variation of vTEC in the preliminary days during a period of 2 years. The objective of our next study will be to observe causality between a variation of vTEC and telluric events in a given area.

REFERENCES

[1] N. N. Ambraseys, "Value of Historical Records of Earthquakes," Nature, vol. 232, no. 5310, pp. 375–379, Aug. 1971.

[2] B. Gutenberg and C. F. Richter, "Magnitude and energy of earthquakes," Ann. Geophys., pp. 7–12, 2010.

[3] B. Gutenberg and C. F. Richter, "Magnitude and Energy of Earthquakes," Nature, vol. 176, no. 4486, p. 795, Oct. 1955.

[4] [4]S. Orjuela, D. A. Molina, and F. Y. Zapata, "en la ionosfera durante un terremoto," p. 6, 2015.

[5] "Catalogo Sismico De La Zona Comprendida Entre Los Meridianos 5o E. Y 20o W. De Greenwich Y Los ParalelOS 45o Y 25o N. Tomo II | José Galbis Rodriguez | Comprar libro mkt0003098317." [Accessed: 25-Aug-2018].

[6] I. P. Dobrovolsky, S. I. Zubkov, and V. I. Miachkin, "Estimation of the size of earthquake preparation zones," Pure Appl. Geophys., vol. 117, no. 5, pp. 1025–1044, Sep. 1979.

[7] M. T. C. Pérez and V. M. M. Castro, "Climatología, Cambios Climáticos Y Atmósfera," vol. 30, no. 1, p. 10, 2010.

[8] Z. Fuying, W. Yun, Z. Yiyan, and L. Jian, "A statistical investigation of pre-earthquake ionospheric TEC anomalies," Geod. Geodyn., vol. 2, no. 1, pp. 61–65, Feb. 2011.

[9] M. H. Sarachaga and F. Sánchez-Dulcet, "Efectos De La Actividad Sísmica En La Ionosfera; Características Y Posibles Aplicaciones," p. 13.

[10] M. Wyss, "Evaluation of Proposed Earthquake Precursors," Eos Trans. Am. Geophys. Union, vol. 72, no. 38, pp. 411–411, Sep. 1991.

[11] F. Zhu, Y. Wu, J. Lin, and Y. Zhou, "Temporal and spatial characteristics of VTEC anomalies before Wenchuan Ms8.0 earthquake," Geod. Geodyn., vol. 1, no. 1, pp. 23–28, Jan. 2010.

[12] "IRIS Earthquake Browser." [Online]. Available: http://ds.iris.edu/ieb/index.html?format=text&nodata=404&starttime=1970-01-01&endtime=2025-01-01&minmag=0&maxmag=10&mindepth=0&maxdepth=900&orderby=time-desc&limit=1000&maxlat=89.21&minlat=-89.21&maxlon=180.00&minlon=-180.00&zm=1&mt=ter. [Accessed: 25-Aug-2018].

[13] C. Sotomayor-Beltran, "Ionospheric anomalies preceding the low-latitude earthquake that occurred on April 16, 2016 in Ecuador," Journal of Atmospheric and Solar-Terrestrial Physics, vol. 182, pp. 61-66, 2019

[14] S. Pulinets, "Low-latitude atmosphere-ionosphere effects initiated by strong earthquakes preparation process," Int. J. Geophysics, vol. 2012, pp.1-14, 2012

[15] V.M Sorokin, "Plasma and electromagnetic effects in the ionosphererelated to the dynamics of charged aerosols in the lower atmosphere," Russian Journal of Physical Chemestry B, vol 1, pp.138-170, 2007

[16] S. Pulinets, "Radon and ionosphere monitoring as a mean for strong earthquakes forecast," Nuovo Cimento, p. 621-626, 1999.

# Fast and Efficient In-Memory Big Data Processing

Babur Hayat Malik[1], Maliha Maryam[2], Myda Khalid[3], Javaria Khlaid[4], Najam Ur Rehman[5], Syeda Iqra Sajjad[6]
Tanveer Islam[7], Umair Ahmed Butt[8], Ali Raza[9], M. Saad Nasr[10]
Department of CS and IT, University of Lahore, Chenab Campus, Gujrat Pakistan

*Abstract*—With the passage of time, the data is growing exponentially and the mostly endured areas are social media networks, media hosting applications, and servers. They have thousands of Tera-bytes of data and the efficient systems, however, they are as yet confronting issue to oversee such volume of information and its size is growing each day. Data systems retrieve information with less time of In-memory. Instead of each factor data systems are required to define good usage of cache and fast memory access with help of optimization. The proposed technique to solve this problem can be the optimal indexing technique with better and efficient utilization of Cache and having less overhead of DRAM with the goal that energy can also be saved for the high-end servers.

*Keywords—Big data processing; indexing techniques; R-tree; B-tree; X-tree; hashing; inverted index; graph query tree*

## I. INTRODUCTION

### A. Big Data

Big Data is large datasets whose scale, diverseness, and complexity involve new strategic technologies for algorithmic computer program, and analytics to manage it and excerpt value and hidden knowledge from it. With the passage of time the data is growing exponentially and the mostly domains that hold thousands of Terabytes of data are social media networks, media hosting applications, and servers [1]. Big data is basically an umbrella term for data that are too large to handle. Characteristics of data involve 4V's.

1) Volume (Scale)
2) Varity (Complexity)
3) Velocity (Speed)
4) Veracity

Day by day data volume is increasing exponentially and it needs to be processed fast. Due to increase in the volume of data, it is also complex to handle because one application is generating different types of data containing several formats, types, and structures, text, numerical, images, audio, video, sequences, time series, social media data, multi-dim arrays, etc. [2]. The veracity of data is measured by checking its accuracy and types of data structured and unstructured and it is difficult to check the veracity of unstructured data.

Growth of data has become an issue for the past few decades because it is not easy to manage that data and store it, perform computation and extract any information from that size of data. There are many solutions for that type of data while managing it some says to classify it to make the retrieval fast. Some says to store it in efficient manner in most upgraded machines with High performance processors, Main memory and also secondary memories. In Memory Systems are better than other disk based systems because performance of main memory is better than hard drive [3]. As field of computing is very vast and progressing, each day is seen with new inventions and innovations in every domain and sub-domain, like the SSD. A secondary storage drive and PCM is very fast than the old HDD. Number of Computations per second in these No-Volatile memories is far greater than the simple hard drive. It is very integrated in size and very efficient access time [4]. Modern Servers in these days normally have hundreds of gigabytes of DRAM and tens of cores while the fastest of them have TB's of DRAM and Peta-bytes of secondary storage with hundreds of cores to process the gigantic size of data [5]. In-memory systems have been discussed in 1980s [6], and after that is has not been studied. However, recent progress in the computing era has changed the previous work entirely and developed an interest to host the entire data in Main-memory to perform faster access data analytics [7]. In-memory Data processing has been widely used to support large-scale applications totally in DRAM. Different aspects can be considered for optimizing in-memory system such as indexing, data layouts, Parallelism, Fault-tolerance, Data-overflow, concurrency control and query processing [8], [9]. Indexing deals with cache utilization and parallelism deals with data compression [10] and fast processing means packing multiple values in single processor word.

Parallelism is of two types scale-up and scale-out and both of these optimize performance by partitioning data. Concurrency Control is another aspect that has great impact on the performance of data analysis [11]. Different mechanism of concurrency control has different demerits like Heavy-weight mechanism consists of too much locks/semaphores which are the key ingredient of degrading the system performance. Similarly, other mechanisms are Light-weight Intent Lock (LIL) and very lightweight locking (VLL) simplifies the data structure by compressing all the lock states and some are based on time stamps [12], [13]. Other efficient techniques are MVCC "Multi-Version Concurrency Control" Search and Retrieval of this type of content has become a challenging task now a-days because the data volume is very high and it is not easy to manage such size of data and performing computation on that type of data.

Using Distributed Environment is one of the solution so that performance of computation on that data can be optimized [14]. In memory data storage defines performance improvement attribute because the latency of DRAM is much less than that of permanent Storage and the data in RAM can be retrieved fast because the processor can better access DRAM than the hard drive. But this still needs some kind of

optimization that can make this process space and time efficient and can retrieve data optimally.

Many Indexing techniques have been proposed such as Hash-based, Tree-based but all of them have deficiencies of any kind such as some of them have performance issues and others are not time efficient. Few enhance cache by performing huge computations while other fails to correctly compute cache. So there should be an optimal solution that can search huge data optimally. All the work about it involves in improving the indexing technique which could be to make a new and different algorithm that is best among the others implemented before or it can be the combination of any two best indexing techniques means.

### B. Optimizing Big Data Processing

With the Increase of Data volume and size, it has become difficult to manage and process it efficiently. Although many techniques have been proposed in this era, as indexing, parallelism, data layouts management, concurrency control and fast query processing but it still needs optimization so that its processing time can be decreased. As the project is about In-memory big data processing so it should be reminded that in this case, a copy of database should be kept in main-memory which is DRAM all the time so that it should be processed fast but it still needs an efficient indexing technique to process more fast with less cost and time. The whole project is about optimizing indexing technique of big data to make it process fast and the end its energy should be measured because all the work overhead depends upon DRAM that's why the energy should be kept in mind.

This paper is organized as follows: Section II presents the comprehensive literature of previous work. Section III presents a comparisons based on related literature. Finally, Section IV concludes the paper and future work is also discussed.

## II. LITERATURE REVIEW

### A. Research Paper Summaries

Author of paper [1] described that for data mining, creative technologies, and analysis of data and prediction the word "Big data" appears on most of the specialized meetings around the world. Where the organization of huge sums of information is toward the best real work this word is applied in these zones. To characterize these zones, a constant rise of information brooks in the administrative procedure, whether it is an economy, savings, making, selling, broadcastings, treatment, etc. The knowledge of huge information place administration inside an equivalent treating of assets for best technology background is mainly related. The training of in memory knowledge application structures for large information groups on heavy technology which is based on SAP HANA knowledge's used as information storage and presented in this paper. SAP HANA Colum store method read information fast and data compression methods efficiency rise. The compacted information procedure lets decomposition process resources cheaper. With a similar engineering of possessions, the prices economy for large information saving and dispensation with the practice of in-memory method founded on SAP HANA let's decrease bulks of information

stored in the main system information. Different to old-style files shaped for hard disk connecting, the communication of RAM and the processor is the initial aim of SAP HANA. A method [1], like a column store is applied in that situation, which moreover takes a quick action of data reading; open the ways so that data compression mechanism is applied in a good way. Deprived of outgoings wealth on the decompression procedure, SAP HANA works with compacted information in a direct way. IT facilities management that is incorporated into the IT setup of an information founded boldness, now takings residence with the service concerned with communal material schemes usage.

Author of [2], described that the use of technology is greatly improved by the installation of sensors. As time has increased the requirements have also increased. In addition to that, the number of problems has also increased which are related to the equipment of an industry. Therefore, it is required that new systems are introduced for obtaining the best performance and storage of data. Industries use a large number of sensors to monitor, process and storage of data. In the past, the systems having particular needs were used in the industries that could not have the capability to process the data in a fast way and store it. For the storage of data in the middle of working, fast processing and storing high volume the In-memory technology is the most suitable one. In this paper [2] the analysis of data at high speed and having high volume is done by developing the prototype having in-memory at the core. So, the time series data is stored continuously and simultaneously analyze the data. The case is discussed in which the memory is stored at the rate of 10000 data points/sec. Also, the analytical collection is implemented for IMDG and the data is updated in a continuous manner. The data of 3 days is stored in memory in addition to duplication in which the fault is tolerated. In this process, five terabyte of Ram is consumed [2]. While this work was performed, numerous challenges were faced like in which way the performance can be maximized. In addition to the selected data storage and processing capability, it is important to have an additional memory so that overhead and processing tasks are carried out. In the future work, it is proposed that the CEP framework will be employed which has the capability to process millions of data points/sec. Other methods like Apache storm, Hadoop file systems, and other such techniques will be implemented.

In [3] with the increase of technology, the usage of computers as well as the computation is increased. Moreover, the computation speed at higher rate is also required. For this purpose, either the approach of utilizing the memory processing or the disk-based approach is used. The main demerit of disk-based approach is that the performance is sacrificed. This paper presents a novel approach in which the scaling of graphs into sub graphs is made by RAM-Disk approach. The graphs are divided into sub graphs which are suitable for RAM. So, there is no modification in algorithm and the approach deals with small memory to large memory machines. The approach of processing in terms of out of core is identical to that of paging. However, it is applied to the sub graph which can also be thought as logical partitions. The logical partitions are formed so that they are fit in the main

memory and then combined with asynchronous push model. For the implementation of algorithms, the use of graph-processing engine has been made. It is helpful in executing the algorithms by which accessible resources can be utilized for quick as well as scalable processing of graphs. For the implementation, STAPL frame work is used. The algorithms of STAP GL are also used without any modification. The running time on various platforms having different processing capabilities is illustrated in the form of diagram in the paper. Therefore, the selected platforms include a 2-core tablet having only 1 GB Ram, another PC having 4-core with the memory of 8 GB and supercomputer of CRAY XE6. The running times for various graph mining as well as the graph analytics on PCs having memories of 4 GB and 8 GB RAM are also shown in the paper. This helps in extending the approach to the machines that are based on distributed memory.

In [4], the Big Data figuring is one of the problem areas of the web of things and distributed computing, whose exploration substance are securing, administration, handling, appear, thus on of monstrous information. In the handling join, how to process effectively on the Big Data is the key to enhancing execution. By methods for dispersed registering or memory figuring, numerous organizations and foundations give a few advancements and produce. The Hadoop innovation is run of the mill illustrative of disseminated registering. In Hadoop, MapReduce is a circulated programming system proposed by Google. It is a framework for parallel handling of extensive informational collections. In MapReduce, the undertaking can be part into subtask and the disseminated parallel registering can be acknowledged effortlessly. In memory processing, the Big Data are circulated stacked in various PCs as indicated by some standard. The SAP HANA database can possibly give execution enhancements to existing SAP applications. The circulating memory is registered in which the Hadoop innovation and memory registering innovation mix together. The PC amount of group isn't constrained, so we can stretch out group vastly to store substantial information. Yet, they are invalid in the scene in which there are continuous requests in the low-arrange group. To manage the issue, this paper gives a conveyed registering and memory processing based viable arrangement (Objectification Parallel Computing, OPC). The OPCA is made out of Client Proxy, Protest Manage Server and Object Server. There are three programming in the OPC: Object Manage Server Soft (OMSS), Object Server Soft (OSS), what's more, Client Proxy Soft (CPS). In the arrangement, the information can be organized into protest. At that point, the items are conveyed put away in the PC recollections and parallel process to finish assignments. The OPC is connected to the Electric Asset Quality Supervision Oversee System (EAQSMS) of State Grid of China, the outcome demonstrates that with PCs the framework is proficiently accessible, solid, furthermore, adaptably expansible. There are a few inquiries to investigate and settle, for instance, information pressure, checking the bunch, hot backup, et cetera. These issues require additionally look into. At present, innovation advancement changes rapidly. We require constantly center around new innovation to unravel the issue.

In [5], it is mentioned that in these days, the capacity of information and the administration of interconnected gadgets has turned into an incredible test. In this way, putting away such a lot of information and its calculation requires a unique record system known as MapReduce display which stores, process, oversees and executes the information. Apache Spark is one of the computational frameworks and its apparatuses and related systems are accessible as an open source permit. Several different associations of the world have made a utilization of Hadoop document framework. These computations are very fast in nature. But the main factor which may cause a failure is the storage or memory. Therefore, tuning of memory can be done by the Spark. For this purpose, the excessive knowledge is required. In this paper, the memory selection methods to which the first step is taken is discussed. In this way, the failure of memory and the execution which is 25 percent is decreased. The load of work to be selected is 20 gigabytes of data. The experiment is performed by the use of one thousand and four thousand on 15 gigabytes and 2 gigabytes of memory. This paper deals with the detection as well as the automation [5]. Various graphical representations are shown in the paper such as execution time against the caching strategies, a sort of algorithm which represents the selection process of cache and its configuration. Various kinds of caching modifications have been observed and are experimentally assessed. In this way, the footprints, as well as the working performance, are improved. However, the Spark's strategy for caching is observed to be inefficient in such cases where the data is not adjustable in memory. Therefore, various strategies discussed will result in some tradeoffs. Because of the reduction of footprints, a collection of garbage and its frequency is reduced.

Author of [6], addresses the problem of performing operations on bulk data is quite a time consuming and had bad impact on system performance too. There are three reasons due to which operations on bulk data reduces its latency, bandwidth and effect the system performance and energy efficiency. First, present systems perform bulk data operations on byte/line at a time due to which system has to face performance issues due to high latency. Second, these tasks require a substantial amount of information to be exchanged over the memory channel. Henceforth, they are in a roundabout way that influenced the execution of simultaneously running applications over the memory data transmission. Third, the information transfer across the memory channel speaks to a significant portion of the vitality required to perform operations on bulk data.

In this paper basic goal is too focused on optimizing two classes of bandwidth-intensive memory operations one is bulk data copy which means transferring of data in physical memory, the second is bulk data initialization to reduces the latency, bandwidth, and energy consumed by bulk data operations. Row clone is basically a simple technique to perform operations on bulk data within DRAM. This methodology takes out the need to exchange information over the memory channel to play out a bulk information activity, and subsequently can possibly moderate the related inactivity, data transfer capacity and vitality issues.

Author of [7] tells us the blast of computerized information and the regularly developing requirement for big data investigation has made in-memory big data processing progressively imperative. Due to fast-growing data, processing large graphs in main memory is still a problem. Because of memory bandwidth restrictions, it is difficult to build systems whose performance increases relatively with the increase in graphs. To accomplish a target processing-in-memory (PIM) can be a practical solution. With the end goal to take advantage of another innovation to empower memory – capacity-proportional performance Tesseract is designed for a large scale graph processing. It is made up of a new architecture that completely uses the memory bandwidth. It is a proficient strategy for correspondence between various memory partitions. It also includes two hardware pre-fetcher which works based on hints. By using this strategy author demonstrated that Tesseract works well in both performance and energy effectiveness. Tesseract accomplishes memory capacity proportional performance, which is the main objective in dealing with a large amount of data. This new plan can be a proficient and useful substrate to execute emerging data-intensive applications with memory bandwidth requests.

In [8], it is described as balanced Binary Trees (B Trees) are not optimal for indexing on modern hardware because they cannot better utilize the Cache. These shortcomings have been improved in Adaptive Radix Tress that is also a space efficient. Its lookup performance is very time efficient and supports optimal insertions and deletions as well as deals with the worst-case space consumption, which plagues most radix trees, by adaptively choosing compact and efficient data structures for internal nodes. Even though Adaptive Radix Trees performance is almost similar in time to hash tables but it stores the data in the sorted form, which requires additional computations like range finding and making it in order. But still, ART has performance issue due to a lot of computations. Later on, more work can be done like synchronizing simultaneous updates.

Author of [9] discuss that the development of unstructured information in social media gives an open challenge to cloud database network. In this era, the way we deal with big data processing becomes a serious issue. MapReduce is one of the techniques used for big data analytics. There is a number of indexing techniques like Hadoop++, HAIL, LIAH, and Adaptive Indexing. These techniques are still not efficient or an optimized way to process big data. To solve this problem author proposes a solution that is basically HDFS indexing techniques named as Low-Index and High-Index. The goal of these new approaches is to provide a platform to index in Hadoop Distributed File System and MapReduce frameworks without changing the current Hadoop structure. Low-index gives an index that contains text and it tells the Hadoop to filter only those which contain the related terms. Low-Index likewise upgrades the throughput (limits reaction time) and defeats the problem of long inactive time for list creation. This new approach is better in performance than the Lucene but not efficient in response time. To overcome this problem, High-Index is proposed which is found better than Low-Index in computation and response time. We compare the execution

time of both Lucene and new suggested approaches Low-Index and High-Index. In the beginning, these two took more time for creating an index but after that, they perform better the Lucene approach. In future, we can improve these approaches by working on composite queries in a huge cluster setup.

In [10], the author discusses that for big data top-k queries are a big challenge. Top-k systems depend on positioning functions with the end goal to decide a general score for every one of the items over all the relevant attributes being examined. This positioning capacity is given by the client at query time. Bit-sliced indices (BSI) were proposed to answer these queries proficiently. MapReduce and key-values stores are strategies for investigating huge information; we set up to assess the execution of BSI. Indexing is implemented over Apache Spark for both row and column stores and appeared to beat Hive when running on Map-decrease, and Tez for top-k queries. This methodology is strong and useful for high dimensional information. Top queries are executed over a dataset with 8,000 dimensions. In performing experiments, when expanding the quantity of CPU from 24 to 48 query time reduced by around half while diminishing the number of bit-slices per measurement, the index size and the query time is also reduced. On the bases of this result, we concluded that proposed techniques performed better over Hadoop MapReduce and Hive over tez. For future work, the author intends to explore the correct measure of information rearranging done by different vertical and horizontal partition sizes and how this influences the query time. This can help on deciding the ideal number of bit-slices that ought to be gathered together during map-reduce aggregation. Additionally, it can also be a plan to research more impacts of the BSI attribute section size.

In [11], author tell us that using old strategies like storage of data on disk is now become one of the problems due to growth in data. They don't scale smoothly to address the issues of huge scale Web applications, and disk space limit have far exceeded enhancements in access latency and data transfer capacity. This paper contends for another way to deal with data center capacity called RAM Cloud, where data is kept totally in DRAM. We trust that RAM Clouds can give long-lasting and available storage with 100-1000x the throughput of disk-based systems. RAM Clouds use different techniques like storing the copy of data in DRAM for the fast retrieval of data and to provide durability for data. The two important points in this approach are that they have low latency and they have the ability to combine the resources of a large number of commodity servers. RAM Clouds also have some drawbacks like high cost and high use of energy. Later on, both innovation patterns and application necessities will manage that a bigger and bigger portion of online information is kept in DRAM.

In [12], the author described that Due to the increase in data recently we are facing emerging security and protection challenges. Huge information, since it can dig new learning for monetary development and specialized advancement. New mined data will be unconvincing; while if security isn't all around tended to, individuals might be hesitant to share their information. Since, security has been explored as another

measurement, "veracity," is enormous information. In this article, commit our consideration toward privacy in the big data era. First, formalize the general design of enormous information investigation, distinguish the relating protection prerequisites, and present a productive and security saving cosine computing protocols as an example in response to privacy requirements. There are different existing privacy-preserving techniques. One is privacy-preserving aggregation, second is operations over encrypted data and third is de-identification techniques. To evaluate the proposed PCSC protocol it is compared with direct cosine similarity computation and the HE-based protocol. Based on JAVA language, both used to evaluate results with the same output on the PC with Intel Pentium CPU B980 running at 2.40 GHz, and with 6 Gbytes of RAM. The experiment results show that the proposed PCSC protocol is also efficient along with privacy preserving. Further research can be done by addressing unique privacy issues in big data analytics.

In [13], author addresses because of different difficulties and significant issues it is difficult to manage big data. Map Reduce is the technique for handling the huge amount of information. In this 3-layer traffic, aware clustering algorithm is proposed as the best solution for traffic aware partition and aggregate to minimize the cost of network traffic. One problem that is faced in network traffic is difficult to process data in a given time. The basic goal is to reduce network traffic cost. As a result of applying this technique helps in reducing response time and simulation experiment result revealed that this proposal can reduce the network traffic.

Author in [14] is investigating a lot of healthcare information is as yet a challenge because of the absence of sufficient information management techniques that empower responsive information investigation and examination. A single patient record consists of multiple attributes. Doctor keeps this record for the future prediction about the patient. To resolve this, issue the author compares two in-memory New SQL database technologies Mem SQL and VoltDB. For doing this author uses Medicare claims synthetic data. The goal is to explore data faster and to enable real time predictions. This proposed solutions shoes that most of the queries reply back within 10 seconds. Further, it is planned to continue this research by using different new SQL databases like SQL Fire and to develop a new tool for real-time analysis of data using in-memory database system.

In [15], author address a problem that due to the increase in data it is difficult to use a DRAM for big data processing. As DRAM is facing capacity issues and become a reason of main power drain in modern computers. In this paper, the author proposed an effective page management technique for vast scale NVDIMM memory and gives an effective management technique by keeping in mind both TLB performances and page fault rates. Page fault rates become small with the increase in NVDIMM capacity. NVDIMM-based memory models are getting huge importance due to shifting of the desktop to cloud environment. The goal of this technique is, to be helpful in designing new applications for big data.

In [16], the author mentioned a problem of the large performance gap between processor computation and memory access. To solve the memory wall problem 3D stacked technology is proposed. It is combined with NVM. The proposed solution has some advantages like big capacity low cost and non-volatility. In future NVM materials can be studied along with 3D memory systems.

In [17], the author tells us that because of less expensive and faster processing in DRAM in-memory data management systems have gained a lot of attention. The author proposed an innovative in-memory data management system named as MemepiC. It brings together both online data queries and data analytics functionality permitting low-latency and proficient in-situ information analytics. For conveying message inside the MemepiC RDMA-based communication protocol is designed. Different experiments are performed by to show how efficient MemphiC is, in terms of both storage and data analytics services.

In [18], author tells us that nowadays storing and managing a big data is not only a challenge but also extracting a useful information from that data. The main purpose of this paper is to analyze unstructured data. There are many techniques to solve this problem. MapReduce in connection with the HDFS and HBase database as part of the Apache Hadoop project is a modern approach to evaluate unstructured data. Hadoop methods are used for handling big data and can be enhanced with the right approach.

Author in [19], addresses the problem of gathering a large amount of data. The organizations have to face issues of big data efficient performance and the raised infrastructure cost with the data processing. The new architecture is shifted from centralized to distributed architecture and with the help of these changes organizations are able to defeat with the problem of getting related information from a large amount of data. Apache Hadoop is a proposed technique to solve this problem. The basic goal is to facilitate user and provide them a useful information in less time with least effort.

In [20], while studying big data there are two problems that occurs one is storage of a large amount of data and second is processing speed. To solve these issues in Grid Technology is used. The Main advantages of using this are capacity abilities and the handling power. The Oracle/Cloudera approach is a successful combination of Cloudera's software tool and the Oracle built frameworks intended to give high performance and adaptable information handling for Big Data.

In [21], author said that B-tree or B + -tree is the most famous index structure in disk-based relational database systems, the T- tree has been broadly used as the good approach of index structure for in-memory databases where the entire database resides in the main memory. However, the research work on T-tree doesn't take into account the concurrency control that is a drawback one can say. Similar to B-tree index is B-link tree outperforms the T-tree if concurrency control is reduced. This is because the concurrency control over a T-tree demands more locking than that of a B-link tree, and the overhead of locking and unlocking is high which results in the performance degradation.

In [22], author describe that the data is not big in terms of their volume but also include the queries that are performed to access that data. There are many strategies through which we can search data using different keywords. In this paper new technique is introduced named as ADAM which allows a Boolean retrieval of data for structural and unstructured data. Using this technique queries are made in in the style MapReduce. Signature-based indexing strategy is supported and minimizes the access time. The accuracy and efficiency of this are measured by performing an experiment and form ImageNet 14 million images are retrieved. In the future, it can be tested on more large-scale data consist of multimedia. The author proposed indexing technique which is called hash. This technique takes the bulk of images and gives images that are same as a result.

In [23], STR (Sort-Title-Recursive) algorithm is modified for indexing technique R-trees. It will take spatial data as a data type. After the implementation of this, it is compared with the previous techniques. It is evaluated in terms of space storing index strategy. To improve STR two methods are presented. One is to collect sorted spatial objects and combine them in each axis in the form of slices. In the second strategy, each object has its own axis and then for every object connects them into suboptimum space filling. This improved STR performed well the previous methods. In future Revised R* - the tree can be implemented.

In [24], the author described that as we all know that it is difficult to handle spatial data. In memory, the database needs a method with the use of which we can easily retrieve data and then update. To fulfill this need author proposed new indexing technique named as R-tree. It is an algorithm that helps in updating and searching for data. This algorithm is implemented in Different experiments are performed in the result of which it is concluded that it is useful for spatial data.

In [25], the author presented an Inverted indexing technique to tackle the problem of processing the growing spatial data. This technique is the partition of inverted and grid indexes. This method is implemented with MapReduce many experiments are done to check the scalability of the technique. This technique is constructing time of index is very much less as compared to the other trees.it is three times faster when compared with Voronoi-based query processing.

## III. COMPARATIVE STUDY OF RELATED LITERATURE

After critically analyzing all the literature, his section present the all techniques discussed. There are two types of indexing techniques. One is Artificial Intelligence Approach [23] and second is Non Artificial Approach. Non Artificial Approach is further divided into three categories [20], [21]. One is Tree Based Indexing, second is Inverted Indexing and third is Hashing [24]. Tree is further divided into more categories. Three are shown in Fig. 1, B-tree-tree and X-tree.

Different Indexing technique use different type of data for processing [25]. Every technique has different ways to perform query. All of these factors affect the complexity of indexing technique. The technical summary of indexing technique is given in Table I.

There are many factors that affect the hashing and tree indexing technique. Table II discusses these factors. A comparison of different indexing techniques is shown in Table III and also the analysis of indexing techniques on the basis of big data characteristics is given in Table IV. Furthermore, on comparing all the indexing techniques from Section II, the advantages and disadvantages of these techniques are analyzed and given in Table V.



Fig. 1. Flow Chart.

TABLE I. TECHNICAL SUMMARY

| Technique | Concerns | Related Work |
|---|---|---|
| Indexing | Space Efficiency | Graph Query Tree |
| | Time Efficiency | X-Tree |
| | Better Use of Cache | B-Tree |
| | | R-Tree |
| | | Hashing |
| | | Inverted Indexing |

TABLE II. FACTOR EFFECTING TECHNIQUES

| Factors | Hashing | Tree Indexing |
|---|---|---|
| Access Latency | Reduce Access Time | Minimum Access Time |
| Space | Efficient in Space Handling | High Space Consumption |
| Indexing Efficiency | Efficient in Balancing Access and Tuning Time | Not Powerful |
| Miscellaneous | N/A | Good for Random Access |
| Time | Hashing Function Reduce Stunning Time | Minimal |

TABLE III. COMPARISON OF INDEXING TECHNIQUES

| Indexing Techniques | Data Type | Query Type | Complexity |
|---|---|---|---|
| R-Tree | Multimedia and Spatial Data | 2 to 3-dimensional Access Method | Worst Time Complexity and Inefficient usage of Time |
| B-Tree | Multimedia and Log Data | 1-Dimensional Access and Range Queries | O(logn) |
| X-Tree | Spatial Data | Multi-Dimensional Access and Range Queries | Linear and Time Complexity O(n) |
| Hashing | Multimedia and Log Data | Point Query | N/A |
| Inverted Index Tree | Multimedia Data and Documents | Keyword Queries | N/A |
| Graph Query Tree | Graph | N/A | N/A |

TABLE IV. ANALYSIS BASED ON BIG DATA CHARACTERISTICS

| Indexing Techniques | Volume | Velocity | Variety | Veracity |
|---|---|---|---|---|
| R-Tree | Yes | N/A | No | N/A |
| B-Tree | Yes | N/A | N/A | N/A |
| X-Tree | Yes | N/A | N/A | N/A |
| Hashing | Yes | No | Yes | No |
| Inverted Index Tree | Yes | N/A | No | N/A |
| Graph Query Tree | Yes | Yes | N/A | N/A |

TABLE V. TECHNICAL ADVANTAGES AND DISADVANTAGES

| Indexing Techniques | Advantages | Disadvantages |
|---|---|---|
| R-Tree | • Less query processing cost.<br>• Query response time depends on buffer size | • Index takes more space |
| B-Tree | • Faster Construction<br>• Fast query response<br>• Less Updating Cost<br>• Index takes less time | • Data increase cause increase in construction cost |
| Hashing | • Efficient query response for large dataset | • More initial setup time |
| Inverted Index Tree | • Index takes less space<br>• Fast query response<br>• Manageable query processing cost | • Require more time to load in memory |
| Graph Query Tree | • Less processing query cost<br>• Index takes less space<br>• Fast index construction<br>• Fast query response<br>• Less update cost and scalable for large data | • More Computational cost for large network |

## IV. CONCLUSION AND FUTURE WORK

As the main memory is used as disk, In-memory data management has become interesting for industries. Shifting the data towards main-memory has improved the access time and throughput at a very great extent. Shifting of data has also developed the interest in different aspects to perform optimized results while performing computation on that size of data. Modern Systems has reduced the problem of management of that volume of data by using efficient memory and performance with cache sensitive indexing techniques to better utilize the cache and perform faster Calculations.

As big-data is very vast area of technology, the shortcomings and problems in that field are also at a large scale and all that is in case of management and processing of this size of data. Space efficiency is a factor that should be considered for hash-based indexing techniques. Hashed tree approach can be further improved by working on binary codes to save space and also creating a unique index for each component of data.

REFERENCES

[1] M. Brusakov and G. Botvin, "In-memory technology integration features for work with big data on high-tech enterprises", in Soft Computing and Measurements (SCM), 2017 XX IEEE International Conference on, IEEE, 2017, pp. 697– 698.

[2] J. W. Williams, K. S. Aggour, J. Interrante, J. McHugh, and E. Pool, "Bridging high velocity and high volume industrial big data through distributed inmemory storage & analytics", in Big Data (Big Data), 2014 IEEE International Conference on, IEEE, 2014, pp. 932–941.

[3] N. M. Amato, L. Rauchweger, et al., "Processing big data graphs on memoryrestricted systems", in Proceedings of the 23rd international conference on Parallel architectures and compilation, ACM, 2014, pp. 517–518.

[4] Z. Yang, C. Zhang, M. Hu, and F. Lin, "Opc: A distributed computing and memory computing-based effective solution of big data", in Smart City/SocialCom/SustainCom (SmartCity), 2015 IEEE International Conference on, IEEE, 2015, pp. 50–53.

[5] A. Koliopoulos, P. Yiapanis, T. Tekiner, G. Nenadic, and J. Keane, "Towards automatic memory tuning for in-memory big data analytics in clusters", 2016.

[6] V. Seshadri, Y. Kim, C. Fallin, D. Lee, R. Ausavarungnirun, G. Pekhimenko, Y. Luo, O. Mutlu, P. B. Gibbons, M. A. Kozuch, et al., "Rowclone: Fast and energy-efficient in-dram bulk data copy and initialization", in Proceedings of the 46th Annual IEEE/ACM International Symposium on Microarchitecture, ACM, 2013, pp. 185–197.

[7] J. Ahn, S. Hong, S. Yoo, O. Mutlu, and K. Choi, "A scalable processing-inmemory accelerator for parallel graph processing", ACM SIGARCH Computer Architecture News, vol. 43, no. 3, pp. 105–117, 2016.

[8] V. Leis, A. Kemper, and T. Neumann, "The adaptive radix tree: Artful indexing for main-memory databases", in 2013 IEEE 29th International Conference on Data Engineering (ICDE), IEEE, 2013, pp. 38–49.

[9] A. B. Mathew, P. Pattnaik, and S. M. Kumar, "Efficient information retrieval using lucene, lindex and hindex in hadoop", in Computer Systems and Applications (AICCSA), 2014 IEEE/ACS 11th International Conference on, IEEE, 2014, pp. 333–340.

[10] G. Guzun, J. E. Tosado, and G. Canahuate, "Scalable preference queries for high-dimensional data using map-reduce", in Big Data (Big Data), 2015 IEEE International Conference on, IEEE, 2015, pp. 2243–2252.

[11] J. Ousterhout, P. Agrawal, D. Erickson, C. Kozyrakis, J. Leverich, D. Mazières, S. Mitra, A. Narayanan, G. Parulkar, M. Rosenblum, et al., "The case for ramclouds: Scalable high-performance storage entirely in dram", ACM SIGOPS Operating Systems Review, vol. 43, no. 4, pp. 92–105, 2010.

[12] R. Lu, H. Zhu, X. Liu, J. K. Liu, and J. Shao, "Toward efficient and privacypreserving computing in big data era", IEEE Network, vol. 28, no. 4, pp. 46–50, 2014.

[13] G. Venkatesh and K. Arunesh, "Map reduce for big data processing based on traffic aware partition and aggregation", Cluster Computing, pp. 1–7, 2018.

[14] M. Mian, A. Teredesai, D. Hazel, S. Pokuri, and K. Uppala, "Work in progressin-memory analysis for healthcare big data", in Big Data (BigData Congress), 2014 IEEE International Congress on, IEEE, 2014, pp. 778–779.

[15] S. M. Kwon and H. Bahn, "Efficient memory page management for nvdimmbased big data processing environments", in Information Science and Control Engineering (ICISCE), 2017 4th International Conference on, IEEE, 2017, pp. 283– 287.

[16] C. Qian, L. Huang, P. Xie, N. Xiao, and Z. Wang, "Efficient data management on 3d stacked memory for big data applications", in Design & Test Symposium (IDT), 2015 10th International, IEEE, 2015, pp. 84– 89.

[17] Q. Cai, H. Zhang, W. Guo, G. Chen, B. C. Ooi, K.-L. Tan, and W. F. Wong, "Memepic: Towards a unified in-memory big data management system", IEEE Transactions on Big Data, 2018.

[18] K. Bakshi, "Considerations for big data: Architecture and approach", in Aerospace Conference, 2012 IEEE, IEEE, 2012, pp. 1–7.

[19] J. Nandimath, E. Banerjee, A. Patil, P. Kakade, S. Vaidya, and D. Chaturvedi, "Big data analysis using apache hadoop", in Information Reuse and Integration (IRI), 2013 IEEE 14th International Conference on, IEEE, 2013, pp. 700–703.

[20] D. Garlasu, V. Sandulescu, I. Halcu, G. Neculoiu, O. Grigoriu, M. Marinescu, and V. Marinescu, "A big data implementation based on grid computing", in Roedunet International Conference (RoEduNet), 2013 11th, IEEE, 2013, pp. 1–4.

[21] H. Lu, Y. Y. Ng, and Z. Tian, "T-tree or b-tree: Main memory database index structure revisited", in adc, IEEE, 2000, p. 65.

[22] I. Giangreco, I. Al Kabary, and H. Schuldt, "Adam-a database and information retrieval system for big multimedia collections", in Big Data (BigData Congress), 2014 IEEE International Congress on, IEEE, 2014, pp. 406–413.

[23] B. C. Giao and D. T. Anh, "Improving sort-tile-recusive algorithm for r-tree packing in indexing time series", in Computing & Communication TechnologiesResearch, Innovation, and Vision for the Future (RIVF), 2015 IEEE RIVF International Conference on, IEEE, 2015, pp. 117–122.

[24] A. Guttman, R-trees: A dynamic index structure for spatial searching, 2. ACM, 1984, vol. 14.

[25] C. Ji, T. Dong, Y. Li, Y. Shen, K. Li, W. Qiu, W. Qu, and M. Guo, "Inverted grid-based knn query processing with mapreduce", in 2012 Seventh chinaGrid annual conference, IEEE, 2012, pp. 25–32.

# Comparing Hybrid Tool for Static and Dynamic Object-Oriented Metrics

Babur Hayat Malik[1], Javaria Khalid [2], Hafsa Arif [3], Ayesha Sadiqa[4] ,Amara Tanveer[5], Asia mumtaz[6]

Zartashiya Afzal[7], Samreen Azhar[8], Muhammad Numan Ali[9]

Department of Computer Science and Information Technology
University of Lahore, Chenab Campus, Gujrat Pakistan

*Abstract*—Software metrics are created and used by the distinctive programming associations intended for assessing, guaranteeing program excellence, activity, and software recovery. Software metrics have turned into a basic part of programming growth and are utilized in each period of the product development life cycle. Software metrics essentially measure programming items like plan source code and help us in taking technical and administrative choices. The desire of this examination is to play out the relative investigation of static and dynamic metrics. In any case, software quality characteristics, for example, performance, execution time and dependability rely upon the dynamic exercises of the product artifact. Due to every one of these variables, we favor dynamic metrics instead of customary static metrics. With the assistance of customary static metrics, we are not capable to analyze different actualities of programming. There are various types of this OO static and dynamic equipments. In this paper we have played out a similar investigation of different OO static and dynamic metrics tools and find out the hybrid too is counted as best one extraction of both, static and dynamic characteristics from mobile Android applications. The source code and a Docker compartment is utilized by open source tool in only three phases pre-static, static and dynamic examination.

*Keywords*—*Software metrics; static metrics; dynamic metrics; Object Oriented (OO)*

## I. INTRODUCTION

A software metric is fundamentally a software engineering track which relates to the various software developments and dimensions. One effective tool used for software product analysis is software metrics [1] [2] [3]. It plays a major role in the analysis and improvement of software quality along with measurement of software complexities [4]. An appropriate software model is required for the development of reliable software. ISO 9126 is one of the quality models that uses software metrics [5] [6]. Several tools are required for making of software quality models which intends to do metrics calculations. Though, these tools are also required to produce accurate data [7]. Software metrics are categorized into three parts: product metrics, process metrics, and project metrics, as shown in Fig. 1.

Results are specified by a standard unit known as "Metric". It is used for evaluation of software processes, products, and services. Different authors have proposed several object-oriented (OO) metrics which are quite famous in the present software development environment [9]. These are different from standard metrics as they use objects instead of

algorithms as a key object [10]. Traditional metrics are not eligible in determining the quality as intricate projects are enforced through OOD design practices, so they are required [11]. Somerville [12] described metrics in two types known as static and dynamic. Static metrics analyze code before executing it whereas dynamic metrics analyze code during code execution. In this research, static metrics is more focused on the understanding of procedural and object-oriented programming languages [4]. In this paper comparison of Static and dynamic OO tools are proposed. They are more emphasized for finding object-oriented metric tools on the basis of several parameters.

This paper is written in several sections. Firstly, Section II describes the literature work of various Object-oriented Static and dynamic metrics tools. Then, in Sections III is discussed the differentiation between Static and Dynamic Metrics. Various types of object-oriented Static and dynamic Metrics are presented in Section IV. In Section V, the comparative study of OO Metric Tools is performed. Lastly, Section VI, presents the conclusion of this article.



Fig. 1. Software Metrics [8]

## II. LITERATURE REVIEW

Various OO metrics are developed until now which differ in their properties and features. The main purpose of this paper is to find out huge OO metric computational tools on the basis of their properties. Complex metrics to be resolved are still an issue whereas in traditional OO some metrics like CK and MOOD are considered quite helpful in the development of software [13].

Munson and Hall [14] identified the program complexity level along with three processes of functional, fractional, and operational complexity. Mayo et al. [15] discussed the quality attribute of the interface which calculates modules complexity and dynamic metrics when it's executed.

Honglei et al. [16] presented metrics definition, types, and history. Measurement of software complexity is one important factor and it's also related to software development price factor.

Hassoun et al. [17] proposed Dynamic Coupling Metric (DCM) for object level coupling that considers program execution as it is used to measure objects coupling during runtime. Though it also estimates the runtime complexity and system comparison at meta-level along with those systems which have no reflective features.

Singh and Singh [18] presented four class-level dynamic couplings for identifying object-oriented systems quality. They are more determined in finding key coupled classes consisting of most active classes during runtime. Gupta [19] presented three dynamic coupling metrics which consists of foremost relations between objects during runtime, i.e. aggregation, inheritance, etc.

Mayo et al. [20] defined both automated Interface and Dynamic Metrics. The first one is used for identifying modules complexity whereas dynamic metric calculates quality factor during execution. Hays in [21] identified OO systems testing and compared them with conventional programming language testing.

Mohsin, Shaikh, and Zeeshan Kaleem [22] presented the idea of code comprehension with a combination of Software metrics and techniques called Program Slicing. It is basically coded automation analysis for coupling, cohesion, and complexity.

Debbarma, Mrinal Kanti et al. [23] described the comparison of static and dynamic metrics and analyzed them in terms of regression testing that helps in effort and time estimation used during testing.

## III. TYPES OF METRICS

In various real-time applications, there is a small number of the most eminent metrics that are analyzed. There are different categories of metrics that are presented below:

### A. Traditional Metrics

In an object-oriented system, traditional metrics are commonly applied to the methods that include the class operation. "A method is a component of an object that operates on data in response to a message and is defined as part of the declaration of a class". Methods reveal how a problem is fragmented into different sections. Two traditional metrics are Cyclomatic complexity and size (line counts) [24].

### B. Object-Oriented Metrics

Object-oriented software metrics emphasis on measurements that are functional to the conceptions of classes, coupling, and inheritance. Encapsulation metrics are applied for classes, not for modules. Information Hiding is measured & enhanced due to Inheritance complexity is additional, the level of abstraction can be measured by Object Abstraction metrics. These are as follows:

- Metrics correlated with Class
- Metrics associated with Methods
- Metrics Encapsulation
- Measurement of Cyclomatic complication
- Metrics used for Inheritance [25, 26].

*1) Static metrics:* This Metric is the outcome of non-executable code. Static metrics describe system features from design through maintenance. Earliest Metric used for Static is [27] (LOC/KLOC) examine the throughput of a software package. In earlier 1990, McCabe was the most powerful metric for examining the intricacy of cyclomatic [28] complexity. Complexity is evaluated from the graphical representation and various mathematical equalities. In 1976 McCabe [29] demarcated the cyclomatic complexity metric. It measures the total numbers of independent routes over a software component.

*2) Dynamic metrics:* These are resultant of source code investigation. When code is running it evaluates what is really happening. Dynamic metrics comprise complication events and processes beneficial in consistency demonstrating at the same time [30]. When software is executing its values are reliant on the involvement or experimental information. From coding to maintenance system aspects are classified by dynamic metrics [8]. The comparison of static and dynamic metrics with its merits, demerits are shown in Tables I and II.

TABLE I.    STATIC VS. DYNAMIC METRICS

| *Static Metrics* | *Dynamic Metrics* |
|---|---|
| 1. Its nature is always static. | 1. Its nature is always dynamic. |
| 2. It is simpler and easier to collect. | 2. It is difficult and tough to gather. |
| 3. OO software attributes are difficult to examine. | 3. Different characteristics are easy to inspect like Inheritance, polymorphism, coupling, cohesion, and difficulty. |
| 4. It takes less time as compared to dynamic analysis of software. | 4. It takes more time to perform dynamic analysis of a program. |
| 5. It is available at the early stages of the software development life cycle. | 5. It is accessible late in the software development life cycle. |
| 6. For software quality prediction its results are less accurate. | 6. For software quality prediction its results are more accurate. |
| 7. More Tools are effortlessly available to accomplish this examination. | 7. Only a few tools are available for this analysis. |
| 8. Its implementation is done on the code. | 8. Its implementation is performed while code is being run. |
| 9. It deals with structural aspects of the system. | 9. It deals with the behavioral aspects of the system. |
| 10. It identifies vulnerabilities in a runtime environment. | 10. It can find weaknesses in the code at the exact location. |

TABLE II.    COMPARISON OF STATIC VS. DYNAMIC METRICS

| Serial No. | Static Software Metric | Description | Merits | Demerits | Equations |
|---|---|---|---|---|---|
| 1 | SLOC (Source lines of code ) [4] | It evaluates total lines in the program to measures its size. When software is developed it determines the productivity of the program. | Measuring automation possibilities | Inaccuracy in Accountability. | For (i = 0; i < 100; i++) printf("hello"); /* How many lines of code is this? */ Above case illustrate the following information: • 1(LOC), • 2(SLOC) (for statement and printf statement), • 1 comment line. |
| 2 | LOC (Line of Code)[4] | It consists of any number of lines, consist of source, whitespace, and comments. | Universal measure. | • Several languages issues • GUI tools Starter | 1 (LOC) as stated in the above example |
| 3 | AMLOC(Average lines per method) [32] | It defines the average size of the method. | Method Size can detect simply. | Less clear and additional code statement. | Average Method Size= (The Total number of LOC) / (Number of Methods) |
| 4 | ACLOC(Average lines per class) [32] | It determines the moderate size of class according to LOC. | It is simple to define the number of code lines for each class therefore accurately determine the size of the class | More code lines can't be verified and can't be altered safely. | Average Class Size= (The Total number of LOC) / (Number of Methods) |
| 5 | NCLASS [32] | These metrics calculate the number of classes in the project. | Main Characteristics are undone or round-trip engineering. | In general UML figure categories, it supports class diagram. | ---- |
| 6 | Cyclomatic complexity [33] | Indicate the program difficulty areas. | • It assesses AI semantic complexity. • Useful in geographical and landscape environmental inquiry. | Positive correlation among cyclomatic complexity and defects. More errors in maximum complexity functions and methods. | $M = E - N + 2P$, $E$ = Graph edges. $N$ = Graph nodes. $P$ =Connected components. |
| 7 | Function point [34] | It is a measurement element to examine the business functionality that delivers to a customer. | • An end-user business function maps to functional consumer requests like data entry. • Function points plot easily into user-oriented requests. | lbrecht perceived in his research that Function Points were extremely associated with code lines and increase complexity. | • Define the number of data functions (ILFs and EIFs) • indicative size (fp) = 35 x number of ILFs + 15 x number of EIFs. |
| 8 | Bug Counting [34] | Program inaccuracy results in improper or unpredicted result act in unintentional ways. | • Failure count models • Error seeding models | • Involved more in program performance, does not concentrate on a number of program bugs. • Most requests of customers define according to functional reliability and not in terms of errors. | Bugs.Count Bugs.SUM(Effort) Bugs.SUM(CustomValues.Number("Cost")) UserStories.SUM(CustomValues.Number("Bugs Count")) + Bugs.Where(UserStory.Feature == null \|\| Feature.Id != UserStory.Feature.Id).Count |
| 9 | Halstead complexity [4][34] | Recognize computable software properties and the associations between them. | • These are traditional metrics but they can evaluate projects like C, C++, and Java. • It calculates the bugs, project length, size, and validity period. | • Modularity • All-Depth • Operator Type • Database Impact and Declaration | Program Vocabulary: N=$n_1$ + $n_2$ Program Length: N=$N_1$+$N_2$ Calculated Program Length: N= $n_1$ log2 $n_1$+ $n_2$log2 $n_2$ |
| 10 | Continuous Value Metrics[24] | In numerous circumstances it innate incorrectness: A straight line in a diagram can have the equivalent general average as a slanting line. | • Define a best, fewer bugs metric where single value metrics is possibly imprecise. • Value Metrics extension | Secondary metrics are frequently insufficient to actually define the dissimilarity in performance, demanding further tertiary metrics. | ---- |

## IV. CURRENT ISSUES AND CHALLENGES

- After negotiations upon dynamic metrics, it has definitely perceived that currently not any metrics available for testability at execution time of the software systems.

- Its benefit includes accuracy and precision; however, they are more difficult in evaluation to static ones. Therefore, a good hybrid approach is required.

- For the analysis of different software aspects pseudo dynamic metrics is another auspicious research prospect readily accessible to researchers.

- It can be certainly observed from the survey of many research studies conducted by different authors that dynamic metrics are examined and tested using a project that is not bulky [31].

## V. RESULTS

We have to concern together static along with Dynamic Metrics to realize the deviation. After comparing both of these metrics we concluded that dynamic metrics analysis gives result at execution time of programs whereas static analysis at rest of the SDLC process. So, for dynamic analysis data is collected with the help of tool based on either Java or C++ based application, then apply a statistical tool to measure the quality of the product. Dynamic analysis can give a better result than static analysis.

AndroPyTool executes different tools in order to extract wide-ranging features from an input set of Android samples. All these features and the evidence that they symbolized are organized in three dissimilar classes (pre-static, static and dynamic), both the features and how they are extracted.

In Pre Static it comprises extracting information without inspection of code and permits to categorize and to track the sample. It also includes the package name and the main activity name, which are found with Andro-guard. In Static analysis, it contains those features that are regained by analyzing the application at the code level. In this category, features such as API calls, activities, opcodes or permissions can be originated. In Dynamic Analysis, it includes Droid Box tool for this purpose, which allows to dynamically find dissimilar information in real time. The information gathered by the Droid Box tool includes: the use of cryptographic functions, loaded DEX classes in run time and the kind of operation, network connections, SMS, phone calls, started services, enforced permissions and information leaks detected. The detail diagram of AndroPyTool is shown in Fig. 2 [38].

## VI. DISCUSSION

Various OO Metrics tools their description, merits and demerits are studied in this research paper. These tools are tabulated under various attributes that would be of interest to developers and researchers using the tools as elaborated in Table III. Our study has further pointed out the work and research findings that has been done till now to use of hybrid approach of static as well as dynamic metrics, although they have tremendous scope. Based on the analysis of existing dynamic metrics, we have tried to reveal potential research challenges and opportunities existing in the field of dynamic metrics. Best methodology that is suitable for pre-static, static and dynamic metrics is hybrid approach and its tool that is AndroPyTool.



Fig. 2. AndroPyTool [38].

TABLE III.     COMPARISON OF STATIC VS. DYNAMIC METRICS TOOLS

| Tool Name | Description | Language | Availability | Authors | Tool Type |
|---|---|---|---|---|---|
| CheckStyle [35] | Java Checkstyle is an improvement tool to enable designers to compose Java code that clings to a coding standard. Presently Checkstyle gives checks that discover class plan issues, copy code, or bug designs like twofold checked to bolt. | Java | Open source | Oliver Burn | Static |
| FindBugs [35] | This is the Static Analysis tool and is open source that checks and study class files or JAR libraries for probable problems adjacent to a list of bug patterns by matching the byte code [5]. | Java | Open source | David Hovemeyer and William Pugh | Static |
| StyleCop [35] | For plugins and customs rules, *StyleCop* provides an extensible framework to write down custom rules which match up to our requirements. | C# | Free | Andy Reeves, Chris Dahlberg | Static |
| JMT [36] | It only associates the Metrics with Java language. | Java | Free | Politecnico di Milano and Imperial College London | Dynamic |
| QMOOD++ [37] | QMOOD++ is easy and free of cost accessible in the runnable application and source code form. It handles the 30+ Metrics. QMOOD++ is an inclusive, multi-handler, multiprocessing, incorporated software tool. | C++ | Free | Bansiya, Jagdish, and Carl Davi, | Dynamic |
| JMetric [11] | JMetric only works with Java. Its information is presented through tables and charts. | Java | Free | Commercial Tool | Dynamic |
| AndroPyTool [38][39] | AndroPyTool incorporate different analysis tools and Android applications Processing tools, in order to convey fine-grained reports drawing their individual performance and features. | Python | ---- | Alejandro Mart, Raul Lara-Cabrera, David Camacho | Hybrid |

## VII. CONCLUSION

A correlation of diverse software metrics and its major tools are presented in this comparative study. On the base of their major types like static and dynamic metrics, these are differentiated. At early stages of software development life cycle (SDLC), Static metrics are reachable easily. These metrics manage the overall structural qualities of the product framework and very simple to assemble. The unpredictability of static metrics has calculated the measure of exertion expected to create and keep up the code. In the latter stage of the software development life cycle, dynamic metrics are easily reachable.

These metrics confine the dynamic conduct of the framework and difficult to acquire and got from hints of code. After a virtual study of various static and dynamic tools are performed and broke down that hybrid tool is best in the greater part of the android applications. AndroPyTool, the primary objective is to furnish scientists and malware examiners with an incredible and coordinated device for extracting multi-source highlights from Android applications. In future work, more tools and features can be add on into AndroPyTool tool for better analysis and to improve the data analysis stages, in order to give more functionalities to the users.

REFERENCES

[1] M. Sharma, Dr. G. Singh, "Analysis of Static and Dynamic Metrics for Productivity and Time Complexity," IJCA, vol. 30, issue.31, September 2011.

[2] H. F Li and W. K Cheung, "An Empirical Study of Software Metrics," Software Engineering IEEE Transactions,vol.13, issue. 6, pp. 697-708, 1987.

[3] N. E Fenton "Software Metrics," Conference Proceedings of on the future of Software engineering ICSE vol. 8, issue: 2,2000.

[4] M. Sharma, A. Bhardwaj, L. Singh, N.Singh and C. Sharma, "Comparative study of static metrics of procedural and object oriented programming languages," International Journal of Computers & Technology, Volume 2 No.1 February 2012.

[5] ISO ISO/IEC 9126-1, Software engineering–Product Quality - Part 1: Quality model., 2001.

[6] ISO. ISO/IEC 9126-3, Software engineering–Product Quality - Part 3: Internal metrics., 2003.

[7] J. Novak and G. Rakić, "Comparison of software metrics tools for .net," University of Novi Sad, Faculty of Sciences, Department of Mathematics and Informatics, 2011.

[8] BM. Goel and S. Bal Gupta, "A Comparative Study of Static and Dynamic Object Oriented Metrics," International Journal of Information Technology & Systems, vol. 5, issue. 1, 2016.

[9] J. Chawla and A. Agarwal, "Object-Oriented Design Metrics to Predict Fault Proneness of Software Applications," (IJCSIT) International Journal of Computer Science and Information Technologies, vol. 5 (3), 2014.

[10] A. Albrecht and J. Gaffney: Software Function, Source Lines of Code, and Development Effort Prediction: A Software Science Validation; in IEEE Trans. Software Eng., pp. 639-648,2008.

[11] Kayarvizhy N, "Systematic Review of Object Oriented Metric Tools," International Journal of Computer Applications vol. 135 issue.2, February 2016.

[12] Somerville "Software Engineering", 6th Edition, Editor: Addison Wesley.

[13] Y. Ma, K. He, D. Du, J. Liu, and Y. Yan , "A Complexity Metrics Set for Large-scale Object-oriented Software Systems," IEEE International Conference on Computer and Information Technology (CIT'06).

[14] John C. Munson and Gregory A. Hall, "Estimating test effectiveness with dynamic complexity measurement," Empirical Software Engineering Journal.

[15] Kevin A. Mayo, Steven A. Wake and Sallie M. Henry, " Static and Dynamic Software Quality Metric Tools," Department of computer Science, Virginia Tech, Blacksburg.

[16] B. Mohan Goel and S. Bal Gupta, "Dynamic Coupling Based Performance Analysis of Object Oriented Systems," International Journal of Advanced Research in Computer Science, vol. 8, issue. 5, May-June 2017.

[17] Y. Hassoun, R. Johnson and S. Counsell, "A Dynamic Runtime CouplingMetric for Meta Level Architectures," In Proceedings of Eighth EuromicroWorking Conference on Software Maintenance and Reengineering, pp. 339, 2004.

[18] P. Singh, H. Singh, "Class-level Dynamic Coupling Metrics for Static and Dynamic Analysis of Object-Oriented Systems," International Journal of Information and Telecommunication Technology, pp. 16-28, 2010.

[19] V. Gupta, "Validation of Dynamic Coupling Metrics for Object-Oriented Software." ACM SIGSOFT Software Engineering Notes,vol.36(5), 2011.

[20] Kevin A. Mayo, Steven A. Wake, Sallie M. Henry, "Static and Dynamic Software Quality Metric Tools," Department of computer Science, Virginia Tech, Blacksburg.

[21] J. Huffman Hayes, "Testing of Object-Oriented Programming Systems (OOPS): A Fault-Based-Approach," Science Applications International Corporation, 1213 Jefferson-Davis Highway, Suite 1300, 22202 Arlington, Virginia.

[22] M. Shaikh, and Z. Kaleem. "Program Slicing Based Software Metrics towards Code Restructuring," In Computer Research and Development, Second International Conference on, pp. 738-741. IEEE, 2010.

[23] Debbarma, M. Kanti, N. Kar, and A. Saha, "Static and dynamic software metrics complexity analysis in regression testing," In Computer Communication and Informatics, International Conference on, pp. 1-6. IEEE, 2012.

[24] S. Pasupathy and R. Bhavani, "Object Oriented Metrics Evaluation," International Journal of Computer Applications (0975 – 8887) vol.78,issue.1, September 2013.

[25] S. Morasca, "Software Measurement: State of the Art and Related Issues, slides from the School of the Italian Group of Informatics Engineering," Rovereto, Italy, September 2008.

[26] J. Alghamdi, R. Rufai, and S. Khan, "Oometer: A software quality assurance tool. Software Maintenance and Reengineering 2009," 9th European Conference on, pp. 190, March 2010.

[27] Li, Cheung, W.K, "An Experimental investigation of software metric and their relationship to software development effort," IEEE Transaction on software engineering 649-653, Piscataway, NJ, USA.

[28] Thomas J. McCabe, "A Complexity Measure, IEEE Transaction on Software Engineering," vol.2 issue. 4, pp. 308-320.

[29] S. Singh, K.S. Kahlon, "Static Analysis to Model & Measure OO Paradigms," SAC, ACM.

[30] K. Kaur, K. Minhas, N. Mehan, and N. Kakkar, "Static and Dynamic Complexity Analysis of Software Metrics," World Academy of Science, Engineering and Technology International Journal of Computer and Systems Engineering vol.3, issue.8, 2009.

[31] Chhabra JK, Gupta V, "A survey of dynamic software metrics. JOURNAL OF COMPUTER SCIENCE AND TECHNOLOGY," vol.25(5),pp.1016–1029 Sept. 2010.

[32] G. Singh, M. Sharma, "A Comparative Study of Static Object Oriented Metrics," International Journal of Advancements in Technology, 21 January 2018.

[33] G. K. Gill, C. F. Kemerer, "Cyclomatic Complexity Density and Software Maintenance Productivity", IEEE Transactions on Software Engineering, 1981, pp. 1284-1288.

[34] A. Versa, Rahul, "A Study of Various Static and Dynamic Metrics for Open Source Software," International Journal of Computer Applications (0975 – 8887) vol. 122, 10 July 2015

[35] J. Novak, A. Krajnc and R. Zontar, "Taxonomy of Static Code Analysis Tools," 16 March 2015.

[36] M. Bertoli G. Casale and G. Serazzi, "JMT: Performance engineering tools for system modelling," Giuliano Casale, 04 June 2014.

[37] J. Bansiya, C. Davis, Using QMOOD++ for object-oriented metrics, Dr. Dobb's Journal , 1997.

[38] A. Martin, R. lara-cabrera and D. Camacho, "A new tool for static and dynamic Android malware analysis," 24 September 2018.

[39] A. Martin, R. LaraCabrera and  D. Camacho, "Android malware detection through hybrid features fusion and ensemble classifers:the AndroPyTool framework and the OmniDroid dataset," 05 February 2019.

# Comparison of Agile Method and Scrum Method with Software Quality Affecting Factors

Muhammad Asaad Subih[1], Babur Hayat Malik[2], Imran Mazhar[3], Izaz-ul-Hassan[4], Usman Sabir[5]
Tamoor Wakeel[6], Wajid Ali[7], Amina Yousaf[8], Bilal-bin-Ijaz[9], Hadiqa Nawaz[10], Muhammad Suleman[11]
University of Lahore, Gujrat Campus, Gujrat, Pakistan

*Abstract*—The software industry used software development lifecycle (SDLC) to design, develop, produce high quality, reliable and cost-effective software products. To develop an application, project team used some methodology which may include artifacts and pre-defining specific deliverables. There are different SDLC process models such as waterfall, iterative, spiral and agile model available to develop a quality product. In this paper we focus only on agile software development model, and Scrum model and their techniques. There are many papers and books written on agile methodologies. We will also use their knowledge in this paper. To collect data for comparison of agile method with software quality affecting factors, an online questionnaire survey was conducted. The survey sample consisted of software developers with several years of industry experience using agile methodologies. The main purpose of this study is to compare soft-ware quality affecting factors with agile and scrum model.

*Keywords—Component; SDLC; Software Quality Affecting Factors; Agile methodologies; Scrum*

## I. INTRODUCTION

Agile methodologies have played a vital role in the development of software as compared to other methodologies. Because many companies want to implement good-quality systems, and they want to do it in a minimum period of time, at a less cost [1][2][5]. Therefore, many companies have started to follow agile methods to develop software [1], and it has been found that the extent of the organizational team's skills [1], culture of the organization [1], nature of project [1], and project constraints must be given in-depth consideration [1] when selecting an Agile method. Agile Software Development Methodologies [1][2][3] are based on both incremental and iterative development.

There are different agile methodologies such as scrum, Kanban, Extreme programming, Dynamic System Development, Feature Driven Development, [1][2][3][5]. As our research aim is focus only on agile software development model, and Scrum model and their techniques. The main purpose of this study is to compare software quality affecting factors with agile and scrum model [1]. Nowadays, many organizations are using agile methodologies because these methodologies support flexibility, changes at any stage, and light documentation.

Scrum [1][5][6] is an agile framework focuses software development. This method is based upon consistent installments, and regular collaboration among self-organizing cross-functional departments. This method is for 3 to 9

members team who break their work into actions that can be completed within time boxed iterations, called "sprints" [1]. The sprints are not more than one month and most commonly two weeks, re-plan in 15min standup meetings, called "daily scrums". Scrum rules are product owner, scrum master and team. Scrum is easy with changes; it accommodates with changes [1]. Agile methodologies address quality issues and focusing on the role of quality.

## II. RELATED WORK

It is important to work on product quality and to measuring software development projects to get a better understanding about the progress of project [2]. Agile methodologies address quality issues repeatedly and continuously. Agile method focuses on software quality through customers [2], minimum error rate, faster development and welcome to changing environment.

As Oberscheven [1] research aim is focus on agile software development model, and Scrum model and their techniques. The main purpose of this study is to compare software quality affecting factors with agile and scrum model. Nowadays, many organizations are using agile methodologies because these methodologies support flexibility, changes at any stage, and light documentation. There are different agile methodologies such as Scrum, Kanban, Extreme Programming, Crystal, Dynamic System Development Method, Feature Driven Development, and Adaptive Software Development. A Hossain, MA Kashem: [2] describe that the agile techniques are used to minimize risk factors by developing software in short period of time. During the software development the changes are welcome in agile method. In the end, this paper concludes that agile techniques are using to increase the quality of software all the way through increased customer value. It de-scribes how we can increase the quality of software by using agile techniques.

Sheetal Sharma [4] describe the factors that affect the quality of software [4] are correctness, reliability, usability, extensibility, reusability, testability, portability, maintainability, and efficiency. This research paper is about SDLC models and different scenarios which are using by developers for developing a well-engineered software. Sheetal Sharma describes some advantages and disadvantages of agile methodologies in our research paper. According to this, there are least documentation in agile method and ensure the customer satisfaction but meanwhile least documentation is also disadvantage of agile methodology.

## III. Software Qualty Affecting Factors

Software quality is related to customer satisfaction and low error levels in software. The external and internal quality criterion are used to evaluate the quality of a software. External quality is relating to the functionality of the software. Internal quality relates to coding, and that are not visible to the end-user [2][6][8]. These Software Quality factors are shortlisted by the brief literature study and are supported by the literature review as in [2][6][8]. The factors under the study are supported by above mentioned literature.

Generally, the software qualities are of three types:

- Quality of design
- Quality of performance
- Quality of Adoption

These three qualities can be further divided into other quality attributes. Table I shows the factors affecting the software quality and their attributes.

TABLE I. Shows the Software Quality Affecting Factors

| Software Quality Factors | |
| --- | --- |
| Quality of Design | **Description** |
| Correctness | If software doesn't work correctly as required then it is wasteful. |
| Maintainability | If software is not able to add new features or to remove error, then it has no worth. |
| Quality of Performance | |
| Efficiency | It is a factor relating to all issues in the execution of software. |
| Reliability | It defines how well the software meets his requirements. |
| Usability | If software is not user friendly then it is hard for user to use. |
| Testability | If testing is not done properly then it makes software errors. |
| Quality of Adoption | |
| Extensibility | If software coding is not extendable by adding new features then it has no worth. |
| Portability | It is the effort required to transfer the software from one configuration to another. |
| Reusability | If software is not reusable then it is limited product. |

## IV. Some Popular Agile Methodologies

After All of the agile methodologies acknowledged to produce higher quality software and more significantly the satisfaction of customer [1][7]. Here, some of agile methodologies are listed below.

### A. Extreme Programming (XP)

XP [3][4][6][7] is one of the most successful agile method. This method focuses on customer satisfaction. It divides the SDLC into various number of short development cycles. At any phase of SDLC, it allows changes or requirements from the customers. First phase of Extreme programming is collecting user requirements, and then these requirements are divided into various small no of cycles. Now the upcoming phase is iteration planning. If any new user requirement may come during the development phase then according to that, iteration plan should be tuned [4]. Next step is testing and errors will be removed in the next iteration.

### B. Feature Driven Development (FDD)

FDD [3][6][7] is one of the agile development methods. Designing of the domain of software is the main feature of this method and it also focus on the building phases of the software. The first phase of the method is to get user requirements and constructing the overall model of the project. Next phase is to list down features relates to the user-valued functions [3] [6]. For example, 'calculation of company's each employee', 'calculation of tax each company's employee'. There are different groups of features are made based on their domains i.e. related features are combined into a single group. In another step make a plan for developing and assigned these tasks to development team.

### C. Scrum Method

Scrum [1][5][6] is an agile framework focuses software development. This method is based upon consistent installments, and regular collaboration among self-organizing cross-functional departments. This method is for 3 to 9 members team who break their work into actions that can be completed within time boxed iterations, called "sprints" [1]. The sprints are not more than one month and most commonly two weeks, re-plan in 15min standup meetings, called "daily scrums". Scrum rules are product owner, scrum master and team. Scrum is easy with changes; it accommodates with changes. Scrum [1][5][6] is a simple framework used to organize teams and get work done more productively and with higher quality. Scrum is easy with changes; it accommodates with changes. Some key scrum practices are discussed below [1][3][4][5]. These key factors are taken in account for description one by one.

- Product Backlog – The software development team identified all tasks and makes a list called the Backlog.

- Sprints – Sprint is 3 to 9 members team who break their work into actions that can be completed within time boxed iterations, called "sprints" [1]. The sprints are not more than one month and most commonly two weeks.

- Sprint planning meeting – There are different stakeholders are involved in sprint planning meeting. They decide the functionality of the system. The stakeholders are customers, product owner and scrum team.

- Sprint Backlog – When a list of tasks is completed than a new iteration of the software product is delivered.

- Daily Scrum – These daily scrum meetings are not more than 15minutes long [9].

## V. Evaluation of Software Qualities with Agile Techniques

### A. System Metaphor

The system metaphor is all about that how system works. In system metaphor customers, programmers, and managers are involved. The System metaphor [2] is used to facilitate communication between customer and developer. It helps the agile development team [2] in the development of software by in-creasing communication between developers and users. So, by using system metaphor maintainability, efficiency, reliability and flexibility of the system enhance.

### B. Architectural Spike

An architectural spike is technique that can re-duce the technical risk factor from the product and it comes from Extreme Programming (XP). The main purpose of this technique is to reduce the risk of a technical problem [2]. The spike is between product owner and development team.

### C. Onsite Customer Feedbacks

Customer are only involved while the development of software. In agile methodologies, communication with the customers are required, which is in-tended to improve productivity [2]. Agile methodologies emphasize a lot on customer feedback.

### D. Refactoring

During the development of software, we can improve the internal structure of the software without effect on external behavior. So, by using refactoring, efficiency, reliability, intra-operability and interoperability [2], testability of the system enhances.

### E. Pair Programming

It is a technique in which two programmers are involved on same code. One programmer writes code while the other monitor the code and gives reviews. So, by using Pair programming correctness, verifiability [2], testability of the system enhances and reduces defects.

### F. Stand-Up-Meeting

This method is for 3 to 9 members team who break their work into actions that can be completed within time boxed iterations, called "sprints" [1]. The sprints are not more than one month and most commonly two weeks, re-plan in 15min standup meetings, called "daily scrums". Stand-up-meeting is very useful for improving the quality of the software like reliability and flexibility.

### G. Continuous Integration (CI)

The programmers are sharing the code in CI technique [2]. In continuous integration like the auto-mated compilation, unit test execution, and source control integration are configure by agile teams.

Table II shows the software qualities in the agile development. The software qualities contain Maintainability, Verifiability, Efficiency, Integrity, Reliability, Usability, Testability, Expandability, Flexibility, Portability, Reusability, Interoperability and interoperability. Table II further describes their use in the agile development technique.

TABLE II.     Showing the Software Qualities in Agile Development Team

| Software Qualities | Agile Technique |
|---|---|
| Maintainability | Continuous Integration |
| Verifiability | Continuous Integration |
| Efficiency | Pair programming System metaphor |
| Integrity | Continuous Integration |
| Reliability | Refactoring System metaphor |
| Usability | Continuous Integration |
| Testability | Pair programming Acceptance testing Unit testing Refactoring |
| Expandability | Continuous Integration Onsite customer feedback |
| Flexibility | Stand-up meeting System metaphor |
| Portability | Pair programming Acceptance testing |
| Reusability | Refactoring Continuous Integration |
| Interoperability | Refactoring System metaphor |
| Intra-operability | Continuous Integration |

## VI. Questionnaires

The following is a questionnaire developed to compare the software quality affecting factors with agile methodologies. The data on agile methodology was collected with the help of questionnaire survey. A summary of the survey and their results are given below.

### A. Questionnaire Format

The questions were divided into three parts: First part deals with the respondent's position in the organization and experience on agile methodologies. Second section related to the agile methodologies and scrum method. Last section deals with the software quality affecting factors of agile methodologies.

### B. Questions

Following the questionnaire survey approach, we formulated questions which help us to assess goals. The defined questions are presented in this section.

Table III shows the questions being asked and their relation to the agile development and software qualities. Table III is divided into three main sections.

First section is about the personal information of the respondent, his positions in the organization and his experience about the agile development. Second section queries about the agile methodology and the Scrum development technique and the last most and the third section inquiries about the software quality affecting factors with respect to agile methodology.

TABLE III.    SHOWS THE QUESTIONS RELATING TO AGILE METHOD AND SOFTWARE QUALITIES

| Questions | |
|---|---|
| Respondent's position in the organization and experience on agile methodologies. | Q1.1 What is your name? ___<br>Q1.2 What is your company name? _<br>Q1.3 What is your job title? __<br>Q1.4 How long have you Experiences with Agile Methods?<br>• <1-year<br>• 1-2 years<br>• 3-5 years. |
| Agile methodologies and Scrum method | Q2.1 What are the key reasons to start agile?<br>• Quality<br>• Productivity<br>• Predictability<br>• Team health<br>Q2.2 Starting Scrum Meeting Practice Difficulty?<br>• Easy<br>• Hard<br>Q2.3 What is the frequency of using Scrum meeting?<br>• Never<br>• Irregularly<br>• Frequently<br>• Daily<br>Q2.4 Frequency of Using agile methodologies?<br>• Never<br>• Rarely<br>• Frequently<br>• Always |
| Software quality affecting factors with agile methodologies | Q3. By using agile methodologies for the development of software, then the software has following qualities: give answer of every option as "agree, disagree, strongly agree, strongly disagree".<br>• Software does work correctly as required.<br>• Software is error free.<br>• Software is user friendly for use.<br>• Software coding is extend-able for adding new features.<br>• Software can be reused in other related applications.<br>• Software has testability quality.<br>• Software can be transfer from one configuration to another.<br>• Software is able to add new features or to remove error.<br>• Software is efficient. |

## C. Questionnaire Distribution

The questionnaire was distributed through e-mail. Due to shortage of time, questionnaire was distributed among 15 different software developers. Then 10 of the software professionals responded within 2 days. One of the main reasons to select these persons which are working in soft-ware development companies was, easy to access and easy communication.

## D. Questionnaire Results

First questionnaire section deals with the respondent's position in the organization and experience on agile methodologies. Fig. 1 represent that 28% respondents have longer experience in agile, 27% of the respondents have 3 to 5 years' experience and only 18% respondents have less than 1-year experience.



Fig. 1.    Shows Agile Methods and Software Quality Factors Questionnaire Results Experiences with Agile Methods.

Second section related to the agile method and scrum method (see Fig. 2, 3, 4 and 5). Fig. 2 shows that what are the key reasons that developers are using agile methodologies. Graph show that 10 developers are using agile methodologies because of quality, 6 developers are using agile method due to team health and five dues to productivity.

Fig. 3 show that 10 respondents state that starting scrum meeting practice is not difficult.



Fig. 2.    Agile Methods and Software Quality Factors Questionnaire Results– Key Reasons to Start Agile methodologies.



Fig. 3.    Shows Agile Methods and Software Quality Factors Questionnaire Results–Starting Scrum Meeting Practice.

Fig. 4 shows that 75% of software professionals are using scrum meeting frequently and only 25 % irregularly. Fig. 5 shows that 70% respondents are always using agile methodologies for the development of software.

Last section (Fig. 6) deals with the software quality affecting factors of agile methodologies. Fig. 6 shows that the comparison of software quality affecting factors with agile methodologies. The questionnaire result indicates that almost all developers use agile methodologies to fulfill the qualities of software. In the survey, 10 respondents answer that they use agile methodologies for the development of software and the reason behind that software qualities like correct-ness, reliability, usability, extensibility, reusability, testability, portability, maintainability and efficiency.



Fig. 4.   Shows Agile Methods and Software Quality Factors Questionnaire Results–Percentage of using Scrum Meeting.



Fig. 5.   Shows Agile Methods and Software Quality Factors Questionnaire Results–Frequency of using Agile Methodologies.



Fig. 6.   Comparison of Software Quality Affecting Factors with Agile Methodologies.

### E.  Data Validation

For the confirmation of data validation, we interact with the respondents and asked them about the surety of their answers. Then, they inform us that, they also consulted with other team members in the company who had higher experience in the agile field. So, after sure about the opinion of the professionals, they filled in the questionnaire. That's why, respondents were confident about their answers.

## VII. CONCLUSIONS AND FUTURE WORK

In this paper, a questionnaire survey was conducted. By this survey we collect data from different software development companies' employees for the comparison of software quality affecting factors with agile and scrum method. We have identified software quality affecting factors such as correctness, reliability, portability, testability, efficiency and extensibility. The main advantage of agile technique is customer satisfaction and its welcome user requirement changing at any phase. By using agile method, the software has almost all software qualities. In future work the questionnaires should be repeated with additional respondents to the results from the user feedback presented in this study. We will enhance our area of research about agile methodologies and software quality affecting factors. In current research we only target the problem that is comparison of agile methodologies with soft-ware quality affecting factors, but in future we will try to identify the all identified issues of agile methodologies and their effect on software product and customer's trust. According to the questionnaires the percent-age of software qualities increased by using agile methodologies. In future work we should work on customers satisfaction by using agile methodologies. How well agile methodologies fulfill customers' satisfaction for the development of software.

REFERENCES

[1] Oberscheven, Falk Martin. "Software Quality Assessment in an Agile Environment." Faculty of Science of Radboud University in Nijmegen. (2013).

[2] Hossain, Amran, Md Abul Kashem, and Sahelee Sultana. "Enhancing software quality using agile techniques." IOSR Journal of Computer Engineering 10.2 (2013): 87-93.

[3] Awad, M. A. "A comparison between agile and traditional software development methodologies." University of Western Australia (2005).

[4] Sharma, Sheetal, Darothi Sarkar, and Divya Gupta. "Agile processes and methodologies: A conceptual study." International journal on computer science and Engineering 4.5 (2012): 892.

[5] Khalane, Tiisetso, and Maureen Tanner. "Software quality assurance in Scrum: The need for concrete guidance on SQA strategies in meeting user expectations." Adaptive Science and Technology (ICAST), 2013 International Conference on. IEEE, 2013.

[6] Sirshar, Mehreen, and Fahim Arif. "Evaluation of Quali-ty Assurance Factors in Agile Methodolo-gies." International Journal of Advanced Computer Sci-ence 2.2 (2012): 73-78

[7] Ullah, Malik Imran, and Waqar Ali Zaidi. "Quality Assurance Activities in Agile: Philosophy to Practice." (2009).

[8] Sharma, Mohit Kumar. "A study of SDLC to develop well engineered software." International Journal of Advanced Research in Computer Science 8.3 (2017).

[9] Rodríguez-Hernández, V., et al. "Assessing quality in software development: An agile methodology approach." Journal of Advanced Computer Science & Technology 4.2 (2015): 225-230.

# Reengineering Framework to Enhance the Performance of Existing Software

Jaswinder Singh[1]

Department of Computer Application
IK Gujral Punjab Technical
University
Kapurthala, Punjab, India

Kanwalvir Singh[2]

Department of Computer Science
and Engineering
BBSB Engineering College
Fatehgarh Sahib, Punjab, India

Jaiteg Singh[3]

Department of Computer
Applications
Chitkara University
Rajpura,Punjab,India

*Abstract*—**Term reengineering refers to improve the quality of the system. Continues maintenance and aging degrade the performance of the software system. Right approach and methodology must be adapted to perform reengineering. With lack of right approach and methodology, reengineering itself will be costly and time-consuming. For the process of reengineering main concerns include when to reengineer, how to estimate cost, the right approach for reengineering, and how to validate software enhancement. This research paper proposed a framework to identify the need for reengineering, to estimate the cost of reengineering, and to validate software quality improvement. Research work used the agile methodology to perform tasks of reengineering. Reengineering needs are identified using prediction based decision tree approach. Reengineering is applied using the agile Scrum methodology. Cost estimation is done using story point estimation. Performance analyses are done using complexity measures analysis of the internal design metrics and mean time to execute metric. The research used various automated tools like CKJM ver1.9, Rapid Miner studio ver7.1, and Net beans7.3 framework.**

*Keywords—reengineering; maintenance; decision tree; agile methodology; scrum*

## I. INTRODUCTION

Maintenance is one of the most critical phases of software development. Continues maintenance degrades the software quality and increases the maintenance cost. The software reengineering plays a vital role to improve the quality of software. Reengineering is required to upgrade the existing system and to reduce the maintenance cost. Many researchers [1, 2] proposed a framework for reengineering identification and reengineering cost estimation. But these frameworks are not able to handle the ever-changing behavior of customer needs and requirements. These existing frameworks lack the flexibility to adopt the changes and as well as to estimate cost. Earlier approaches are based upon conventional engineering methods. Since a few decades, we have also seen changes in software development approach, especially with the use of agility in software development. So need is to use a comprehensive approach to provide a new framework for software reengineering that is flexible as well as an interactive model to adopt the customer requirements and able to perform the cost estimations. This research work proposed a framework to identify the need for reengineering, estimate the cost of reengineering, uses an agile approach, reduce the maintenance

cost of the reengineered system and finally evaluate the performance of reengineering system. Proposed research work provides a vision for developers to quickly identify reengineering needs and able to apply to reengineer in a people-centric environment using agile. Research work in this paper organized under different sections. Related work discussed in the literature review section. Section 3 describes the research methodology used in this paper. Another section identifies whether the software is required to be reengineered or maintained. Reengineering agile model and estimations discussed in Sections 5 and 6. Performance evaluation is given in the last section.

## II. LITERATURE REVIEW

The existence of a reengineering approach is not new. It has been observed that due to continuing changes in the existing software, software quality deteriorates [3] and reengineering must be performed to adapt the changing requirements of end-user. Researchers identified [4] the importance of reengineering and stated the importance of information technology in software reengineering. Reengineering performs preventive maintenance for the software system [5]. Reengineering includes three important subtasks named reverse engineering, restructuring or alteration, and forward engineering. Reengineering tasks are shown in Fig. 1 [6]. Researcher [7] also identified various benefits like better software quality, fewer maintenance efforts, and ease of software testing and a better understanding of the software. Sneed observed the impact of reengineering over maintenance [8]. Researchers [9] proposed a cost model for reengineering using the conventional approach. Agile methodology has proven to be a successful approach to software development for the last few years [10]. Agile is integrated with the field of reengineering by many researchers. The researcher proposed N-Process model [11]. N-process model is N-shaped reengineering structure to perform various tasks of reengineering. Tasks are mapped in N shaped structure. Other work gives the idea of service-oriented software reengineering [12]. Service-oriented computing paradigms applied to enhance the legacy systems. Work is also done to provide prototypes at the initial stages of reengineering [13]. Researchers also worked on aspect-oriented reengineering [14]. In aspect-oriented reengineering, reengineering work is validated by applying various object-oriented metrics, and tasks were performed in short iterations of agile.

Fig. 1. Software Reengineering [6].

## III. RESEARCH METHODOLOGY

The case study includes twenty open sources, Java-based software systems. The complexity of the Java-based system is measured using Chidamber and Kemerer metric popularly known as CK Metric [15] of object-oriented software. Six basic metric sets of CK metric suit include Depth of the Inheritance Tree (DIT), Number of Children (NOC), Response for a Class (RFC), Lack of Cohesion of Methods (LCOM), Weighted Methods per Class (WMC) and Coupling between Object Classes (CBO). CKJM tool is used to measure the basic set of CK metric. Using CK metric, internal design complexity of the software system can be determined. To identify the need for reengineering, prediction based decision tree approach [16] is used for the software systems. Once the software got categorized for reengineering or maintenance requirements [17], agile development approach is used to get the software reengineered and also to estimate the cost of reengineering. Performance of the reengineered system is evaluated by comparing the design complexity of reengineered and old software. Classes having complex design are the candidate for reengineering.

## IV. IDENTIFY THE NEED FOR REENGINEERING

The decision among maintenance and reengineering is made using prediction based decision tree approach. Data set consist of twenty software systems divided into two parts. Each part is a mix of varying lines of codes and complexity. Research work considered fifteen software projects as a training data set, and five projects as model data set. Attributes of the training data set will be applied to predict model data sets. For implementing the predictions using a decision tree approach, an average of internal design complexity and size of software systems act as two main metrics. Table I shows Java-based software systems considered under the training data set. Software belongs to different size and having different average internal design complexity. Internal design complexity is measured using a basic set of CK metric suit.

Five projects are considered under the model data set. Table II shows various software systems for model data set. Training data set will be applied to the mode data set to predict reengineering and maintenance requirements.

The process of Applying and executing a decision tree using rapid minor tool is as followed.

- Import training data set having fifteen Java-based software projects.

- Roles are used to selecting attributes. First role operator is used to choosing a category attribute. Average complexity and size are chosen as two parameters.

- The second role is used to skip project names from the analysis part.

- Predictions are made using a decision tree with Decision Tree operator.

- Training data set will be input to the Decision Tree operator.

- The classification model is the output of Decision Tree that will be used for decision making

- Model data set is imported using a retrieve operator. New role is applied to the data set setting parameter 'Category' which is required to be predicted.

- There are two outputs of Apply model. One output is the prediction of attributes applied to model data using training data, and other is training data itself.

- The complete design is presented in Fig. 2.

- Finally, decision tree and predictions can be viewed by executing the designed scenario.

TABLE I.     TRAINING DATA SET COMPLEXITY MEASURE [17]

| SrNo | Software | SLOC(Size) | Mean Complexity |
|---|---|---|---|
| 1 | PongGame Software | 713 | 31.3 |
| 2 | Software ChessGame | 150 | 29 |
| 3 | Battle City Software | 563 | 77.2 |
| 4 | Software Customer Info System | 1139 | 120.3 |
| 5 | Parser Software | 143 | 13.8 |
| 6 | Software Scheduling and dispatch | 203 | 82.7 |
| 7 | Dictionary Software | 337 | 24.7 |
| 8 | Software ChatServer | 284 | 24.3 |
| 9 | My Notepad Project | 290 | 2 |
| 10 | Trigonometric Function Software | 634 | 362.7 |
| 11 | SoftwareCricketAnalyzer | 234 | 16.7 |
| 12 | Diary App Software | 431 | 26.3 |
| 13 | Software TicTacToe | 276 | 12.7 |
| 14 | FIFO Software | 637 | 75 |
| 15 | Software BounceBall | 160 | 12.1 |

TABLE II.     MODEL DATA SET COMPLEXITY MEASURE [17]

| Sr. No | Software | SLOC(Size) | Mean Complexity |
|---|---|---|---|
| 1 | E-library Software | 323 | 55 |
| 2 | Shopping Cart Software | 154 | 24.7 |
| 3 | Code Level Security Software | 201 | 144.5 |
| 4 | Point of Sale Software | 1082 | 526.5 |
| 5 | SmartFileConverter Software | 440 | 39.7 |

Fig. 2 represents the design interface. Apply Model operator is required to apply a decision tree on the model data set.

Once executed, the decision tree will appear, as shown in Fig. 3. A decision tree is made up of nodes and edges. The root of tree denotes prominent predictor. Thus it is observed that Average complexity is our best predictor of deciding reengineering requirements. It predicts whether or not the Java project requires reengineering. The predicted value for Average complexity comes out to be 25.5. The second node is of size attribute. Thus best predictor at second level is source line of code SLOC (Size). The tree from root to the leaf can be interpreted as if Average complexity >25.5 and SLOC (Size)>176.5 the software undergoes reengineering. Thus except shopping cart software of model data set given in Table II, all other software are the candidate for reengineering.



Fig. 2.    Decision Tree Modeling in Rapid Miner [18].



Fig. 3.    The Decision Tree Structure for Model Data Set [18].

## V. AGILE REENGINEERING MODEL

Once the software is chosen for performing reengineering, an approach to perform reengineering is required. Among various software development approaches, one of the most popular and acceptable methods for development is agile [19, 20]. Development of software in agile include active participation among various stakeholders of software. Many agile frameworks exist like Scrum, Extreme Programming, Lean programming, United Process, Kanban, FDD (Feature-Driven Development), Crystal, DSDM (Dynamic Systems Development Method). Among these frameworks, Scrum is one of the most useful approaches by IT professionals. In a report of Scrum alliances [21], 89% of agile users used the scrum approach. Major Scrum activities include Scrum planning; Daily Scrum, sprint review, and sprint retrospective shown in Fig. 4. Sprint represents a single iteration in fix time. Many sprints can be used to develop the required product. Requirements are analyzed in terms of user stories, and estimation is performed by assigning story points to each user story. Whole requirements are collected as a product backlog. Requirements of high priority assembled in the sprint backlog. In sprint planning, the work required to perform decided. A product backlog is analyzed, and sprint backlog is prioritized in this phase. The team meets every day to evaluate the progress of the sprint. The team reviews the work and changes required. The team finally discusses goals achieved and if anything went wrong, ways of improvement.

Because of the flexible and interactive approach of software development, it is decided to perform reengineering using agile methodology. The inclusion of reengineering tasks with agile scrum methodology is shown in Fig. 5. Proposed agile reengineering model retains the essence of reengineering and agility. Three tasks of reengineering are performed using an Scrum methodology. All three reengineering tasks are enclosed in one sprint of three-week iteration.

Agile reengineering model works as follows:

- Ensure planning of release of reengineering software, planning of iterations (time allocation for iteration, team members required, etc.), and estimation of cost. All requirements are prioritized in the product backlog.

- Requirements required to implement in one sprint are assigned to the sprint backlog. Planning is done by the Scrum team, including all stakeholders.

- Analysis of Reengineering Requirements in terms of user stories and allocation of story points.

- Execution of sprint with 3-week iteration to accommodate forward, alteration and reverse engineering

- Retrospective action to confirm the implementations of required objectives. After iteration, estimation, and speed of requirement implementations (velocity) is verified.

- Daily planning is performed every day.

- One sprint perform reverse, alterations and forward engineering

- Integration for final complete System.



Fig. 4. Scrum Activities.



Fig. 5. Agile Reengineering Model.

## VI. ESTIMATIONS

Estimations of effort and cost are essential for any projects. One important aspect of proposed framework is to estimate efforts and cost of existing systems. Proposed work measures effort and cost estimations of reengineering with the help of an agile approach.

### A. Efforts Estimations

Reengineering efforts are estimated by assigning story points to the required tasks. Planning poker also called Scrum poker is highly acceptable techniques for assigning story points to reengineering requirements. As stated by Cohen [22], "Planning poker is a proper mix of expert opinion, analogy, and disaggregation techniques which can successfully give quick and reliable estimates.

Benefits of planning poker include

- Scrum team, including Scrum master, product owner, and development team (developer, testers, Analyst) sit together to perform estimations.

- Both high and low estimation points for user story are discussed. Meeting avoids the problem of conflict for the future.

- Work starts when all members are agreed upon the same consensus so commitment for the project increases.

As everyone gets a chance to justify himself and everyone's opinion is welcomed, so no chance of dominance of individual arises.

### B. Cost Estimations

Reengineering cost is estimated considering the cost of human resources, Time required to complete the tasks, cost of other resources required (hardware, software licensing, etc.).Sprint is planned, and the time of one sprint is estimated. Formulations of Various cost estimations are as follows:

- Let Ann. Sal represents the annual salary of the Scrum team member.

- Acc.Sal denotes accumulated salary, which is the sum of the annual salary and other expenses. For reengineering process, we can include other expenses as half of the annual salary of employee as suggested by Cohen [22] for software development. Character 'i' denotes n number of members of the Scrum team.

$$(Acc.Sal)i \; \forall \; (i = 1,2 \ldots n) = (Ann.Sal)i \; \forall \; (i = 1,2 \ldots n) + \left(\frac{1}{2}\right) * \left((Ann.Sal)i \; \forall \; (i = 1,2 \ldots n)\right) \quad (1)$$

- Let K denotes the number of weeks per iterations then salary per iteration (Sal.Iter) is

$$(Sal.Iter)i \; \forall \; (i = 1,2 \ldots n) = \left(\frac{K}{52}\right) * (Acc.sal)i \; \forall \; (i = 1,2 \ldots n) \quad (2)$$

- Let P denotes the estimated number of days required for an employee to work on the project then the percentage of time spent by employees will be

$$(Time.Spent)i\forall \; (i = 1,2 \ldots n) = \left(\frac{1}{5*K}\right) * (Pi\forall \; (i = 1,2 \ldots n) * 100) \quad (3)$$

- Accumulated cost per time spent (Acc.Cost.Time.Spent) for each member of scrum team is

$$(Acc.Cost.Time.Spent)i\forall \; (i = 1,2 \ldots n) = ((Sal.Iter)i * (Time.Spent)i))\forall \; (i = 1,2 \ldots n) \quad (4)$$

The initial cost may be estimated in a long time, and estimations can be reviewed after each sprint.

## VII. CASE STUDY

As discussed, reengineering is performed using the Scrum methodology. For our case study, software named CodeLevelSecurity from Table II is chosen for reengineering. Table III shows all the classes of this software. Complexity measures for software are determined using a basic set of CK metric.

Three classes are selected for sprint backlog depending upon their usage and importance in the project. Login, IDE, and UserDetail classes have been chosen to perform reengineering. Several reengineering tasks performed in one sprint are discussed in Table IV. Sprint iteration of 3 weeks is estimated for implementing the required reengineering. Story points assigned to Login, IDE, and UserDetails are 2, 8 and 5 respectively. Each task in sprint backlog is estimated on an hourly basis.

TABLE III. CANDIDATE SOFTWARE FOR REENGINEERING [19].

| Sr No | Classes | Design Metrics | | | | | |
|---|---|---|---|---|---|---|---|
| | | WMC | DIT | NOC | CBO | RFC | LCOM |
| 1 | Login | 12 | 6 | 0 | 9 | 8 | 60 |
| 2 | IDE | 17 | 6 | 0 | 17 | 21 | 70 |
| 3 | UserDetail | 23 | 5 | 0 | 12 | 09 | 183 |
| 4 | program access report | 12 | 6 | 0 | 8 | 9 | 62 |
| 5 | Profile detail' | 11 | 5 | 0 | 6 | 4 | 25 |
| 6 | User report | 8 | 6 | 0 | 6 | 7 | 24 |
| 7 | Saved program report | 14 | 6 | 0 | 10 | 8 | 73 |
| 8 | User maintenance | 2 | 1 | 0 | 1 | 5 | 1 |
| 9 | Program update report | 12 | 6 | 0 | 9 | 94 | 48 |
| 10 | Main frame | 22 | 6 | 0 | 19 | 9 | 233 |
| 11 | program report | 8 | 6 | 0 | 6 | 4 | 24 |

TABLE IV.     VARIOUS REENGINEERING TASKS PERFORMED IN ONE SPRINT

| S. No | Reengineering Tasks | Hours allotted |
|---|---|---|
| 1. Reverse Engineering | | |
| 1.1 | Generating Documentation/re-documentation | 6 |
| | Design Recovery | |
| 1.2 | High-level design analysis | 4 |
| 1.3 | Low-level design analysis | 8 |
| 1.4 | Analysis of restructuring requirements. | 12 |
| 2. Alterations and forward engineering | | |
| 2.1 | Classes Remodeling | 6 |
| 2.2 | Design Complexity reduction in classes through Alterations | 18 |
| 2.3 | Performing Unit test | 6 |
| 2.4 | Performing Regression test | 12 |
| 2.5 | Increment Integration | 6 |
| 2.6 | Testing | 6 |
| 2.7 | Retrospective | 6 |
| | Total Sprint Time | 90 Hrs |

## A. Cost Estimation

Cost is estimated using the equations (1), (2), (3), and (4). Scrum team includes three members named Scrum master, programmer, and tester. Consider annual Salary of Scrum master, programmer and tester as $1,50,000, $60,000 and &60,000 respectively. Accumulated salary (Acc.Sal) as given in equation (1) is calculated as $225000, $90000, $90000 for each employee. For the reengineering process with three-week iteration (putting K=3 in equation (2)) accumulated salary per iteration (Sal.Iter) is approximately $12981, $5192 and $5192 corresponding to all scrum team. Time Estimation is performed for each Scrum team member. The estimated value of P is15 days for Scrum master, 12 days for programmers and eight days for tester then the percentage of time spent by every member as calculated by equation (3).

For scrum Master, estimated days (P) are 15.

By putting the value of P in equation (3),

(1/(5*3)) * (15*100) that is 100 % time.

For programmer, estimated days (P) are 12.

By putting value of P in equation (3),

(1/ (5*3)) * (12*100) that is 80 % time.

For Tester, estimated days (P) are 8.

By putting the value of P in equation (3),

(1/ (5*3)) * (8*100) that is approximately 53% time.

Accumulated cost per time spent as given in equation (4) for Scrum Master is $12981.Similarly, by putting values in equation (4), Accumulated cost per time spent for the programmer is approximate $4154 and for the tester is $2752. So the total cost of project per iteration is $19887.Thus we can estimate the actual cost of reengineering. Still, we can assume the uncertainty factor which can reduce or increase the actual cost of reengineering. As suggested by Cohen [22], the actual cost in an agile environment can be + or – 25% of estimated values.

## B. Evaluating Complexity Reduction and Performance Improvement of Reengineered Software

Once reengineering is performed, software is analyzed for complexity reduction and performance up gradation. Outcomes of reengineering interpreted in three ways.

- Complexity in terms of the Basic set of CK metrics has been reduced in reengineered classes.

- Reduction of software complexity results in an improvement in maintainability.

- Improvement in the overall mean time to execution (MTTE) of the project, due to CK metric value reduction.

*1) Complexity in terms of the basic set of CK metrics reduced in reengineered classes:* It has been observed that by applying reengineering tasks, the inherent design complexity of classes measured in terms of CK metrics has been reduced to a reasonable extent, as shown in Table V. For all the three classes of the project, there is a reduction in WMC, CBO, RFC, and LCOM. Due to reengineering, classes are restructured, and alterations are done at the function level. The numbers of functions and dependencies in each class have been reduced. Comparisons of reengineered and old classes are shown in Table V.

TABLE V.     CK METRIC COMPARISON BEFORE AND AFTER REENGINEERING

| Metrics & Software Classes | | Design Metrics | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | WMC | DIT | NOC | CBO | RFC | LCOM | Total |
| Login Class | Reengineered | 4 | 6 | 0 | 5 | 48 | 2 | 65 |
| | Before Reengineering | 12 | 6 | 0 | 9 | 78 | 60 | 165 |
| IDE Class | Reengineered | 4 | 6 | 0 | 12 | 102 | 0 | 124 |
| | Before Reengineering | 17 | 6 | 0 | 17 | 121 | 60 | 221 |
| UserDetail Class | Reengineered | 6 | 6 | 0 | 4 | 67 | 0 | 83 |
| | Before Reengineering | 23 | 5 | 0 | 12 | 109 | 183 | 332 |

*2) Reduction in the software complexity results in an improvement in maintainability of software system:* As stated [23], larger the values of CK metric more will be the software complexity, and hence the software will be more error-prone. Total reduction of CK metric values for all the three classes are shown in Fig. 6. The Fig. 6 shows CK metric analysis for both reengineered and existing candidate classes. On the x-axis, there are three classes, and on the y-axis, CK metric complexity is depicted.

Once the software got reengineered, the maintenance cost of the reengineered project will be undoubtedly low. As suggested by Chaudhary and Ugrasen [24], maintenance can be estimated based on story points. Before reengineering, story points for three classes were fifteen, then after reengineering they are reduced to six only. Story points assigned to Login, IDE, and UserDetails are 1, 3 and 2 respectively. That means

- More classes can be accommodated (if required) in one iteration of the Scrum

- Will results in the reduction of the cost

- less time spent to perform changes

- fewer complexity results in less possibility to induce more errors

So the system once reengineered can survive longer and further can adapt changes (undergo maintenance) with less cost and time.



Fig. 6.    CK Metric based Maintainability Comparisons.



Fig. 7.    MTTE Values for Existing and Reengineered Project.

*3) Improvement in the overall mean time to execution (MTTE) of the project:* Another improvement in the reengineered project is in the meantime to execute (MTTE). For all the three old and reengineered class modules, samples of 35 executions are taken. Net beans7.3 is used with system configuration of i5-4th gen processor, 8GB RAM, HDD 1TB and Java7. MTTE is 290.6 milliseconds for classes of the old project and 271.7 milliseconds for reengineered project classes. MTTE analysis is shown in Fig. 7.

## VIII.    RESULTS AND DISCUSSION

The proposed research work is discussed from Sections IV to VII. Except shopping cart software of model data set given in Table II, all other software is the candidate for reengineering. These software systems are CodeLevelSecurity, PointofSale, E-Library, and SmartFileConvertor. Among these four candidate systems, CodeLevelSecurity is chosen to reengineer. Three classes of the software are reengineered. CK metric suit is used to measure the design complexity of software. Agile Reengineering Model is proposed to perform estimations and to apply reengineering tasks. Two main objectives of the proposed model included:

- To apply agile-reengineering development approach to perform reengineering on the candidate system.

- Performing effort and cost estimations for reengineering.

After performing reengineering, the reengineered system is evaluated for maintainability and performance up gradation. It is validated that the reengineered system performs much better than the existing candidate system. Results for average complexity and MTTE are shown in Table VI.

It is important to note that the candidate software gone through reverse, alteration and forward engineering. After reengineering, the numbers of functions in the three classes are also reduced from thirty three to thirteen. So in place of refactoring, the reengineering process is applied to the software to inculcate requirements of reducing complexity and increasing performance. Not only the complexity of the software is reduced but the performance of the system is also improved.

TABLE VI.      REENGINEERED SYSTEM PERFORMANCE MEASURES

| Sr No | Software Type | Average Complexity of three classes | MTTE in milliseconds for complete software |
|---|---|---|---|
| 1 | Reengineered Software | 272 | 271.7 |
| 2 | Existing Candidate Software | 718 | 290.6 |

## IX.    CONCLUSION

Proposed work introduced a framework that identifies reengineering requirements for software using prediction based decision tree approach. Agile Reengineering model uses features of the agile development approach with a reengineering approach. Cost estimation is done using story

point technique. Complexity and performance analyses are performed using CK metric and MTTE metric. Using agile reengineering approach, reduction in the cost of maintenance and improvement in maintainability observed. After reengineering of classes of software, complexity is reduced to a greater extent. As a result indicates, performing reengineering by using agile methodology is beneficial in terms of implementing requirements, estimating cost, and enhancing performance. Cost estimations are realistic and involve a consensus of all stakeholders. Software complexity in terms of the internal design of software calculated using CK metric. This research can be a benchmark to the software development companies to identify whether software needs maintenance or reengineering. Also, cost estimations can easily be measured for the software to be reengineered.

Further, software complexity can be validated using other software metrics like cyclomatic complexity, reliability, etc. For a generalization of the framework and to make it industry ready more and more software systems can be considered to make an extensive training data set. Larger the training data set, more accurate will be the predictions. The experience of Scrum team will play crucial role to successfully implement Agile Reengineering Model.

### REFERENCES

[1] H. M. Sneed, "Estimating the Costs of a Reengineering Project," Proceeding of 12th Working Conference Reverse Engineering. IEEE CS Press, Pittsburgh, USA , 2005, pp. 111–119.

[2] S. Sood, Software reengineering-A metric set based approach, Himachal Pradesh University.2012.

[3] M.M. Lehman, "Programs, life cycles, and laws of software evolution," Proceedings of the IEEE. Vol. 68, No. 9, 1980, pp.1060-107.

[4] M. Hammer and J. Champy, Reengineering the Corporation: A Manifesto for Business Revolution. New York: HarperCollins Publishers, 1993.

[5] S. Ian, Software Engineering, 9th ed., Pearson publication. 2014.

[6] E.J. Byrne, "A conceptual foundation for software re-engineering," Proceedings of Conference on Software Maintenance 1992, Orlando, FL, USA, 1992, pp. 226-235.

[7] R.S. Arnold, " Software Restructuring," Proceeding of IEEE, vol. 77, no. 4, April 1989, pp. 607–617.

[8] H.M. Sneed and A.A. Kaposi, "Study on the effect of reengineering upon software maintainability," Proceedings. Conference on Software Maintenance 1990, San Diego, CA, USA. 1990, pp. 91-99.

[9] P. Kumawat and N.Sharma, "Design and Development of Cost Measurement Mechanism for Re-Engineering Project Using Function Point Analysis," In R. Kamal, M. Henshaw and P. Nair (eds.), International Conference on Advanced Computing Networking and Informatics. Advances in Intelligent Systems and Computing, vol. 870. Springer, Singapore,2019.

[10] J. kisielnicki and A.M. Misiak, "Effectiveness of agile compared to waterfall implementation methods in it projects: analysis based on business intelligence projects", Foundation of management, vol. 9, No. 1, pp. 273–286, 2017.

[11] A. Sahoo. D. Kung, and S. Gupta, "An Agile Methodology for Reengineering Object-Oriented Software," Proceeding of 28th International Conference on Software Engineering & Knowledge Engineering: SEKE. California; USA, 2016.

[12] S. Chung, D.H. Won, S. H. Baeg and S. Park, "A Model-Driven Scrum Process for Service-Oriented Software Reengineering:mScrum4SOSR," Proceeding of 2nd International Conference on Computer Science and its Applications, Korea (South), 2009, pp. 1-8.

[13] M.I. Cagnin, J. C. Maldonado and R.D. Penteado, "PARFAIT: Towards a framework-based agile reengineering process," Proceedings of the Agile Development Conference (ADC): USA. 2003.pp. 22-31

[14] P.O. Adrian, "Aspect-Oriented Reengineering of an Object-oriented Library in a Short Iteration Agile Process," Informatica, vol. 35, No.4, 2012.

[15] S.R. Chidamber and C.F. Kemerer, "A metrics suite for object-oriented design," IEEE Transactions on Software Engineering, Vol. 20, No. 6, 1994, pp.476-493.

[16] M. North. Data mining for the masses. Global Text Project. August 2012.

[17] J. Singh, A. Gupta, and J. Singh, "Identification of requirements of software reengineering for JAVA projects," Proceeding of International Conference on Computing, Communication, and Automation (ICCCA). Greater Noida; India, 2017, pp.931-934.

[18] J. Singh, K. Singh, and J. Singh, "Reengineering framework for open source software using decision tree approach," International Journal of Electrical and Computer Engineering (IJECE), vol. 9, No. 3, 2019, pp.2041-2048.

[19] P. Serrador and J.K. Pinto," Does Agile work? - A quantitative analysis of agile project success," International Journal of Project Management, vol. 33,No. 5, 2015,pp.1040-1051.

[20] J. Kisielnicki and A.M. Misiak, "Effectiveness of agile compared to waterfall implementation methods in its projects analysis based on business intelligence projects," Foundation of management, vol. 9 , No.1, 2017, pp. 273–286.

[21] scrumalliance.org. Scrum Alliances report-2017 [cited 2019 March 9] available from https://www.scrumalliance.org/learn-about-scrum/state-of-scrum/2018-state-of-scrum,.

[22] M. Cohan. Estimating and planning. Pearson Education.USA. 2006.

[23] V.R. Basili, L.C. Briand, and W.L. Melo, "Validation of object-oriented design metrics as quality indicators," IEEE Transactions on Software Engineering, vol. 22, No.10, 1996, pp.751-761.

[24] J. Choudhary and U. Suman, "Story Points Based Effort Estimation Model for Software Maintenance," Procedia Technology. Elsevier.vol.4, 2012, pp.761-765.

# Dynamic Bandwidth Allocation in LAN using Dynamic Excess Rate Sensing

Muhammad Abubakar Muhammad[1], Muhammad Azhar Mushatq[2], Abid Sultan[3], Muhammad Afrasayab[4]

Department of Computer Sciecne & IT, University of Sargodha Sub- Campus Bhakkar, Bhakkar, Pakistan[1, 2, 3]
Department of Computer Sciecne, Govt. College University, Faisalabad, Faisalabad, Pakistan[4]

*Abstract*—**Today human and information processing system both need rapid access to anything they want on the internet. To fulfill these needs more and more internet service providers with a large amount of bandwidth are introducing themselves in the market. For these providers, a lot of bandwidth is free during off-peak hours while during peak hours the total available bandwidth might be insufficient. The primary purpose of our research is to divide and distribute the excessive bandwidth among the users during off-peak hours to attain the maximum user satisfaction. In order to do this dynamic excess rate (DER) scheme and its frame work is proposed in this paper.**

*Keywords—DER; ISPs; PIR; DBA; MRT; CIR*

## I. INTRODUCTION

After wheel and electricity, the Internet has transformed the world into an information village. The provision of the internet is carried out by specialized operators called "Internet Service Provider" (ISP). These ISPs operate on a different level ranging from continental to metropolitan local or ISPs. Internet service providers, regardless of the level on which they are operating, provide both network infrastructure and the bandwidth to both individuals and the corporate world.

The users want the maximum bandwidth at lowest billing, while on other hand the provider has bulk of bandwidth which can be free at off-peak hours. To distribute the bandwidth among a large number of users ISP uses different mechanism most commonly as static allocation at subscribed peak information rate (PIR). In static allocation, the excessive bulks of the bandwidth cannot produce advantages for both the users and providers. Where on the other hand bandwidth eager users want to browse and download at maximum data rates within their budgets. The above scenario arises as a challenge for the local and regional ISPs, trying to full fill the needs of their users with the best data rates.

In order to create a better user experience, dynamic allocation of the excessive available bandwidth through sensing the real-time utilization is the need of time. In this paper, we have proposed a methodology called DER, which will serve the purpose of dynamic allocation of bandwidth through sensing real-time available bandwidth. Moreover, it will also improve the quality of services (QoS), user experience and the business opportunities for providers of a network [1].

In dynamic bandwidth allocation scheme, allocation can be changed dynamically depending on different factors. So far Dynamic Bandwidth Allocation (DBA) is completely adaptable only on OFC (Optical Fiber Cable) networks. In DBA the link capacity is enhanced with extra bandwidth (if available) on per-flow per request basis which can be implemented only by the access router itself [2]. In other word allocation schemes, link capacity is enhanced for a very short period of time, on request, and on a temporary basis [3]. The major benefits of the DBA may include [4].

- Efficient performance in case of over subscription.

- Low latency on both downstream and upstream.

- Strong QoS.

- Increased user satisfaction.

- Fairness in distribution among different users belonging to the same priority/group/class.

However, the two major drawbacks of most of the DBA algorithms are as:

- The allocation of excessive rates is only for a very short period of time even may be for a single flow request.

- Lack of supportability for shared networks like asymmetric digital subscriber line (ADSL), ADSL+, wireless local area network (WLAN) and Internet cables.

In this research, we have tried to address and resolve the above drawbacks of the DBA approach by designing a scheme named as DER that will be able to allocate the excessive bandwidth for a longer period of time and will work fine for all networks.

This paper is ordered as follows. Section II covers the related work whereas Section III presents the proposed methodology. In Section IV, flow charts of the proposed method lay are presented and finally in Section V conclusion and future work is presented.

## II. RELATED WORK

In literature, different schemes have been proposed by authors for bandwidth management of local area network (LAN) or WLAN. The main objective is to make efficient use of bandwidth while achieving good QoS. In this section, we have covered some of the existing bandwidth management schemes.

An AIFS or contention window based scheme for LAN is proposed by authors in [5]. The scheme is mainly based on

priority level. This priority level is implemented by using a priority table at the access point, and priority will be assigned to each user on the basis of physical address of the node. Moreover, bandwidth utilization will also be monitored by this technique. In case a user is making improper usage of bandwidth, the bandwidth from that user will be allocated to other high priority users [5].

For users who spend most of their time for downloading videos, a Weighted Fair Intelligent Bandwidth Algorithm (FIBA) has been proposed by Xiaomei Yu et al. FIBA algorithm is computationally less costly, prevent unwanted time delays and make usage of reporting instrument for efficiently running videos traffic. Once the connection is triggered by a switch, FIBA monitors the total available bandwidth and fairly re-allocate the bandwidth. The major fault of FIBA is that it has unlimited downloading and require a lot of configuration [6].

As described in cluster-based Bandwidth Allocation Algorithm, bandwidth can be allocated to clusters and specific users. In a cluster, bandwidth is allotted to the cluster head only, this happens when cluster head obtains content from the source station. The main objective is to obtain the highest throughput in a network connection [7].

In Smart Clustering Based technique, a clever clustering approach is designed for wireless network nodes through which bandwidth will be allocated dynamically to the users. Moreover, this method is useful for wireless nodes to adapt bandwidth allocation according to the changing number of users over time [8].

In order to increase benefits for the user, the authors in [9] designed the Dynamic Bandwidth Allocation algorithm to supervise and control traffic of the network. This algorithm ensures QoS by reducing network congestion and divide bandwidth equally to active network users as well as those users who are not currently active in the network. The main benefit of the algorithm is all users are treated equally within the network. In order to deal with hypermedia data that needs a large range of bandwidth, it uses a multiplexing technique. The major problem of the DBA algorithm is that it has absences of appropriate feedback and reporting instruments. In a situation where a large amount of data is downloaded bandwidth costs will be increased and this method only supplies bandwidth according to the demand [9].

The authors in [10] proposed mesh topology network architecture, based on reformative digital video broadcasting geostationary satellite (DVB-S GEO), which has the ability to connect heterogeneous terrestrial domain name server (DNS) with each other, within lowest delay. In proposed satellite network architecture an appropriate dynamic bandwidth management instrument is applied, that is used to enable the delivery of interactive internet protocol (IP) dependent hypermedia services according to the defined level of QoS [10].

## III. Proposed Methodology

This section covers the proposed DER framework for dynamic bandwidth allocation, see Fig. 5. The components of the framework are:

### A. Scheduler

The scheduler in the described flowcharts is a predefined process. The scheduler will perform the following function:

*a) Activation of DER:* As the aim of designing the DER flowcharts is to distribute the excessive bandwidth when the resources are free, or to withdraw when the network is over populated. The system support personals will decide the peak hours, off-peak hours, and the scheduler will be set accordingly. So the scheduler will automatically activate and deactivate the appropriate action.

*b) Running the DER*: The second function of the scheduler is to run the pre-activated DER after regular intervals of time. The duration of the intervals will be set by the network administrator, keeping in view the number of users and the hardware capabilities as the running of DER will consume the processing and other resources. The intervals should not be too short to choke the whole system.

### B. Active user Count

The DER mechanism starts with the count of active online subscribers or users. The active user means that the user is using the resources of the network at that time. The total users can be counted by their active IP addresses and also by their usernames.

### C. Count Total Bandwidth usage

After counting the total number of active users the total usage or the maximum usage of the bandwidth during certain time will be calculated. This will be done by analyzing the wide area network (WAN) interfaces in which backhaul link is terminated. The peak reading of the Multi Router Traffic (MRT) graph can also be taken as an alternative. However, DER can also be applied on the basis of only the active or alive user count. The current usage will be decided from the maximum data rate variable on last 15 minutes (or any other period of time) using Multi Router Traffic Grapher (MRTG) of WAN interfaces. If more than one WAN interface is present the system will simply sum up the average rates.

### D. Check the Actual Available Bandwidth

The actual bandwidth is not that has been written on your Service Level Agreement (SLA). It is that you actually get when the backhaul link is being fully utilized. If the provider has Committed Information Rate (CIR) backbone links the condition might be similar as written in the agreements but in case of shared links the condition will be different than the actual bandwidth will be remarkably low as written in SLA. In both scenarios, the system will take the maximum MRT value and will consider it the actual network bandwidth. The second option is more suitable in which service provider will manually define the total consumable bandwidth that he has. In both cases, the system will read the variable in which this is stored.

### E. Estimate Net Excessive Bandwidth

In order to find the excessive bandwidth (EB), the difference between the actual available bandwidth (TB) and the total bandwidth usage during the specific time will be used, see equation (1).

EB= TB – Max.Usage (in defined period of time)　　　　(1)

### F. Calculate Excessive Bandwidth Per user

In this step, the system will calculate the per-user available excessive data rate (NEB). This will be done by a simple equation as shown in (2).

$$NEB \text{ per user} = EB / \text{Total active user count (AUC)} \qquad (2)$$

### G. Analyze and Match

After taking both counts the algorithm will check the predefined thresholds for the implementation of the DER. These thresholds will be set by the service providers depending upon the following factors.

- The total backhaul bandwidth of network

- Total number of subscriber

- Total subscription for bandwidth

- Sharing Ratio

- Number of corporate users

- Number of home users

- Average, minimum and maximum data rates on WAN interfaces

- Peak hours and off-peak hours

After reading the threshold variable the system will analyze and match the thresholds to the E.B per user value. Moreover, the system will decide whether to distribute the excess network bandwidth or not.

$$Xc = E.B \text{ per user} - \text{Threshold} \qquad (3)$$

By this analysis system will also able to cut down the already allocated data rates in the peak time or in the conditions of over subscriptions.

### H. Decide the Action

After the analysis system will make the allocation decision on the of previously calculated values that decision might be

- Do nothing and wait for next turn

- Deactivate the DER

- Allocate the Excessive rate as per policy

- Withdraw some data rate from active users

### I. Decide the Per-user Allocation

As it is a strong probability that the more users will be online after taking the values from the user count so it is very important to keep them served and not starved. Thus we cannot allocate all excessive bandwidth. The system will allocate the excess data rate to every online user less than the value of the E.B per user as the something must be available for the new users till the next turn of the DER.

### J. Updating the user Queues

As most of the Radius-based authenticated networks maintain the user queues in which the data rate of the user is also defined. The DER will update the all or the selective queues with an incremental or detrimental data rate.

### K. Updating the user Profile

After the decision of the allocation, the user profile of the active user will be updated and the data rates in the databases will also be changed. In this respect in the database, a separate field of Excess rate might be included which will only be updated by the DER scheme. But here the best approach itself is to just update the queues, not the profiles as the queues are made only for the one connection or even for the session.

## IV. FLOW CHARTS OF PROPOSED METHODOLOGY

The proposed methodology of DER consists of four parts. As in Fig. 1 the first part is designed to count the online active users and to estimate the bandwidth usage and availability. It counts the total active subscribers, gets the overall maximum usage of the bandwidth and then performs some calculation to determine the total available excessive rate. This part of the methodology will be used in all other three parts as their action will be decided from the result achieved in part one. There are three possible actions listed as:

    *a)* No allocation.

    *b)* Allocate the excessive bandwidth.

    *c)* Withdraw bandwidth

The second flow chart in Fig. 2 deals with the possibility of no allocation of bandwidth, which occurs when there is not a sufficient bandwidth for the DER to be provided to the active users. This will happen when the number of active user or the usage of the total available bandwidth on the network is very high.

The third flow chart deals with the possibility to allocate the excessive bandwidth. This is the ideal condition or the main goal of our research. On the basis of the result achieved from the first part, the decision will be made of how much additional bandwidth will be allocated to the users. The additional bandwidth will be added up to the already allocated peak information rate (PIR) values in the user queues. This possibility will occur if the number of online users is less than the pre-described threshold and the usage of the bandwidth is low seeing Fig. 3.

The flow chart in Fig. 4 is designed to cope with the worst possibility of withdrawing bandwidth on a network. If the network is over subscribed, then to accommodate the maximum number of active users a small portion of the bandwidth (PIR allocated statically according to user profile) will be cut down. This will avoid the starvation of resources on the network.

Fig. 1. Flow Chart for Active user Count & Bandwidth Estimation.



Fig. 2. Flow Chart for Analyzer and Decision Maker.

Fig. 3.   Flow Chart for Excessive Bandwidth Allocation.



Fig. 4.   Flow Chart for Withdrawal of Bandwidth.

Fig. 5. DER Allocation Frame Work.

## V. CONCLUSION

In this paper, different we have proposed a new and effective DBA scheme, that focuses on the efficient allocation of the excessive network bandwidth among the users, especially during the off-peak hours when network traffic in marginally low. Moreover, the objective is to allocate the excessive, but not to withdraw and increase the user satisfaction by providing some extra other than that they have paid for.

We further aim to practically implement DER on any local area network so that bandwidth monitoring can be observed for efficient bandwidth allocation to increase the user satisfaction. Moreover, by conducting different experiments results are expected to show that total available backbone bandwidth will be increased in great numbers. As far as limitations are concerned, the proposed scheme is only for large local area network.

### REFERENCES

[1] Bisio, M. Marcliese, "Performance studies of bandwidth allocation technologies for QOS-constrained satellite networks", in IEEE international communication conference, Glasgow, 2007, pp.7-12

[2] J. Aracil, D. Morato, 'A-Priori Flow Bandwidth Estimates for Dynamic Bandwidth Allocation in ISP Access Links', in ITC Specialists Seminar on Access Networks and Systems, Gerona, Spain, 2001, p. 161-167.

[3] K. Kyeong Soo, 'On the excess bandwidth allocation in ISP traffic control for shared access networks', Communications Letters IEEE, vol.18, no.4, pp. 692-695, 2014.

[4] O. Haran, A. Sheffer,' The Importance of Dynamic Bandwidth Allocation in GPON Networks', PMC-Sierra Incorporation, 2008.

[5] Satyajit Sarmah, Shikhar Kumar Sarma, "Dynamic Bandwidth Management in 802.11Wireless LAN", International Journal of Computer Sciences and Engineering, Vol.-6, Issue-6, June 2018.

[6] Xiaomei Yu, Doan B. Hoang, David D. Feng,"Weight-Based Fair Intelligent Bandwidth Allocation for Red Adaptive Video Traffic", 2006.

[7] Bo Fan, Supeng Leng, Kun Yang, "A dynamic bandwidth allocation algorithm in mobile networks with big data of users and networks", IEEE Network, Volume: 30, Issue: 1, January-February 2016.

[8] Mohammed Awad, Abdelmunem Abuhasan, "A Smart Clustering Based Approach to Dynamic Bandwidth Allocation in Wireless Networks", International Journal of Computer Networks & Communications (IJCNC) Vol.8, No.1, January 2016.

[9] Elias, F. Martignon, A. Capone, "An Efficient Dynamic Bandwidth Allocation Algorithm for Quality of Service Networks".

[10] G. Xilouris, A. Kourtic ,G. Stefanou, " Dynamic bandwidth allocation for LAN2LAN interconnection using DVB-S satellite transmission", 2019.

# SentiNeural: A Depression Clustering Technique for Egyptian Women Sentiments

Doaa Mohey ElDin[1]

Information System Department
Faculty of computers and Information, Cairo University
Cairo, Egypt

Mohamed Hamed N. Taha[2], Nour Eldeen M. Khalifa[3]

Information Technology Department
Faculty of computers and Information, Cairo University
Cairo, Egypt

*Abstract*—Online Sentiments Analysis is a trending research domain of study which is based on natural language processing, artificial intelligence, and computational linguistics. Negation sentiments usually are not included in sentiment's analysis process. The depression analysis can be improved by negative sentiments processing. The negation sentiments may contribute to classify the depression problems and its causes. The proposed clustering technique can detect female sentiments from the sentiment's text through cause's classification, and the written sentiment style. The combination of sentiment analysis and neural network is a promising solution for creating a new clustering algorithm. According to Egypt Independent Journal in 2018, 7% of Egyptians suffer from mental illness reported by the Public Health Ministry in Egypt. But the real statistics is more than the mentioned percentage which causes major social problems such as divorce, avoiding responsibilities, or non-marriage. This paper will address the real statistics and cluster the depression causes and social status for each sentiment. Online women sentiments are the essential focus of this research. The proposed technique consists of two algorithms clustering for user's sex and classification algorithm for causes and responsibilities of women. The proposed clustering algorithm can recognize automatically for the sentiments user sex (females or males) and the level of depression automatically. The neural network clustering approach will produce accurate analysis results. The hardness of depression analysis implicitly and explicitly demonstrated in the different classifications for sentiments. This paper introduces a new technique for clustering sentiments and evaluating Egyptian women depression based on social sentiments.

*Keywords—Sentiment analysis; negation handling; depression analysis; neural network; clustering*

## I. INTRODUCTION

Recently, women spend a lot time online on social networks and communities. That becomes a main platform for expressing their feelings such positive or negative. The depression is one of the hardest text classification and recognition analysis. This research focuses on clustering the sentiment's texts from women or men. This paper takes care of the women sentiments and opinions, so that this technique can cluster the women' sentiments based on the written style, automatic responsibilities identification, and causes classification. The classification of responsibilities consists of a disease, and financial status. According to Egypt Independent Journal in 2018, 7% of Egyptians suffer from mental illness reported by the Health Ministry in Egypt[1].

More solid and real statistics can be obtained by analyzing the women sentiments on social media networks.

Most the depression analysis is based on age or work, but this research based on responsibilities and work for women in Egypt. After the horrific reports announced by the World Health Organization [2], According to Global Terrorism Database, the number of suicide bombers exceeded 400 attacks per year in Egypt [3], in addition to thousands of suicide attempts in the homes and streets of Egypt every year. About half of these attempts and suicides are from ladies or women.

The share of women in the labor force, or as some call the labor market, accounts for almost a third of women's participation: 30.8% are female, 21.2% are male, and 118.9 thousand women in the age group (12-17 years old) have married (1.2 thousand widows, 1.2 thousand divorcees, 111 thousand married women, 5.5 thousand marriages) and will show population estimates for 2017. The rate of unemployment between normal and 23.6% compared to 8.9% among males in 2016, and the proportion of those who work in permanent work 76.5% of the total working factor compared to 71.1% of males among the total working male. 17.6% compared to 4.5% for health insurance versus 64.1% compared to 45.4% for males [4].

Women responsible for children and spend money on them by 16%. The highest divorce rate for divorcees was less than 25 years in (2010) by (22.58%). Percentage of divorce to total divorce certificates in 2010 (3335 cases) 2.2%. More than 77% are workers at home and out home. These responsibilities can keep women in depression although denying more than 20% about depression. But they suffer deep problems in their work, home, and relationships [4]. So, this paper presents a classification algorithm for identifying the causes, and responsibilities recognition and sentiments evaluation on Facebook social network.

Sentiment Analysis becomes a recent research trend for improving software development, decision making. Social media can be an essential source of sentiments tweets, app reviews, bloggers comments. Sentiment analysis is defined by the subjectivity study (neutral vs. emotionally loaded) and class polarity (positive vs. negative) of a sentiment word [5].

This paper presents a proposed technique for evaluating the sentiment analysis of Egyptian women. Using Artificial Intelligence [6], machine learning[7], and deep learning

algorithms [8], it also can create a new clustering algorithm for written style for women based on written style and causes classification from Facebook reviews. The proposed technique can provide the identification social status and financial status for each woman that can be a parameter in the relationship between depression and causes.

The rest of this paper is organized as follows: Section 2 presents related works. Section 3 shows the proposed framework. In Section 4, experiment and discussion will be illustrated. Finally, Section 5, conclusion and proposes directions for future work will be presented.

## II. BACKGROUND

This section presents a summary of the essential aspects of sentiment analysis and discusses the initial research efforts in the sentiment techniques and applications.

### A. Sentiment Analysis

It expresses the mining of users' reviews [5 & 9] about product/ topic online as a social network or blogger. The recent research introduces the challenges [10] of sentiment meaning understanding and sentiment parsing sentences and words. Writer's review is a major criterion for the quality of services enhancement to grant deliverables. There are multisource of these reviews as social networks, web portals, and Blogs. Recently, the researches targets to analyze and benefit online sentiments not only on sentiment polarity positive or negative but also on the analyzing the variant feelings as depression, violent, risky, or serious for negative sentiments and wishing, safety, happiness for positive analysis.

### B. Sentiment Classification & Clustering

There are difference between classification [10] and clustering [11] for sentiment 'text. Sentiment Classification is a type of learning models for identifying data classes and pattern recognition based on labeled data. It depends on the predefined data or classes. There are several motivation in classification for supporting the results in various data types (as text, images, or videos) and domains (such as medical, tourism, online reviews, etc…). Clustering is considered a learning technique model also but it is based on unknown classes previously. It can support the objects and relationships automatically based on the input data. It is based on similarity between data and objects. It is known unsupervised technique. This clustering targets determining the relations between different sentiments, distributing them dynamically in natural groups, or discovering the most relevant subjects within their content and expressing them in their own terms

### C. Deep Learning

It is an important branch of machine learning. Deep learning [8, 12 & 13] has several algorithm to improve accuracy and performance. That can support the relationships and discover the hidden features between layers of data.

### D. Sentiment Summarization

The summarization of sentiment aims at generating a summary [14] of entities and opinion analysis that is representative of the average opinion and speaks to its

important aspects. A standard sentiment setting for social networks summarization assumes a set of sentiments S = {s1, . . . , sn} that include reviews about interesting topic or something. The target of any summarization system is to generate a summary S of that topic that is representative of the average review and talks to its significant aspects. The summarization requires to measure with respect to the mentioned aspects: optimization, and intensity.

The optimization is measured by the equation (1):

$$argmax\ L(SM)\ , s.t.:Length(SM) \leq K \tag{1}$$

Where L is many possible scores of summaries, Length (SM) refers to all sentiments summary length and K is the pre-specified length constraint.

The intensity of summarization is declared in equation (2):

$$Intensity\ (SM) = \sum_{t \in SM} |LEX - S\ (t)| \tag{2}$$

$$S\ (SM) = \frac{\sum_{t \in SM} Lex - S\ (t)}{\alpha + intensity\ (SM)} \tag{3}$$

A central function for normalized sentiment,

## III. PROPOSED CLUSTERING TECHNIQUE

This study focuses on four kinds of factors such as sentimental process, personal style process, linguistic style, and negative sentiment depression level. The proposed technique consists of two algorithms clustering for user's sex and classification algorithm for causes and responsibilities of women. There are some features extracted from this classification that can provide results by identifying the social status and financial status for each woman that can be a parameter in the relationship between depression and causes. The proposed classification algorithm is based on six classes (education problems, work problems, family problems, baby problems, disease problems, and shape appearance/personal problems). The social status (single, married, divorced, married and has children, widowed). This algorithm shows the evaluation and classification for each sentiment based on deep neural algorithms (word2vector embedding algorithm and recurrent neural network). Word embedding that is introduced by word2vec is generally used to learn context and produce high-dimensional vectors in a space. These embedding are then classified using the machine-learning algorithm.

These features can be analyzed and detect the depressive information for each sentiment that enable to recognize some new clusters that can detect the writer' sex based on the sentiment written style and some causes classification which extracted from online sentiments data such as (causes of depression, responsibilities from each cause, social and finical status). These extracted sentiments received as Facebook posts or comments. We then apply supervised deep learning approaches to considerate each factor kinds independently. The used classification technique is Word2vec and LSTM. The clustering algorithm can produce the man, and female sentiments (see Fig. 1). This research focuses on women depression so that is not easy to cluster the sentiment based on the written requirements and written style.

Fig. 1. The Proposed Technique of SentiNeural (Holding a Proposed Clustering and Classification Algorithms).

Feature extraction support describing and demonstrating amongst depressive and non-depressive posts. Fetching many features includes linguistic specifications from user's reviews. It is clarified briefly as follows:

- Psychological analysis operation —affective process, social process, cognitive process, time orientations, relativity, personal concerns

- Part-of-speech Linguistic operation such as number of words on sentiments sentence, word, pronoun, prepositions, adverbs, conjunctions, Negations.

- Other types of language grammar such as verbs, adjectives, adverbs, comparisons, numbers, or quantifiers.

This research applies text analysis algorithms in sequenced processes as follow:

### A. Translation

This algorithm is suitable to variant language but this paper focuses on the Arabic language dataset. Translate Arabic into English process: the proposed algorithm works for English sentiments, so we will translate the Arabic language into English and works on the English sentiment. The translation has the challenge to detect language and slang words. The slang words interpretation in Arabic Language requires several motivations for reaching better results. That is considered a new research opportunity.

### B. Preprocessing Phase Remove stop Words

A stop word is a commonly used word (such as "the", "a", "an", "in") that a search engine has been programmed to ignore, both when indexing entries for searching and when retrieving them as the result of a search query. Part-of-speech tagging (POST) is considered grammatical tagging for sentiment text. Stemming: A stemming algorithm is a process of linguistic normalization, in which the variant forms of a word.

### C. Hierarchal Classification Sentiments

This clustering is a new proposed approach for sentiment analysis. It creates a hierarchal clustering model for detecting women sentiments automatically based on the sentiment style of written and a new depression classification for causes and responsibilities. For each sentiment S that requires to identify the polarity first such as (positive or negative sentiments). The technique avoids the positive sentiments and focuses on negative sentiments. The proposed technique splits negative sentiments into three levels of negative sentiment (high, medium, and low levels of depressions). Clustering [11, 12] can be roughly featured as Hard clustering: each object belongs to an exact cluster or not. Soft clustering (also: fuzzy clustering): each object belongs to each cluster to a certain degree (for example, a likelihood of belonging to the cluster).

It also can predict causes using machine learning and deep learning algorithms (word2vec embedding algorithm and Recurrent neural network). Word2vec is known aslo Word embeddings are basically a form of word representation that bridges the human understanding of language to that of a machine. Word embeddings are distributed representations of text in an n-dimensional space. These are essential for solving most NLP problems. The main goal for W*ord2vec* algorithm is applying on big data of text and generating a vector space, typically of several hundred dimensions, with each unique word in the corpus being assigned a corresponding vector in the space. Word vectors are positioned in the vector space such that words that share common contexts in the corpus are located in close proximity to one another in the space.

This cluster algorithm can recognize automatically for the sentiments user sex (females or males) and the level of depression. This research takes care of women/ female sentiments on Facebook social network which becomes one of the *biggest platform for expressing feelings and opinions simultaneously.

### D. A Depression Clustering

This cluster algorithm can recognize automatically for the sentiments user sex (females or males) and the level of depression. It is based on LSTM neural network as the following in Fig. 2.

Fig. 2.   The Neural Network Structure for the Depression Analysis.

### E. Sentiment Evaluation

That can detect the sentiment polarity to can determine the sentiment is positive or negative. We focus on this research on the negative polarity sentiments. The level of negative refers to the level of depression. Sentiment evaluation: this process requires the text analysis, interpretation, evaluate the words, sentence, and sequence of words. Sentiment polarity: is based on positive or negative polarity. Sentiment depression level: is based on negative sentiment polarity.



(a)



(b)

Fig. 3.   The Proposed Depression Detection Technique (a) The Sequenced Processes of the Technique (b) The Relationship between Sentiment Analysis, Summary and the Results of Depression Level.

TABLE I.        THE PROPOSED CLUSTERING SENTIMENT BASED ON NEURAL ALGORITHM

| A Proposed Algorithm of : for Arabic Sentiments | |
|---|---|
| **Input** | A corpus S of N number of Sentiments {s1, s2, . . . , sN} |
| **Output** | Assign a Positive OR Negative label For each sentiment si ∈ S, (i =1, 2, . . . , m) |
| **Pre-processing:** | 1: for all sentiments sj ∈ S do<br>2: for all each token ti ∈ Sj do  (token refers to the word)<br>3: Tag ti with part of speech tagging ptj<br>4: if ptj == %x OR ptj == x%, x donates searching to the infinitive verb from the token whether noun, verb, adjective and adverb respectively<br> then<br>5: Keep pti<br>6: Add pti to F, F is the features set<br>7: else 8: Remove ti<br>9: end if<br>10: end for<br>11: end for |
| | **Clustering**:<br>12: Set the clusters number C = 2<br>13: for all matrix files Mi , (i = 1, 2, . . . , n), do<br> 14: using Deep neural network algorithm, learning the features and properties of women feelings problems based on word2vec embedding algorithm.<br>15. Determine the negative polarity classification for the translated sentiment.<br>16: if sj ∈ Ns do, Ns refers to Negative sentiment<br>17: Cluster Lj into three levels of depression based on recurrent neural network algorithm L1,L2,L3 |
| **Evaluation** | 18. Evaluate the Sp into three classifications of negative sentiment polarity from -1, to 0, -1 refers to high depression, -.05 medium depression, and 0 low depression. |
| **A proposed clustering algorithm for Facebook 'users** | 19. Cluster sj into two user's sex  male M or female F<br>20: Predict sj cluster having depression or not<br>21: Predict depression level into three clusters (high, medium, low)<br>22: endif<br>23: end for |
| **A proposed classification algorithm for reasons of depression** | 24. Classify sj into five social clusters G1, G2,G3,G4, and G5<br>25: Predict sj class three financial for each sentiment<br>26: Predict class the cause problem ci (from six classes)<br>27: Predict the responsibility class<br>28: endif<br>29: end for |
| **Poisson Distribution** | 30. counting data using poison regression, which refers to the probability of events for a Poisson distribution<br><br>$P(X = x) = \frac{\lambda^x e^{-\lambda}}{x!}$           (1)<br><br>, λ is the mean number of events x=0,1, 2.... |

### F. Sentiment Summarization

It targets a summary of causes of depression based on the level of depression and responsibilities.

### G. Poisson Distribution Phase

This distribution targets the determination of the number of depression women from the collected dataset. It aims to discuss the number for occurrences the depression in other conditions (multi-clustering) e.g. high depression, married, and cancer. A simple flowchart for the phases of the proposed framework is presented in Fig. 3, while Table I shows the pseudocode of the proposed framework.

## IV. EXPERIMENT AND RESULTS

This experiment applies on Arabic sentiments for women on social media.

### A. Preparing the Dataset

The dataset is manually collected from women groups on Facebook. We observe the users' reviews for depressive behavioral inspection and detection. The social network includes the other faces and depth emotions of people that can express about them in a freeway. Preparing of social network data, in particular Facebook writer's reviews is one of the main obstacles which give information on whether or not they could contain depression producing content. That is not easy to use any extractor program for recognizing the Arabic sentiment from women. So the preparing the data is based on the clustering sex male or women online. We depend on Facebook 'groups in Egypt.

After gathering the raw data from Facebook, the analyzing reviews process is applied on the sentiment analysis word level. It can interpret text and prepare it for evaluation. Our primary dataset contains 10.000 sentiments (7000 negatives and 3000 positives) from Facebook from groups to ensure clustering women from without examining their profiles.

Table II shows a sample of women sentiments in the first phase (Translation from Arabic into English), and Table III shows SentiNeural results in processing in the Multi-clustering based on deep neural networks for achieving the causes of depression based on clustering (financial status, social status, and responsibilities).

### B. Accuracy of Sentiment Analysis Polarity

The women sentiment will be evaluated accuracy in two times. First time to know the positive or negative sentiments is shown in Table IV. The second time is evaluated for the classification levels of negative as shown in Table V.

$$Accuracy = \frac{tp+tn}{tp+tn+fp+fn} \tag{4}$$

TABLE II. WOMEN SENTIMENTS TRANSLATION FROM ARABIC TO ENGLISH

| No | Sentiment | Translate |
|----|-----------|-----------|
| 1. | المدير متحكم وودني وكل شوية يقلل من شغلى مش عارفة اثبت نفسي أزاي | The manager is controlled and Woody and every angle reduces my job not knowing how to prove myself |
| 2. | الشغل في الشركة صعب جدا مش بعرف اخد بريك خالص وبتاخر ل 9 او ل كل يوم103 بالي | The work in this company is very heavy I can't take a break and must to go out around 9 or 10 daily |
| 3. | انا تخينة اوي وكل الناس بتقول عليا وحشة وزني فوق ال130 كيلو مش عارفة اعمل ايه | I am overweight more than 130 kilo, I don't know what can I do |
| 4. | لو سمحتوا ابني تعبان اوي وسخن اعمل ايه مش معايا فلوس اروح لدكتور | If you allow my son is very tired and hot I do not know how to make money |
| 5. | من وقت ما اتطلقت وامي تضربني وتهيني كأني عالة عليهم | From time to time my mother tasted me and humiliated me as if I was a burden on them |

TABLE III. MULTI-CLUSTERING FOR PREDICTING THE CAUSES OF DEPRESSION INTO TO LEVELS OF CLUSTERING FROM EXTRACTING FEATURES IN WOMEN SENTIMENTS FOR THE PREVIOUS EXAMPLES IN TABLE I

| No. | Social status cluster | Financial status cluster | Causes | Responsibilities |
|-----|----------------------|--------------------------|--------|------------------|
| 1. | Single | Medium | Manager problem | Work |
| 2. | Single | Medium | Work dates and delays | Work |
| 3. | Single | Medium | The shape and weight | Non |
| 4. | Married and has child | Low | Baby patient | Baby |
| 5. | Divorce | Low | Divorce and mother problems | Divorce |

TABLE IV. THE ACCURACY OF POSITIVE OR NEGATIVE POLARITY

|  | Precision | recall | F1-score |
|--|-----------|--------|----------|
| **Negative** | 0.97 | 0.85 | 0.91 |
| **Positive** | 0.67 | 0.80 | 0.83 |
| **Average** | 0.85 | 0.82 | 0.83 |

TABLE V. THE ACCURACY FOR NEGATIVE LEVELS

|  | Precision | recall | F1-score |
|--|-----------|--------|----------|
| **High negative** | 0.77 | 0.65 | 0.71 |
| **Medium negative** | 0.41 | 0.52 | 0.46 |
| **Low negative** | 0.47 | 0.60 | 0.53 |
| **Average** | 0.65 | 0.62 | 0.63 |

Tp refers to true positive, Fn: false negative, Tn: true negative, and Fn: false negative

### C. Clustering Machine Learning and Poisson Distribution

The Poisson regression algorithm [15] is used to predict numbers based on regression models. The response variable has a Poisson distribution. The distribution of regression refers to a discrete distribution that is a method with non-whole numbers. According to Fig. 4 and Fig. 5, the results are found: Mean Absolute Error= 49958.764234, Root Mean Squared Error = 52002.038613, Relative Absolute Error =3.996701, Relative Squared Error= 12.980218, and Coefficient of Determination= -11.980218.

These results discuss the relationship between responsibilities, or causes, financial status/level, and diseases. Note: without caring the age because there are responsibilities can be in a different age, (for example, marriage age, work age, experience age, born children age, poor age, and disease age).

(a)



(b)

Fig. 4. (a) The Poisson Regression Illustrates the Relationship between the Levels of Depression and the Causes.(b) The Poisson Regression Illustrates the Relationship between the Levels of Depressions and the Predicting Social Status.



(a)



(b)

Fig. 5. (a) The Poisson Regression Distribution Shows for the Hierarchal Clustering of Translated Sentiments into Causes Clustering Dimension (Education Problems, Work Problems, Family Problems, Baby Problems, Disease Problems, and Shape Appearance Problems). (b) The Poisson Regression Distribution Shows for the Highest Rate of Family Problems Clustering of Translated Sentiments (Single, Married, Widow, and Divorce).

## V. CONCLUSION

Online sentiments have a big effect in decision making in business. There are several challenges in analyzing and evaluating sentiments. More than 60% of sentiments face a negative polarity challenge. This paper proposes SentiNeural which is a new clustering and evaluating online women sentiments from Facebook. This technique targets clustering the user's sentiments based on the text of sentiments style and clustering the level of depression based on cause's classification algorithm. SentiNeural introduces a new classification algorithm for extracting the causes, responsibilities, financial status, and social status for each women using deep neural algorithms. It also includes a

translation from Arabic to English languages and shows a summary for each sentiment. The experiment of the proposed framework relies on Poisson regression distribution that can determine the number of sentiments in (high, medium, and low) depression according to the different causes (family, education, work, diseases, baby, shape/personal, problems) with respect to the prediction of two levels finical and social status (married, single, widow, divorce). The accuracy reaches in average between 85% into 91% that is based on translation and summarization results. For future work, improving the translation and summarization algorithms for achieving 98% in accuracy. Further, another future work targets Appling the same algorithm in various languages and providing some solution for Arabic language challenges.

REFERENCES

[1] Egypt Independent, "7% of Egyptians suffer from mental illness: Health Ministry," Egypt Independent, 2018. [Online]. Available: https://egyptindependent.com/7-of-egyptians-suffer-from-mental-illness-health-ministry/. [Accessed: 18-May-2019].

[2] "Depression and Other Common Mental Disorders : Global Health Estimates," Organization, World Health, 2018. [Online]. Available: https://apps.who.int/iris/bitstream/handle/10665/254610/WHO-MSD-MER-2017.2-eng.pdf?sequence=1. [Accessed: 18-May-2018].

[3] N. E. M. Khalifa, M. H. N. Taha, S. H. N. Taha, and A. E. Hassanien, "Statistical Insights and Association Mining for Terrorist Attacks in Egypt," 2019, pp. 291–300.

[4] E. Ortiz-Ospina and S. Tzvetkova, "Working women: Key facts and trends in female labor force participation," Our World in Data, 2017. [Online]. Available: https://ourworldindata.org/female-labor-force-participation-key-facts. [Accessed: 18-May-2019].

[5] T. Chen, D. Borth, T. Darrell, and S.-F. Chang, "Deepsentibank: Visual sentiment concept classification with deep convolutional neural networks," arXiv Prepr. arXiv1410.8586, 2014.

[6] W. B. Gevarter, "Introduction to Artificial Intelligence.," Chem. Eng. Prog., 1987.

[7] Y. Baştanlar and M. Özuysal, "Introduction to machine learning," Methods Mol. Biol., 2014.

[8] F. Q. Lauzon, "An introduction to deep learning," in 2012 11th International Conference on Information Science, Signal Processing and their Applications, ISSPA 2012, 2012.

[9] M. De Choudhury, M. Gamon, S. Counts, and E. Horvitz, "Predicting depression via social media.," in Proceedings of the Seventh International AAAI Conference on Weblogs and Social Media 128, 2013, pp. 128–137.

[10] G. Shen et al., "Depression Detection via Harvesting Social Media: A Multimodal Dictionary Learning Solution," in Proceedings of the 26th International Joint Conference on Artificial Intelligence, 2017, pp. 3838–3844.

[11] Y. Cao, M. Huang, and X. Zhu, "Clustering sentiment phrases in product reviews by constrained co-clustering," in Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics), 2015.

[12] N. Eldeen, M. Khalifa, M. Hamed, N. Taha, and A. E. Hassanien, "Aquarium Family Fish Species Identification System Using Deep Neural Networks," 2018.

[13] N. E. Khalifa, M. Hamed Taha, A. E. Hassanien, and I. Selim, "Deep galaxy V2: Robust deep convolutional neural networks for galaxy morphology classifications," in 2018 International Conference on Computing Sciences and Engineering, ICCSE 2018 - Proceedings, 2018, pp. 1–6.

[14] P. Beineke, T. Hastie, C. Manning, and S. Vaithyanathan, "Exploring Sentiment Summarization," Proc. AAAI Spring Symp. Explor. Attitude Affect Text Theor. Appl., 2004.

[15] P. Sedgwick, "Poisson regression," BMJ (Online). 2014.

# Bayesian Network Analysis for the Questionnaire Investigation on the Needs at Fuji Shopping Street Town under the View Point of Service Engineering

Tsuyoshi Aburai[1]
Tokushima University
Tokushima, JAPAN

Akane Okubo[2]
University of Shizuoka, Japan
Shizuoka, JAPAN

Daisuke Suzuki[3]
Fujisan Area Management Company
Shizuoka, JAPAN

Kazuhiro Takeyasu[4]
College of Business Administration, Tokoha University,
Japan, Shizuoka, JAPAN

*Abstract*—**Shopping streets at local city in Japan became old and are generally declining. In this paper, the area rebirth and/or regional revitalization of shopping street are handled. Fuji city in Japan is focused. Four big festivals are held at Fuji city (two for Fuji Shopping Street Town and two for Yoshiwara Shopping Street Town). Many people visit these festivals including residents in that area. Therefore a questionnaire investigation to the residents and visitors is conducted during these periods in order to clarify residents and visitors' needs for the shopping street, and utilize them to the plan building of the area rebirth and/or regional revitalization of shopping street. There is a big difference between Fuji Shopping Street Town and Yoshiwara Shopping Street Town. Therefore Fuji Shopping Street Town is focused in this paper. These are analyzed by using Bayesian Network. These are analyzed by sensitivity analysis and odds ratio is calculated to the results of sensitivity analysis in order to obtain much clearer results. The analysis utilizing Bayesian Network enabled us to visualize the causal relationship among items. Furthermore, sensitivity analysis brought us estimating and predicting the prospective visitors. Sensitivity analysis is performed by back propagation method. These are utilized for constructing a much more effective and useful plan building. Fruitful results are obtained. To confirm the findings by utilizing the new consecutive visiting records would be the future works to be investigated.**

*Keywords*—*Fuji city; area rebirth; regional vitalization; Bayesian network; back propagation; service engineering*

## I. INTRODUCTION

Shopping streets at local city in Japan are generally declining. It is because most of them were built in the so-called "High Growth Period (1954-1973)". Therefore they became old and area rebirth and/or regional revitalization are required everywhere.

There are many papers published concerning area rebirth or regional revitalization. Author in [1] has pointed out the importance of tourism promotion. Author in [2] developed the project of shutter art to Wakkanai Chuo shopping street in Hokkaido, Japan. Author in [3] has made a questionnaire research at Jigenji shopping street in Kagoshima Prefecture,

Japan and analyzed the current condition and future issues. For about tourism, many papers are presented from many aspects as follows.

Author in [4] designed and conducted a visitor survey on the spot, which used a questionnaire to investigate the activities of visitors to the Ueno district in Taito ward, Tokyo. Author in [5] analyzed the image of the Izu Peninsula as a tourist destination in their 2003 study "Questionnaire Survey on the Izu Peninsula." Author in [6] conducted tourist behavior studies in Atami city in 2008, 2009, 2014 and in other years.

In this paper, the area rebirth and/or regional revitalization of shopping street are handled. Fuji city in Japan is focused. Fuji city is located in Shizuoka Prefecture. Mt. Fuji is very famous all around the world and its beautiful scenery from Fuji city can be seen, which is at the foot of Mt. Fuji. There are two big shopping streets in Fuji city. One is Yoshiwara shopping street and another one is Fuji shopping street. They became old and building area rebirth and regional revitalization plan have started. Following investigation was conducted by the joint research group (Fuji Chamber of Commerce & Industry, Fujisan Area Management Company, Katsumata Maruyama Architects, Kougakuin University and Tokoha University). The main project activities are as follows:

- Investigation on the assets which are not in active use

- Questionnaire Investigation to Entrepreneur

- Questionnaire Investigation to the residents and visitors

After that, area rebirth and regional revitalization plan were built.

In this paper, above stated C is handled.

Four big festivals are held at Fuji city. Two big festivals are held at Yoshiwara Shopping Street Town and two big festivals at Fuji Shopping Street Town.

At Yoshiwara Shopping Street Town, Yoshiwara Gion Festival is carried out during June and Yoshiwara Shukuba (post-town) Festival is held during October. On the other hand,

Kinoene Summer Festival is conducted during August and Kinoene Autumn Festival is performed during October at Fuji Shopping Street Town. Many people visit these festivals including residents in that area.

Therefore questionnaire investigation of C is conducted during these periods.

Finally, 982 sheets (Yoshiwara Shopping Street Town: 448, Fuji Shopping Street Town: 534) were obtained.

Basic statistical analysis and Bayesian Network analysis are executed based on that. This is really a quite new approach in this field and there is no related paper on this theme as far as searched.

In recent years, the Bayesian network is highlighted because it has the following good characteristics (Neapolitan, 2004).

- Structural Equation Modeling requires normal distribution to the data in the analysis. Therefore, it has a limitation in making analysis, but the Bayesian network does not require a specific distribution type to the data. It can handle any distribution type.

- It can handle the data which include partial data.

- Expert's know-how can be reflected in building a Bayesian Network model.

- Sensitivity analysis can be easily performed by settling evidence. The prospective purchaser can be estimated and predicted by that analysis.

- It is a probability model having a network structure. Related items are connected with directional link. Therefore, understanding becomes easy by its visual chart.

The field of service marketing generally handles the shapeless.

Therefore it is often the case that it is hard to catch the influence to consumers.

Bayesian Network analysis enables to visualize the relationship and/or influence of shapeless products to consumers which is the field of service marketing.

These are also applied to service engineering.

In this paper, a questionnaire investigation is executed in order to clarify residents and visitors' needs for the shopping street and utilize them to the plan building of the area rebirth and/or regional revitalization of shopping street. There is a big difference between Fuji Shopping Street Town and Yoshiwara Shopping Street Town. Therefore Fuji Shopping Street Town is focused in this paper. These are analyzed by using Bayesian Network. These are analyzed by sensitivity analysis and odds ratio is calculated to the results of sensitivity analysis in order to obtain much clearer results. By that model, the causal relationship is sequentially chained by the characteristics of visitors, the purpose of visiting and the image of the surrounding area at this shopping street. The analysis utilizing Bayesian Network enabled us to visualize the causal relationship among items. Furthermore, sensitivity analysis brought us estimating and predicting the prospective visitors. Sensitivity analysis was conducted by back propagation method.

Some interesting and instructive results are obtained.

The rest of the paper is organized as follows. Outline of questionnaire investigation is stated in Section 2. In Section 3, Bayesian Network analysis is executed which is followed by the sensitivity analysis in Section 4. Conclusion is stated in Section 5.

## II. OUTLINE AND THE BASIC STATISTICAL RESULTS OF THE QUESTIONNAIRE RESEARCH

### A. Outline of the Questionnaire Research

A questionnaire investigation to the residents and visitors is conducted during these periods in order to clarify residents and visitors' needs for the shopping street, and utilize them to the plan building of the area rebirth and/or regional revitalization of shopping street. The outline of questionnaire research is as follows. Questionnaire sheet is attached in Appendix 1.

| (1) | Scope of investigation | : | Residents and visitors who have visited four big festivals at Fuji city in Shizuoka Prefecture, Japan |
| (2) | Period | : | Yoshiwara Gion Festival: June 11,12/2016 Yoshiwara Shukuba (post-town) Festival: October 9/2016 Kinoene Summer Festival: August 6,7/2016 Kinoene Autumn Festival: October 15,16/2016 |
| (3) | Method | : | Local site, Dispatch sheet, Self writing |
| (4) | Collection | : | Number of distribution 1400 Number of collection 982(collection rate 70.1%) Valid answer 982 |

### B. Basic Statistical Results

Now, the main summary results by single variable are shown.

#### 1) Characteristics of answers
   *a) Sex (Q7):* Male 43.3%, Female 56.7%
These are exhibited in Fig. 1.



Fig. 1. Sex (Q7).

*b) Age (Q8):* 10th 20.6%, 20th 16.7%, 30th 25.3%, 40th 17.0%, 50th 10.1%, 60th 6.9%, More than 70 3.4%

These are exhibited in Fig. 2.

*c) Residence (Q9):* a. Fuji city 82.8%, b. Fujinomiya city 8.8%, c. Numazu city 2.1%, d. Mishima city 0.7%, e. Shizuoka city 0.9%, F. Else (in Shizuoka Prefecture) 2.1%, g. Outside of Shizuoka Prefecture 2.6%

These are exhibited in Fig. 3.

*d)* How often do you come to this shopping street? (Q1)

Everyday 21.2%, More than 1 time a week 17.2%, More than 1 time a month 22.7%,

More than 1 time a year 26.8%, First time 3.0%, Not filled in 4.1%

These are exhibited in Fig. 4.

*e)* What is the purpose of visiting here? (Q2)

Shopping 17.2%, Eating and drinking 13.6%, Business 7.4%, Celebration, event 34.1%,

Leisure, amusement 6.1%, miscellaneous 21.6%

These are exhibited in Fig. 5.



Fig. 4.    How often do you Come to this Shopping Street? (Q1).



Fig. 5.    What is the Purpose of Visiting here? (Q2).

*f)* How do you feel about the image of the surrounding area at this shopping street? (Q3)

Beautiful 51.2%, Ugly 48.8%, of the united feeling there is 44.3%, Scattered 55.7%,

Varied 38.5%,Featureless 61.5%, New 37.1%, Historic 62.9%, Full of nature 37.1%,Urban 62.9%,

Cheerful 44.1%, Gloomy 55.9%, Individualistic 42.0%, Conventional 58.0%, Friendly 57.8%,

Unfriendly 42.2%, Healed 53.3%, Stimulated 46.7%, Open 44.8%, exclusive 55.2%, want to reside 43.6%,

Do not want to reside 56.4%, Warm 55.1%, Aloof 44.9%, Fascinating 42.1%, not fascinating 57.9%,

Want to play 47.1%, Want to examine deliberately 52.9%, lively 36.8%, Calm 63.2%,

Atmosphere of urban 28.0%, Atmosphere of rural area 72.0%

These are exhibited in Fig. 6.

*g)* There are many old building at the age of nearly 50 years. Do you think we can still use them? (Q4)

Can use it 48.7%, Cannot use it 29.2%, Have no idea 22.1%

These are exhibited in Fig. 7.



Fig. 2.    Age (Q8).



Fig. 3.    Residence (Q9).

Fig. 6. How do you Feel about the Image of the Surrounding Area at this Shopping Street? (Q3).



Fig. 7. There are Many Old Building at the Age of Nearly 50 years. Do you think we Can Still use them? (Q4).

## III. BAYESIAN NETWORK ANALYSIS

In constructing Bayesian Network, it is required to check the causal relationship among groups of items.

BAYONET software (http://www.msi.co.jp/BAYONET/) is used. When plural nodes exist in the same group, it occurs that causal relationship is hard to set a priori. In that case, BAYONET system set the sequence automatically utilizing AIC standard. Node and parameter of Fig. 8 are exhibited in Table I.

In the next section, sensitivity analysis is achieved by back propagation method. Back propagation method is conducted in the following method (Fig. 9).



Fig. 8. A Built Model.

$$Pr(X = x) = \alpha\lambda(x)\pi(x)$$

$$\pi(x) = \sum_u P(x|U = u)\prod_{U_i}\pi_{U_i X}(u)$$

$$\lambda(x) = \prod_{Y_j}\lambda_{Y_j X}(x)$$

$$\pi_{XY_j}(x) = \pi(x)\prod_{k \neq j}\lambda_{Y_k X}(x)$$

$$\lambda_{XU_i}(u) = \sum_x \lambda(x)\sum_{k \neq i} P(x|U)\prod_{k \neq i}\pi_{U_k X}(u_k)$$



Fig. 9. Back Propagation Method (Takeyasu et al., 2010).

TABLE I.        NODE AND PARAMETER

| Node | Parameter | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| Gender | Male | Female | | | | | | | | |
| Age | 10th | 20th | 30th | 40th | 50th | 60th | More than 70 | | | |
| The purpose of visiting | Shopping | Eating and drinking | Business | Celebration、event | Leisure, amusement | miscellaneous | | | | |
| The image of the surrounding area at this shopping street | Beautiful | Ugly | Of the united feeling there is | Scattered | Varied | Featureless | New | Historic | Full of nature | Urban |

| Node | Parameter | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
| The image of the surrounding area at this shopping street | Cheerful | Gloomy | Individualistic | Conventional | Friendly | Unfriendly | Healed | Stimulated | Open | Exclusive |

| Node | Parameter | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 |
| The image of the surrounding area at this shopping street | Want to reside | Do not want to reside | Warm | Aloof | Fascinating | Not fascinating | Want to play | Want to examine deliberately | Lively | Calm |

| Node | Parameter | |
|---|---|---|
| | 31 | 32 |
| The image of the surrounding area at this shopping street | Atmosphere of urban | Atmosphere of rural area |

## IV. SENSITIVITY ANALYSIS

Now, posterior probability is calculated by setting evidence as, for example, 1.0. Comparing Prior probability and Posterior probability, the change can be seen and the preference or image of the surrounding area at this shopping street can be confirmed. Evidence is set to all parameters. Therefore the analysis volume becomes too large. In this paper, nearly 1/3 of the total cases are picked up and analysis is executed. Nodes that are analyzed here are "Gender", "Age" and "The purpose of visiting". Another paper for the rest of them is prepared.

As stated above, evidence is set to each parameter, and the calculated posterior probability is exhibited in Appendix 2 which includes the calculation results of odds ratio.

Here, each item is classified by the strength of the odds ratio.

- Very Strong (+++): Select major parameter of which the odds ratio is more than 1.6

- Strong (++): Select major parameter of which the odds ratio is more than 1.3

- Medium (+): Select major parameter of which the odds ratio is more than 1.08

- Weak: Else

Now each of them is examined for Very Strong, Strong and Medium case.

### A. Sensitivity Analysis for "The Purpose of Visiting"

*1) Setting evidence to "Shopping":* After setting evidence to "Shopping", the result is exhibited in Table II.

Those who visit for "Shopping" had come with the purpose of visiting for "Leisure, amusement" of an age of "20th", "60th" or "*More than 70*" in which the gender is "*Female*".

(Very Strong part is indicated by bold character and Strong is indicated by italic.)

*2) Setting evidence to "Eating and drinking":* After setting evidence to "Eating and drinking", the result is exhibited in Table III.

Those who visit for "Eating and drinking" had come with the purpose of visiting for "Business", "Celebration、event" under the image of the surrounding area at this shopping street as "Scattered", "Conventional" or "Exclusive" of an age of "*20th*", "*40th*" or "*50th*" in which the gender is "Male".

TABLE II.        SETTING EVIDENCE TO "SHOPPING" CASE

| | |
|---|---|
| Leisure, amusement | + |
| Female | ++ |
| Age: 20th | + |
| Age: 60th | + |
| Age: More than 70 | ++ |

TABLE III.     SETTING EVIDENCE TO "EATING AND DRINKING" CASE

| | |
|---|---|
| Business | + |
| Celebration、event | + |
| Scattered | + |
| Conventional | + |
| Exclusive | + |
| Male | + |
| Age: 20th | ++ |
| Age: 40th | ++ |
| Age: 50th | ++ |

*3) Setting evidence to "Business":* After setting evidence to "Business", the result is exhibited in Table IV.

Those who visit for "Business" had come with the purpose of visiting for "Eating and drinking", "Celebration、event" under the image of the surrounding area at this shopping street as "Conventional" or "Aloof" of an age of "**20th**", "30th," or "*50th*" in which the gender is "**Male**".

*4) Setting evidence to "Celebration 、 event":* After setting evidence to "Celebration 、 event", the result is exhibited in Table V.

Those who visit for "Celebration、event" had come with the purpose of visiting for "Eating and drinking", "Business" under the image of the surrounding area at this shopping street as "Scattered", "Conventional" or "Exclusive" of an age of "30th", "40th" or "*50th*" in which the gender is "*Male*".

TABLE IV.     SETTING EVIDENCE TO "BUSINESS" CASE

| | |
|---|---|
| Eating and drinking | + |
| Celebration、event | + |
| Conventional | + |
| Aloof | + |
| Male | +++ |
| Age: 20th | +++ |
| Age: 30th | + |
| Age: 50th | ++ |

TABLE V.     SETTING EVIDENCE TO "CELEBRATION 、 EVENT" CASE

| | |
|---|---|
| Eating and drinking | + |
| Business | + |
| Scattered | + |
| Conventional | + |
| Exclusive | + |
| Male | ++ |
| Age: 30th | + |
| Age: 40th | + |
| Age: 50th | ++ |

*5) Setting evidence to "Leisure, amusement":* After setting evidence to "Leisure, amusement", the result is exhibited in Table VI.

Those who visit for "Leisure, amusement" had come with the purpose of visiting for "Shopping" under the image of the

surrounding area at this shopping street as "Unfriendly" of an age of "*60th*" or "*More than 70* "in which the gender is "*Female*".

*B. Sensitivity Analysis for "Gender"*

*1) Setting Evidence to "Male":* After setting evidence to "Male", the result is exhibited in Table VII.

Those who are "Male" had come with the purpose of visiting for "Eating and drinking", "*Business*", or "Celebration 、 event" under the image of the surrounding area at this shopping street as "Gloomy", "Conventional" or "Aloof".

*2) Setting Evidence to "Female":* After setting evidence to "Female", the result is exhibited in Table VIII.

Those who are "Female" had come with the purpose of visiting for "Shopping", or "Leisure, amusement" under the image of the surrounding area at this shopping street as "Beautiful", "New", "Full of nature", "Cheerful", "Individualistic", "Warm" or "Want to play".

TABLE VI.     SETTING EVIDENCE TO "LEISURE, AMUSEMENT" CASE

| | |
|---|---|
| Shopping | + |
| Unfriendly | + |
| Female | ++ |
| Age: 60th | ++ |
| Age: More than 70 | ++ |

TABLE VII.     SETTING EVIDENCE TO "MALE" CASE

| | |
|---|---|
| Eating and drinking | + |
| Business | ++ |
| Celebration、event | + |
| Gloomy | + |
| Conventional | + |
| Aloof | + |

TABLE VIII.     SETTING EVIDENCE TO "FEMALE" CASE

| | |
|---|---|
| Shopping | + |
| Leisure, amusement | + |
| Beautiful | + |
| New | + |
| Full of nature | + |
| Cheerful | + |
| Individualistic | + |
| Warm | + |
| Want to play | + |

*C. Sensitivity Analysis for "Age"*

*1) Setting evidence to "10th":* After setting evidence to "10th", the result is exhibited in Table IX.

Those who are at the age of "10th" had come under the image of the surrounding area at this shopping street as "*Beautiful*", "*Of the united feeling there is*", "*Varied*", "*Full of nature*", "Urban", "*Cheerful*", "**Individualistic**", "**Friendly**",

"*Healed*", "**Open**", "*Want to reside*", "**Warm**", "**Fascinating**", "**Want to play**" or "*Lively*".

*2) Setting evidence to "20th":* After setting evidence to "20th", the result is exhibited in Table X.

Those who are at the age of "20th" had come with the purpose of visiting for "*Shopping*", "*Eating and drinking*" or "**Business**" under the image of the surrounding area at this shopping street as "Beautiful", "New", "Full of nature", "*Cheerful*", "Conventional", "Healed", "Stimulated", "Open", "Want to reside", " Fascinating", "Want to play", "Want to examine deliberately" or "Lively".

TABLE IX.        SETTING EVIDENCE TO "10TH" CASE

| Beautiful | ＋＋ |
|---|---|
| Of the united feeling there is | ＋＋ |
| Varied | ＋＋ |
| Full of nature | ＋＋ |
| Urban | ＋ |
| Cheerful | ＋＋ |
| Individualistic | ＋＋＋ |
| Friendly | ＋＋＋ |
| Healed | ＋＋ |
| Open | ＋＋＋ |
| Want to reside | ＋＋ |
| Warm | ＋＋＋ |
| Fascinating | ＋＋＋ |
| Want to play | ＋＋＋ |
| Lively | ＋＋ |

TABLE X.        SETTING EVIDENCE TO "20TH" CASE

| Shopping | ＋＋ |
|---|---|
| Eating and drinking | ＋＋ |
| Business | ＋＋＋ |
| Beautiful | ＋ |
| New | ＋ |
| Full of nature | ＋ |
| Cheerful | ＋＋ |
| Conventional | ＋ |
| Healed | ＋ |
| Stimulated | ＋ |
| Open | ＋ |
| Want to reside | ＋ |
| Fascinating | ＋ |
| Want to play | ＋ |
| Want to examine deliberately | ＋ |
| Lively | ＋ |

*3) Setting evidence to "30th":* After setting evidence to "30th", the result is exhibited in Table XI.

Those who are at the age of "30th" had come with the purpose of visiting for "Business" or "Celebration、event" under the image of the surrounding area at this shopping street as "Conventional" or "Want to play".

*4) Setting evidence to "40th":* After setting evidence to "40th", the result is exhibited in Table XII.

Those who are at the age of "40th" had come with the purpose of visiting for "*Eating and drinking*" or "*Celebration 、event*" under the image of the surrounding area at this shopping street as "Scattered", "Featureless", "New", "Gloomy", "*Exclusive*", "*Do not want to reside*", "Aloof", "Not fascinating", "Calm", "*Atmosphere of urban*" or "Atmosphere of rural area".

TABLE XI.        SETTING EVIDENCE TO "30TH" CASE

| Business | ＋ |
|---|---|
| Celebration、event | ＋ |
| Conventional | ＋ |
| Want to play | ＋ |

TABLE XII.        SETTING EVIDENCE TO "40TH" CASE

| Eating and drinking | ＋＋ |
|---|---|
| Celebration、event | ＋＋ |
| Scattered | ＋ |
| Featureless | ＋ |
| New | ＋ |
| Gloomy | ＋ |
| Exclusive | ＋＋ |
| Do not want to reside | ＋＋ |
| Aloof | ＋ |
| Not fascinating | ＋ |
| Calm | ＋ |
| Atmosphere of urban | ＋＋ |
| Atmosphere of rural area | ＋ |

TABLE XIII.        SETTING EVIDENCE TO "50TH" CASE

| Eating and drinking | ＋＋ |
|---|---|
| Business | ＋＋ |
| Celebration event | ＋＋＋ |
| Ugly | ＋＋＋ |
| Scattered | ＋＋＋ |
| Featureless | ＋ |
| Urban | ＋ |
| Gloomy | ＋＋ |
| Individualistic | ＋ |
| Conventional | ＋ |
| Unfriendly | ＋＋ |
| Stimulated | ＋＋ |
| Exclusive | ＋＋ |
| Aloof | ＋＋ |
| Not fascinating | ＋＋ |
| Calm | ＋ |
| Atmosphere of urban | ＋ |
| Atmosphere of rural area | ＋ |

*5) Setting evidence to "50th":* After setting evidence to "50th", the result is exhibited in Table XIII.

Those who are at the age of "50th" had come with the purpose of visiting for "*Eating and drinking*", "*Business*" or "**Celebration、event**" under the image of the surrounding area at this shopping street as "**Ugly**", "**Scattered**",

"Featureless", "Urban", "*Gloomy*", "Individualistic", "Conventional", "*Unfriendly*", "*Stimulated*", "*Exclusive*", "*Aloof*", "*Not fascinating*", "Calm", "Atmosphere of urban" or "Atmosphere of rural area".

*6) Setting evidence to "60th":* After setting evidence to "60th", the result is exhibited in Table XIV.

Those who are at the age of "60th" had come with the purpose of visiting for "Shopping", "**Leisure, amusement**" under the image of the surrounding area at this shopping street as "**Scattered**", "**Featureless**", "New", "**Urban**", "**Gloomy**", "*Conventional*", "**Unfriendly**", "**Stimulated**", "*Exclusive*", "**Do not want to reside**", "Aloof", " *Not fascinating*", "**Want to examine deliberately**", "**Calm**" or "**Atmosphere of rural area**".

*7) Setting evidence to "More than 70":* After setting evidence to "More than 70", the result is exhibited in Table XV.

Those who are at the age of "More than 70" had come with the purpose of visiting for "**Shopping**", "Celebration、event" or "**Leisure, amusement**" under the image of the surrounding area at this shopping street as "Ugly", "Featureless", "*Historic*", "Full of nature", "**Gloomy**", "Conventional", "**Unfriendly**", "**Stimulated**", "**Exclusive**", "*Do not want to reside*", "**Aloof**", "Not fascinating", "*Want to examine deliberately*", "**Calm**" or "Atmosphere of rural area".

TABLE XIV.    SETTING EVIDENCE TO "60TH" CASE

| Shopping | + |
|---|---|
| Leisure, amusement | +++ |
| Scattered | +++ |
| Featureless | +++ |
| New | + |
| Urban | +++ |
| Gloomy | +++ |
| Conventional | ++ |
| Unfriendly | +++ |
| Stimulated | +++ |
| Exclusive | ++ |
| Do not want to reside | +++ |
| Aloof | + |
| Not fascinating | ++ |
| Want to examine deliberately | +++ |
| Calm | +++ |
| Atmosphere of rural area | +++ |

TABLE XV.    SETTING EVIDENCE TO "MORE THAN 70" CASE

| Shopping | +++ |
|---|---|
| Celebration、event | + |
| Leisure, amusement | +++ |
| Ugly | + |
| Featureless | + |
| Historic | ++ |
| Full of nature | + |
| Gloomy | +++ |
| Conventional | + |
| Unfriendly | +++ |
| Stimulated | +++ |
| Exclusive | +++ |
| Do not want to reside | ++ |
| Aloof | +++ |
| Not fascinating | + |
| Want to examine deliberately | ++ |
| Calm | +++ |
| Atmosphere of rural area | + |

## V. CONCLUSION

Shopping streets at local city in Japan became old and are generally declining. In this paper, the area rebirth and/or regional revitalization of shopping street are handled. Fuji city in Japan is focused. Four big festivals are held at Fuji city (two for Fuji Shopping Street Town and two for Yoshiwara Shopping Street Town). Many people visit these festivals including residents in that area. There is a big difference between Fuji Shopping Street Town and Yoshiwara Shopping Street Town. Therefore Fuji Shopping Street Town is focused in this paper. A questionnaire investigation to the residents and visitors is conducted during these periods in order to clarify residents and visitors' needs for the shopping street, and utilize them to the plan building of the area rebirth and/or regional revitalization of shopping street. These are analyzed by using Bayesian Network. By that model, the causal relationship is sequentially chained by the characteristics of visitors, the purpose of visiting and the image of the surrounding area at this shopping street. This is really a quite new approach in this field and there is no related paper on this theme as far as searched.

In the Bayesian Network Analysis, model was built under the examination of the causal relationship among items. These are analyzed by sensitivity analysis and odds ratio is calculated to the results of sensitivity analysis in order to obtain much clearer results. The main result of sensitivity analysis is as follows.

Those who visit for "Business" had come with the purpose of visiting for "Eating and drinking", "Celebration、event" under the image of the surrounding area at this shopping street as "Conventional" or "Aloof" of an age of "20th", "30th" or "50th" in which the gender is "Male".

Those who are "Male" had come with the purpose of visiting for "Eating and drinking", "Business", or "Celebration、event" under the image of the surrounding area at this shopping street as "Gloomy", "Conventional" or "Aloof".

Those who are at the age of "10th" had come under the image of the surrounding area at this shopping street as "Beautiful", "Of the united feeling there is", "Varied", "Full of nature", "Urban", "Cheerful", "Individualistic", "Friendly", "Healed", "Open", "Want to reside", "Warm", "Fascinating", "Want to play" or "Lively".

Those who are at the age of "50th" had come with the purpose of visiting for "Eating and drinking", "Business" or "Celebration、event" under the image of the surrounding area at this shopping street as "Ugly", "Scattered", "Featureless", "Urban", "Gloomy", "Individualistic", "Conventional", "Unfriendly", "Stimulated", "Exclusive", "Aloof", "Not fascinating", "Calm", "Atmosphere of urban" or "Atmosphere of rural area".

Those who are at the age of "More than 70" had come with the purpose of visiting for "Shopping", "Celebration、event" or "Leisure, amusement" under the image of the surrounding area at this shopping street as "Ugly", "Featureless", "Historic", "Full of nature", "Gloomy", "Conventional", "Unfriendly", "Stimulated", "Exclusive", "Do not want to reside", "Aloof", "Not fascinating", "Want to examine deliberately", "Calm" or "Atmosphere of rural area".

The analysis utilizing Bayesian Network enabled us to visualize the causal relationship among items. Furthermore, sensitivity analysis brought us estimating and predicting the prospective visitors. Sensitivity analysis was achieved by back propagation method. These are utilized for constructing a much more effective and useful plan building.

Although it has a limitation that it is restricted in the number of researches, the fruitful results could be obtained. To confirm the findings by utilizing the new consecutive visiting records would be the future works to be investigated.

REFERENCES

[1] Inoue, Akiko(2017) "Changes in Local Communities Brought by Municipal Mergers : From the Viewpoint of Tourism Promotion as the Main Industry", Bulletin of the Faculty of Regional Development Studies, Otemon Gakuin University, Vol.2, pp.1-32.

[2] Ingu, Shuzo / Uemura, Miki / Uchida, Yuka / Omiya, Misa / Miura, Taiki / Hironori, Hironori(2017)"A study on the application of geothermal power generation to local revitalization in Obama Town, Unzen City: in consideration of futurability in Obama", Environmental Science Research, Nagasaki University, 20(1), pp.51-63.

[3] Ohkubo, Yukio(2017) "Current status and problems in Jigenji-dori shopping area : from a consumer questionnaire", Bulletin of Local Research, Kagoshima International University, Vol.44 no.2 p.1 -15.

[4] Yoshida, Ituki (2009) "Consideration on the Characteristic of Visitors' Activity and the Research Method for Tourist Visitors in Urban Areas"

[5] Doi, Hideji(2009) "Evaluation of policies to build tourist destinations and statistical analysis" Nippon Hyoron Sha

[6] Kano, Michiko (2011) "Characteristic analysis of Atami tourists: Reconsideration based on data add and modify" Shizuoka Economic Research. 16 (2), p. 61-78，Shizuoka University

[7] Takeyasu, Kazuhiro et al.(2010) "Modern Marketing", Chuoukeizaisha Publishing

APPENDIX 1

Questionnaire Sheet about the Image around the Shopping Street

1. How often do you come to this shopping street?
a. Everyday  b. ( ) times a week  c. ( ) times a month  d. ( ) times a year
e. miscellaneous (                    )

2. What is the purpose of visiting here? (Plural answers allowed)
a. shopping  b. eating and drinking  c. business  d. celebration、event  e. leisure, amusement
f. miscellaneous (                    )

3. How do you feel about the image of the surrounding area at this shopping street?
Select the position

| Beautiful | · | · | · | · | · | Ugly |
|---|---|---|---|---|---|---|
| Of the united feeling there is | · | · | · | · | · | Scattered |
| Varied | · | · | · | · | · | Featureless |
| New | · | · | · | · | · | Historic |
| Full of nature | · | · | · | · | · | Urban |
| Cheerful | · | · | · | · | · | Gloomy |
| Individualistic | · | · | · | · | · | Conventional |
| Friendly | · | · | · | · | · | Unfriendly |
| Healed | · | · | · | · | · | Stimulated |
| Open | · | · | · | · | · | exclusive |

| Want to reside | · | · | · | · | · | Do not want to reside |
|---|---|---|---|---|---|---|
| Warm | · | · | · | · | · | Aloof |
| Fascinating | · | · | · | · | · | Not fascinating |
| Want to play | · | · | · | · | · | Want to examine deliberately |
| Lively | · | · | · | · | · | Calm |
| Atmosphere of urban | · | · | · | · | · | Atmosphere of rural area |

4. There are many old building at the age of nearly 50 years. Do you think we can still use them?
a. Can use it  b. Cannot use it  c. Have no idea

5. Is there any functions or facilities that will be useful?

6. Comments

7. Sex
a. Male b. Female

8. Age
a.10th  b.20th  c.30th  d.40th  e.50th  f.60th  g. More than70

9. Residence
a. Fuji City  b. Fujinomiya City  c. Numazu City  d. Mishima City  e. Shizuoka City  f. Miscellaneous in Shizuoka Prefecture
g. Outside of Shizuoka Prefecture 〔                         〕

APPENDIX 2

Calculated posterior probability

| name_fuji | state | Prior | The purpose of visiting | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Shopping | Shopping_odds | Eating and drinking | Eating and drinking_odds | Business | Business_odds | Celebration、event | Celebration、event_odds | Leisure, amusement | Leisure, amusement_odds |
| The purpose of visiting | Shopping | 0.215 | 1 | – | 0.211 | 0.976 | 0.208 | 0.964 | 0.211 | 0.981 | 0.233 | 1.114 |
| | Eating and drinking | 0.174 | 0.172 | 0.988 | 1 | – | 0.197 | 1.163 | 0.191 | 1.121 | 0.155 | 0.867 |
| | Business | 0.103 | 0.101 | 0.985 | 0.117 | 1.164 | 1 | – | 0.113 | 1.115 | 0.090 | 0.866 |
| | Celebration、event | 0.396 | 0.392 | 0.983 | 0.433 | 1.167 | 0.435 | 1.177 | 1 | – | 0.374 | 0.913 |
| | Leisure, amusement | 0.089 | 0.098 | 1.111 | 0.080 | 0.890 | 0.079 | 0.878 | 0.084 | 0.945 | 1 | – |
| The image of the surrounding area at this shopping street | Beautiful | 0.339 | 0.342 | 1.013 | 0.324 | 0.933 | 0.328 | 0.949 | 0.326 | 0.942 | 0.346 | 1.028 |
| | Ugly | 0.292 | 0.287 | 0.977 | 0.299 | 1.036 | 0.299 | 1.033 | 0.300 | 1.039 | 0.285 | 0.969 |
| | Of the united feeling there is | 0.255 | 0.251 | 0.983 | 0.239 | 0.919 | 0.241 | 0.926 | 0.240 | 0.926 | 0.253 | 0.989 |
| | Scattered | 0.381 | 0.381 | 1.000 | 0.399 | 1.081 | 0.392 | 1.048 | 0.400 | 1.084 | 0.390 | 1.039 |
| | Varied | 0.175 | 0.171 | 0.968 | 0.167 | 0.943 | 0.167 | 0.944 | 0.168 | 0.952 | 0.171 | 0.969 |
| | Featureless | 0.490 | 0.491 | 1.004 | 0.491 | 1.008 | 0.487 | 0.990 | 0.496 | 1.025 | 0.503 | 1.056 |
| | New | 0.124 | 0.128 | 1.039 | 0.129 | 1.047 | 0.124 | 1.002 | 0.127 | 1.026 | 0.128 | 1.036 |
| | Historic | 0.561 | 0.565 | 1.014 | 0.557 | 0.983 | 0.556 | 0.980 | 0.559 | 0.992 | 0.570 | 1.038 |
| | Full of nature | 0.370 | 0.374 | 1.017 | 0.350 | 0.919 | 0.358 | 0.950 | 0.355 | 0.936 | 0.381 | 1.046 |
| | Urban | 0.231 | 0.228 | 0.983 | 0.225 | 0.963 | 0.223 | 0.955 | 0.228 | 0.982 | 0.235 | 1.022 |
| | Cheerful | 0.259 | 0.259 | 1.002 | 0.251 | 0.959 | 0.249 | 0.952 | 0.244 | 0.925 | 0.249 | 0.950 |
| | Gloomy | 0.432 | 0.434 | 1.008 | 0.444 | 1.053 | 0.445 | 1.057 | 0.447 | 1.064 | 0.435 | 1.015 |
| | Individualistic | 0.238 | 0.232 | 0.964 | 0.214 | 0.869 | 0.213 | 0.866 | 0.218 | 0.891 | 0.237 | 0.994 |
| | Conventional | 0.438 | 0.440 | 1.005 | 0.471 | 1.143 | 0.479 | 1.177 | 0.466 | 1.120 | 0.432 | 0.975 |
| | Friendly | 0.443 | 0.434 | 0.966 | 0.413 | 0.883 | 0.416 | 0.897 | 0.417 | 0.900 | 0.435 | 0.967 |
| | Unfriendly | 0.236 | 0.245 | 1.047 | 0.242 | 1.032 | 0.242 | 1.030 | 0.246 | 1.053 | 0.257 | 1.122 |
| | Healed | 0.285 | 0.279 | 0.969 | 0.279 | 0.970 | 0.282 | 0.986 | 0.275 | 0.953 | 0.267 | 0.913 |
| | Stimulated | 0.180 | 0.187 | 1.050 | 0.182 | 1.016 | 0.185 | 1.036 | 0.183 | 1.022 | 0.193 | 1.091 |
| | Open | 0.257 | 0.254 | 0.984 | 0.236 | 0.894 | 0.239 | 0.911 | 0.237 | 0.900 | 0.256 | 0.995 |
| | Exclusive | 0.393 | 0.407 | 1.060 | 0.413 | 1.087 | 0.404 | 1.048 | 0.411 | 1.080 | 0.407 | 1.061 |
| | Want to reside | 0.241 | 0.243 | 1.009 | 0.230 | 0.939 | 0.231 | 0.946 | 0.230 | 0.942 | 0.246 | 1.026 |
| | Do not want to reside | 0.395 | 0.397 | 1.010 | 0.396 | 1.007 | 0.392 | 0.987 | 0.400 | 1.022 | 0.406 | 1.049 |
| | Warm | 0.398 | 0.393 | 0.980 | 0.375 | 0.907 | 0.370 | 0.889 | 0.375 | 0.907 | 0.395 | 0.988 |
| | Aloof | 0.252 | 0.254 | 1.011 | 0.264 | 1.067 | 0.269 | 1.093 | 0.265 | 1.072 | 0.251 | 0.995 |
| | Fascinating | 0.223 | 0.222 | 0.994 | 0.205 | 0.900 | 0.210 | 0.928 | 0.208 | 0.912 | 0.223 | 0.999 |
| | Not fascinating | 0.423 | 0.424 | 1.004 | 0.435 | 1.050 | 0.430 | 1.029 | 0.436 | 1.053 | 0.428 | 1.019 |
| | Want to play | 0.218 | 0.217 | 0.996 | 0.202 | 0.908 | 0.198 | 0.886 | 0.200 | 0.898 | 0.216 | 0.991 |
| | Want to examine deliberately | 0.312 | 0.321 | 1.042 | 0.314 | 1.009 | 0.312 | 0.999 | 0.313 | 1.002 | 0.330 | 1.086 |
| | Lively | 0.181 | 0.178 | 0.982 | 0.175 | 0.960 | 0.176 | 0.967 | 0.173 | 0.948 | 0.174 | 0.949 |
| | Calm | 0.520 | 0.530 | 1.041 | 0.528 | 1.035 | 0.527 | 1.030 | 0.528 | 1.035 | 0.538 | 1.076 |
| | Atmosphere of urban | 0.097 | 0.095 | 0.981 | 0.099 | 1.031 | 0.097 | 1.003 | 0.099 | 1.022 | 0.090 | 0.928 |
| | Atmosphere of rural area | 0.629 | 0.630 | 1.004 | 0.633 | 1.017 | 0.626 | 0.988 | 0.635 | 1.028 | 0.643 | 1.061 |
| Gender | Male | 0.433 | 0.364 | 0.751 | 0.485 | 1.235 | 0.556 | 1.642 | 0.492 | 1.267 | 0.285 | 0.522 |
| | Female | 0.567 | 0.636 | 1.331 | 0.515 | 0.810 | 0.444 | 0.609 | 0.508 | 0.789 | 0.715 | 1.916 |
| Age | 10th | 0.205 | 0.172 | 0.804 | 0.082 | 0.348 | 0.088 | 0.373 | 0.111 | 0.484 | 0.197 | 0.948 |
| | 20th | 0.166 | 0.203 | 1.279 | 0.219 | 1.406 | 0.256 | 1.727 | 0.169 | 1.018 | 0.124 | 0.708 |
| | 30th | 0.251 | 0.229 | 0.886 | 0.263 | 1.064 | 0.286 | 1.191 | 0.277 | 1.143 | 0.261 | 1.051 |
| | 40th | 0.170 | 0.168 | 0.984 | 0.225 | 1.414 | 0.139 | 0.786 | 0.203 | 1.246 | 0.143 | 0.813 |
| | 50th | 0.102 | 0.081 | 0.775 | 0.140 | 1.443 | 0.146 | 1.515 | 0.136 | 1.396 | 0.058 | 0.542 |
| | 60th | 0.070 | 0.079 | 1.129 | 0.051 | 0.712 | 0.053 | 0.735 | 0.066 | 0.933 | 0.133 | 2.025 |
| | More than70 | 0.035 | 0.069 | 2.023 | 0.019 | 0.535 | 0.032 | 0.920 | 0.037 | 1.061 | 0.086 | 2.571 |

The image of the surrounding area at this shopping street

| Beautiful | Beautiful_odds | Ugly | Ugly_odds | Of the united feeling there is | Of the united feeling there is_odds | Scattered | Scattered_odds | Varied | Varied_odds | Featureless | Featureless_odds | New | New_odds |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0.216 | 1.009 | 0.211 | 0.977 | 0.212 | 0.987 | 0.214 | 0.996 | 0.209 | 0.968 | 0.215 | 1.002 | 0.222 | 1.047 |
| 0.167 | 0.950 | 0.178 | 1.029 | 0.163 | 0.928 | 0.183 | 1.060 | 0.166 | 0.946 | 0.175 | 1.005 | 0.181 | 1.047 |
| 0.099 | 0.965 | 0.105 | 1.025 | 0.096 | 0.932 | 0.106 | 1.035 | 0.098 | 0.955 | 0.102 | 0.995 | 0.103 | 1.006 |
| 0.380 | 0.937 | 0.407 | 1.046 | 0.374 | 0.912 | 0.416 | 1.088 | 0.381 | 0.939 | 0.401 | 1.022 | 0.404 | 1.035 |
| 0.091 | 1.029 | 0.086 | 0.965 | 0.088 | 0.991 | 0.090 | 1.015 | 0.087 | 0.977 | 0.091 | 1.029 | 0.092 | 1.041 |
| 1 | – | 0 | 0.000 | 0.347 | 1.036 | 0.328 | 0.952 | 0.347 | 1.036 | 0.336 | 0.983 | 0.336 | 0.984 |
| 0.000 | 0.000 | 1 | – | 0.288 | 0.979 | 0.301 | 1.042 | 0.293 | 1.007 | 0.294 | 1.012 | 0.288 | 0.984 |
| 0.260 | 1.030 | 0.251 | 0.980 | 1 | – | 0 | 0.000 | 0.264 | 1.047 | 0.250 | 0.978 | 0.248 | 0.964 |
| 0.368 | 0.949 | 0.392 | 1.049 | 0 | 0.000 | 1 | – | 0.368 | 0.947 | 0.392 | 1.049 | 0.389 | 1.036 |
| 0.179 | 1.029 | 0.176 | 1.006 | 0.181 | 1.042 | 0.169 | 0.960 | 1 | – | 0 | 0.000 | 0.170 | 0.961 |
| 0.484 | 0.977 | 0.494 | 1.017 | 0.481 | 0.967 | 0.504 | 1.059 | 0 | 0.000 | 1 | – | 0.493 | 1.015 |
| 0.123 | 0.990 | 0.122 | 0.987 | 0.120 | 0.969 | 0.127 | 1.025 | 0.120 | 0.964 | 0.125 | 1.009 | 1 | – |
| 0.564 | 1.010 | 0.563 | 1.008 | 0.562 | 1.001 | 0.560 | 0.994 | 0.564 | 1.009 | 0.561 | 0.998 | 0 | 0.000 |
| 0.380 | 1.046 | 0.367 | 0.987 | 0.379 | 1.038 | 0.359 | 0.954 | 0.379 | 1.039 | 0.367 | 0.985 | 0.364 | 0.976 |
| 0.231 | 1.001 | 0.233 | 1.011 | 0.233 | 1.011 | 0.235 | 1.020 | 0.233 | 1.010 | 0.234 | 1.016 | 0.227 | 0.977 |
| 0.268 | 1.049 | 0.252 | 0.968 | 0.273 | 1.076 | 0.241 | 0.909 | 0.269 | 1.053 | 0.250 | 0.956 | 0.256 | 0.985 |
| 0.421 | 0.956 | 0.438 | 1.024 | 0.419 | 0.950 | 0.448 | 1.069 | 0.420 | 0.952 | 0.439 | 1.030 | 0.435 | 1.014 |
| 0.247 | 1.047 | 0.240 | 1.012 | 0.255 | 1.096 | 0.226 | 0.935 | 0.254 | 1.087 | 0.235 | 0.980 | 0.226 | 0.935 |
| 0.426 | 0.952 | 0.445 | 1.028 | 0.416 | 0.912 | 0.458 | 1.084 | 0.422 | 0.937 | 0.444 | 1.025 | 0.447 | 1.034 |
| 0.456 | 1.057 | 0.439 | 0.986 | 0.465 | 1.093 | 0.421 | 0.914 | 0.464 | 1.087 | 0.435 | 0.967 | 0.427 | 0.938 |
| 0.230 | 0.968 | 0.241 | 1.029 | 0.225 | 0.938 | 0.252 | 1.089 | 0.226 | 0.946 | 0.244 | 1.044 | 0.242 | 1.035 |
| 0.291 | 1.031 | 0.283 | 0.992 | 0.295 | 1.052 | 0.271 | 0.934 | 0.295 | 1.051 | 0.278 | 0.966 | 0.277 | 0.962 |
| 0.178 | 0.986 | 0.183 | 1.021 | 0.176 | 0.976 | 0.188 | 1.053 | 0.174 | 0.964 | 0.184 | 1.027 | 0.182 | 1.015 |
| 0.266 | 1.051 | 0.247 | 0.948 | 0.273 | 1.090 | 0.238 | 0.903 | 0.268 | 1.059 | 0.250 | 0.964 | 0.248 | 0.957 |
| 0.381 | 0.953 | 0.401 | 1.033 | 0.377 | 0.934 | 0.412 | 1.086 | 0.376 | 0.932 | 0.400 | 1.033 | 0.406 | 1.056 |
| 0.246 | 1.029 | 0.239 | 0.987 | 0.250 | 1.052 | 0.233 | 0.958 | 0.247 | 1.031 | 0.239 | 0.986 | 0.237 | 0.978 |
| 0.388 | 0.973 | 0.395 | 1.001 | 0.390 | 0.982 | 0.405 | 1.043 | 0.387 | 0.968 | 0.401 | 1.025 | 0.398 | 1.012 |
| 0.409 | 1.049 | 0.393 | 0.978 | 0.416 | 1.080 | 0.381 | 0.932 | 0.413 | 1.066 | 0.392 | 0.975 | 0.391 | 0.972 |
| 0.244 | 0.956 | 0.259 | 1.040 | 0.243 | 0.951 | 0.263 | 1.062 | 0.244 | 0.957 | 0.256 | 1.023 | 0.253 | 1.008 |
| 0.232 | 1.052 | 0.217 | 0.965 | 0.237 | 1.081 | 0.207 | 0.910 | 0.233 | 1.060 | 0.217 | 0.966 | 0.216 | 0.958 |
| 0.415 | 0.966 | 0.429 | 1.026 | 0.413 | 0.958 | 0.438 | 1.061 | 0.415 | 0.966 | 0.429 | 1.025 | 0.427 | 1.017 |
| 0.228 | 1.064 | 0.207 | 0.941 | 0.233 | 1.090 | 0.198 | 0.888 | 0.229 | 1.068 | 0.210 | 0.953 | 0.215 | 0.985 |
| 0.310 | 0.991 | 0.310 | 0.991 | 0.310 | 0.987 | 0.318 | 1.025 | 0.306 | 0.970 | 0.316 | 1.015 | 0.317 | 1.024 |
| 0.186 | 1.036 | 0.179 | 0.983 | 0.188 | 1.048 | 0.172 | 0.939 | 0.188 | 1.044 | 0.177 | 0.971 | 0.177 | 0.974 |
| 0.514 | 0.976 | 0.521 | 1.006 | 0.511 | 0.965 | 0.533 | 1.056 | 0.508 | 0.954 | 0.526 | 1.027 | 0.527 | 1.031 |
| 0.095 | 0.984 | 0.099 | 1.023 | 0.097 | 0.999 | 0.097 | 1.004 | 0.097 | 1.005 | 0.096 | 0.995 | 0.097 | 1.004 |
| 0.623 | 0.977 | 0.631 | 1.011 | 0.622 | 0.971 | 0.641 | 1.054 | 0.623 | 0.973 | 0.635 | 1.028 | 0.633 | 1.018 |
| 0.390 | 0.837 | 0.444 | 1.047 | 0.425 | 0.969 | 0.442 | 1.036 | 0.416 | 0.935 | 0.431 | 0.993 | 0.380 | 0.804 |
| 0.610 | 1.195 | 0.556 | 0.955 | 0.575 | 1.032 | 0.558 | 0.965 | 0.584 | 1.070 | 0.569 | 1.007 | 0.620 | 1.243 |
| 0.244 | 1.248 | 0.195 | 0.939 | 0.292 | 1.597 | 0.137 | 0.613 | 0.273 | 1.454 | 0.185 | 0.878 | 0.141 | 0.635 |
| 0.177 | 1.079 | 0.136 | 0.788 | 0.168 | 1.014 | 0.133 | 0.770 | 0.148 | 0.873 | 0.142 | 0.833 | 0.189 | 1.167 |
| 0.263 | 1.060 | 0.240 | 0.941 | 0.216 | 0.819 | 0.247 | 0.977 | 0.258 | 1.036 | 0.248 | 0.980 | 0.253 | 1.011 |
| 0.130 | 0.728 | 0.174 | 1.027 | 0.159 | 0.922 | 0.198 | 1.204 | 0.140 | 0.793 | 0.179 | 1.065 | 0.212 | 1.315 |
| 0.089 | 0.864 | 0.160 | 1.685 | 0.079 | 0.755 | 0.146 | 1.508 | 0.108 | 1.068 | 0.115 | 1.143 | 0.089 | 0.864 |
| 0.063 | 0.893 | 0.056 | 0.785 | 0.056 | 0.787 | 0.103 | 1.515 | 0.048 | 0.669 | 0.093 | 1.364 | 0.080 | 1.158 |
| 0.034 | 0.982 | 0.039 | 1.111 | 0.030 | 0.859 | 0.037 | 1.049 | 0.025 | 0.698 | 0.038 | 1.092 | 0.035 | 0.999 |

| Historic | Historic_odds | Full of nature | Full of nature_odds | Urban | Urban_odds | Cheerful | Cheerful_odds | Gloomy | Gloomy_odds | Individualistic | Individualistic_odds | Conventional | Conventional_odds |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0.216 | 1.007 | 0.217 | 1.014 | 0.211 | 0.982 | 0.214 | 0.998 | 0.215 | 1.002 | 0.208 | 0.961 | 0.214 | 0.997 |
| 0.173 | 0.990 | 0.166 | 0.942 | 0.169 | 0.963 | 0.169 | 0.966 | 0.180 | 1.039 | 0.157 | 0.884 | 0.187 | 1.091 |
| 0.102 | 0.991 | 0.099 | 0.966 | 0.099 | 0.964 | 0.099 | 0.961 | 0.106 | 1.042 | 0.092 | 0.888 | 0.112 | 1.104 |
| 0.394 | 0.995 | 0.380 | 0.935 | 0.390 | 0.975 | 0.372 | 0.906 | 0.410 | 1.062 | 0.364 | 0.873 | 0.421 | 1.109 |
| 0.090 | 1.014 | 0.092 | 1.039 | 0.090 | 1.018 | 0.085 | 0.959 | 0.089 | 1.005 | 0.088 | 0.993 | 0.087 | 0.979 |
| 0.341 | 1.006 | 0.349 | 1.044 | 0.339 | 0.999 | 0.352 | 1.056 | 0.331 | 0.963 | 0.352 | 1.058 | 0.330 | 0.959 |
| 0.293 | 1.005 | 0.290 | 0.989 | 0.295 | 1.013 | 0.284 | 0.964 | 0.296 | 1.019 | 0.294 | 1.011 | 0.296 | 1.022 |
| 0.255 | 1.001 | 0.260 | 1.030 | 0.257 | 1.011 | 0.268 | 1.071 | 0.247 | 0.961 | 0.273 | 1.098 | 0.241 | 0.931 |
| 0.380 | 0.996 | 0.369 | 0.953 | 0.386 | 1.025 | 0.354 | 0.890 | 0.395 | 1.062 | 0.361 | 0.919 | 0.398 | 1.075 |
| 0.176 | 1.005 | 0.180 | 1.030 | 0.177 | 1.009 | 0.182 | 1.046 | 0.170 | 0.967 | 0.187 | 1.080 | 0.169 | 0.955 |
| 0.489 | 0.998 | 0.485 | 0.982 | 0.496 | 1.025 | 0.473 | 0.935 | 0.497 | 1.032 | 0.482 | 0.970 | 0.496 | 1.027 |
| 0 | 0.000 | 0.122 | 0.983 | 0.122 | 0.983 | 0.123 | 0.992 | 0.125 | 1.010 | 0.117 | 0.940 | 0.126 | 1.024 |
| 1 | – | 0.564 | 1.012 | 0.560 | 0.995 | 0.561 | 0.999 | 0.560 | 0.994 | 0.564 | 1.010 | 0.559 | 0.992 |
| 0.372 | 1.008 | 1 | – | 0 | 0.000 | 0.382 | 1.050 | 0.362 | 0.967 | 0.386 | 1.069 | 0.360 | 0.956 |
| 0.231 | 0.997 | 0 | 0.000 | 1 | – | 0.228 | 0.984 | 0.233 | 1.008 | 0.239 | 1.044 | 0.229 | 0.988 |
| 0.259 | 1.001 | 0.267 | 1.043 | 0.255 | 0.983 | 1 | – | 0 | 0.000 | 0.277 | 1.097 | 0.244 | 0.928 |
| 0.431 | 0.995 | 0.422 | 0.962 | 0.434 | 1.011 | 0 | 0.000 | 1 | – | 0.413 | 0.928 | 0.445 | 1.058 |
| 0.239 | 1.006 | 0.248 | 1.055 | 0.246 | 1.042 | 0.254 | 1.088 | 0.228 | 0.945 | 1 | – | 0 | 0.000 |
| 0.437 | 0.994 | 0.426 | 0.953 | 0.434 | 0.983 | 0.414 | 0.905 | 0.453 | 1.060 | 0 | 0.000 | 1 | – |
| 0.445 | 1.007 | 0.457 | 1.058 | 0.447 | 1.015 | 0.468 | 1.105 | 0.427 | 0.938 | 0.479 | 1.154 | 0.420 | 0.910 |
| 0.236 | 1.001 | 0.232 | 0.978 | 0.240 | 1.022 | 0.218 | 0.900 | 0.246 | 1.057 | 0.223 | 0.931 | 0.246 | 1.058 |
| 0.285 | 1.001 | 0.291 | 1.029 | 0.283 | 0.991 | 0.301 | 1.083 | 0.276 | 0.955 | 0.300 | 1.074 | 0.275 | 0.954 |
| 0.179 | 0.996 | 0.180 | 1.000 | 0.184 | 1.031 | 0.174 | 0.961 | 0.185 | 1.039 | 0.178 | 0.985 | 0.184 | 1.030 |
| 0.256 | 0.999 | 0.265 | 1.047 | 0.257 | 1.004 | 0.276 | 1.105 | 0.245 | 0.942 | 0.276 | 1.106 | 0.240 | 0.915 |
| 0.393 | 1.001 | 0.383 | 0.959 | 0.391 | 0.994 | 0.373 | 0.920 | 0.407 | 1.060 | 0.370 | 0.908 | 0.408 | 1.066 |
| 0.241 | 1.001 | 0.247 | 1.033 | 0.244 | 1.014 | 0.252 | 1.060 | 0.236 | 0.974 | 0.256 | 1.081 | 0.232 | 0.950 |
| 0.394 | 0.996 | 0.389 | 0.976 | 0.399 | 1.019 | 0.382 | 0.948 | 0.403 | 1.036 | 0.388 | 0.974 | 0.399 | 1.019 |
| 0.398 | 1.003 | 0.409 | 1.046 | 0.402 | 1.016 | 0.422 | 1.107 | 0.383 | 0.941 | 0.427 | 1.130 | 0.377 | 0.916 |
| 0.252 | 0.999 | 0.245 | 0.964 | 0.252 | 1.002 | 0.237 | 0.921 | 0.262 | 1.053 | 0.240 | 0.937 | 0.263 | 1.061 |
| 0.224 | 1.005 | 0.232 | 1.051 | 0.223 | 0.998 | 0.240 | 1.099 | 0.214 | 0.948 | 0.242 | 1.108 | 0.210 | 0.923 |
| 0.422 | 0.998 | 0.416 | 0.970 | 0.425 | 1.009 | 0.406 | 0.933 | 0.432 | 1.039 | 0.410 | 0.948 | 0.433 | 1.043 |
| 0.218 | 1.005 | 0.226 | 1.051 | 0.214 | 0.978 | 0.242 | 1.146 | 0.204 | 0.919 | 0.235 | 1.106 | 0.201 | 0.905 |
| 0.311 | 0.996 | 0.312 | 0.997 | 0.315 | 1.013 | 0.307 | 0.978 | 0.316 | 1.018 | 0.308 | 0.980 | 0.314 | 1.010 |
| 0.181 | 1.000 | 0.186 | 1.032 | 0.180 | 0.995 | 0.193 | 1.081 | 0.175 | 0.956 | 0.191 | 1.064 | 0.174 | 0.954 |
| 0.519 | 0.996 | 0.515 | 0.982 | 0.523 | 1.015 | 0.507 | 0.950 | 0.530 | 1.040 | 0.508 | 0.954 | 0.529 | 1.038 |
| 0.097 | 1.002 | 0.095 | 0.980 | 0.096 | 0.986 | 0.097 | 1.008 | 0.097 | 0.998 | 0.097 | 1.004 | 0.097 | 1.002 |
| 0.628 | 0.997 | 0.624 | 0.978 | 0.633 | 1.019 | 0.615 | 0.942 | 0.636 | 1.031 | 0.621 | 0.967 | 0.635 | 1.025 |
| 0.419 | 0.945 | 0.392 | 0.843 | 0.436 | 1.014 | 0.384 | 0.816 | 0.476 | 1.188 | 0.393 | 0.849 | 0.477 | 1.193 |
| 0.581 | 1.059 | 0.608 | 1.186 | 0.564 | 0.986 | 0.616 | 1.226 | 0.524 | 0.841 | 0.607 | 1.178 | 0.523 | 0.838 |
| 0.208 | 1.016 | 0.248 | 1.278 | 0.236 | 1.194 | 0.278 | 1.492 | 0.163 | 0.752 | 0.343 | 2.025 | 0.117 | 0.515 |
| 0.159 | 0.949 | 0.176 | 1.067 | 0.137 | 0.795 | 0.225 | 1.457 | 0.153 | 0.909 | 0.138 | 0.805 | 0.178 | 1.086 |
| 0.259 | 1.040 | 0.253 | 1.007 | 0.220 | 0.840 | 0.216 | 0.820 | 0.242 | 0.951 | 0.191 | 0.705 | 0.281 | 1.163 |
| 0.165 | 0.967 | 0.119 | 0.661 | 0.158 | 0.918 | 0.156 | 0.902 | 0.186 | 1.118 | 0.133 | 0.751 | 0.175 | 1.038 |
| 0.104 | 1.026 | 0.097 | 0.953 | 0.119 | 1.198 | 0.076 | 0.731 | 0.118 | 1.179 | 0.116 | 1.158 | 0.127 | 1.288 |
| 0.064 | 0.900 | 0.066 | 0.939 | 0.099 | 1.461 | 0.032 | 0.432 | 0.091 | 1.318 | 0.052 | 0.733 | 0.082 | 1.185 |
| 0.041 | 1.176 | 0.041 | 1.168 | 0.030 | 0.858 | 0.017 | 0.469 | 0.047 | 1.363 | 0.025 | 0.713 | 0.039 | 1.110 |

| Friendly | Friendly_odds | Unfriendly | Unfriendly_odds | Healed | Healed_odds | Stimulated | Stimulated_odds | Open | Open_odds | Exclusive | Exclusive_odds | Want to reside | Want to reside_odds |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0.210 | 0.974 | 0.222 | 1.044 | 0.210 | 0.973 | 0.222 | 1.046 | 0.212 | 0.983 | 0.222 | 1.045 | 0.216 | 1.007 |
| 0.163 | 0.922 | 0.178 | 1.026 | 0.171 | 0.977 | 0.175 | 1.010 | 0.160 | 0.904 | 0.183 | 1.062 | 0.165 | 0.940 |
| 0.097 | 0.936 | 0.104 | 1.019 | 0.102 | 0.997 | 0.105 | 1.023 | 0.096 | 0.924 | 0.106 | 1.036 | 0.097 | 0.945 |
| 0.373 | 0.910 | 0.411 | 1.067 | 0.383 | 0.950 | 0.401 | 1.021 | 0.366 | 0.880 | 0.414 | 1.081 | 0.377 | 0.924 |
| 0.088 | 0.988 | 0.096 | 1.086 | 0.083 | 0.936 | 0.093 | 1.057 | 0.088 | 0.993 | 0.092 | 1.037 | 0.090 | 1.017 |
| 0.350 | 1.048 | 0.331 | 0.963 | 0.347 | 1.035 | 0.335 | 0.981 | 0.352 | 1.058 | 0.329 | 0.956 | 0.347 | 1.032 |
| 0.290 | 0.989 | 0.298 | 1.031 | 0.290 | 0.993 | 0.296 | 1.022 | 0.280 | 0.946 | 0.298 | 1.028 | 0.289 | 0.984 |
| 0.267 | 1.067 | 0.243 | 0.937 | 0.264 | 1.050 | 0.250 | 0.973 | 0.271 | 1.089 | 0.244 | 0.946 | 0.264 | 1.051 |
| 0.361 | 0.921 | 0.406 | 1.113 | 0.362 | 0.925 | 0.397 | 1.071 | 0.352 | 0.885 | 0.400 | 1.084 | 0.368 | 0.948 |
| 0.183 | 1.057 | 0.168 | 0.950 | 0.182 | 1.044 | 0.170 | 0.964 | 0.183 | 1.053 | 0.168 | 0.949 | 0.179 | 1.027 |
| 0.480 | 0.964 | 0.506 | 1.067 | 0.478 | 0.953 | 0.500 | 1.044 | 0.476 | 0.949 | 0.499 | 1.040 | 0.484 | 0.980 |
| 0.119 | 0.960 | 0.127 | 1.028 | 0.120 | 0.969 | 0.125 | 1.012 | 0.120 | 0.964 | 0.128 | 1.039 | 0.122 | 0.982 |
| 0.563 | 1.008 | 0.562 | 1.003 | 0.562 | 1.001 | 0.560 | 0.994 | 0.561 | 0.998 | 0.562 | 1.001 | 0.562 | 1.002 |
| 0.382 | 1.052 | 0.364 | 0.973 | 0.378 | 1.034 | 0.370 | 0.998 | 0.383 | 1.057 | 0.360 | 0.960 | 0.379 | 1.038 |
| 0.233 | 1.011 | 0.235 | 1.022 | 0.230 | 0.991 | 0.237 | 1.034 | 0.232 | 1.004 | 0.230 | 0.995 | 0.234 | 1.015 |
| 0.273 | 1.077 | 0.239 | 0.899 | 0.274 | 1.083 | 0.250 | 0.955 | 0.278 | 1.106 | 0.246 | 0.933 | 0.270 | 1.061 |
| 0.416 | 0.938 | 0.451 | 1.080 | 0.418 | 0.946 | 0.445 | 1.057 | 0.413 | 0.925 | 0.447 | 1.065 | 0.423 | 0.963 |
| 0.257 | 1.107 | 0.226 | 0.933 | 0.251 | 1.071 | 0.236 | 0.985 | 0.257 | 1.104 | 0.225 | 0.926 | 0.252 | 1.077 |
| 0.416 | 0.912 | 0.457 | 1.081 | 0.424 | 0.945 | 0.448 | 1.041 | 0.410 | 0.891 | 0.456 | 1.073 | 0.421 | 0.931 |
| 1 | – | 0 | 0.000 | 0.462 | 1.082 | 0.429 | 0.945 | 0.472 | 1.124 | 0.420 | 0.909 | 0.459 | 1.065 |
| 0 | 0.000 | 1 | – | 0.222 | 0.923 | 0.252 | 1.092 | 0.218 | 0.904 | 0.251 | 1.086 | 0.230 | 0.968 |
| 0.297 | 1.062 | 0.268 | 0.918 | 1 | – | 0 | 0.000 | 0.300 | 1.074 | 0.272 | 0.938 | 0.292 | 1.033 |
| 0.174 | 0.962 | 0.192 | 1.085 | 0 | 0.000 | 1 | – | 0.172 | 0.949 | 0.189 | 1.060 | 0.181 | 1.008 |
| 0.273 | 1.091 | 0.237 | 0.902 | 0.270 | 1.071 | 0.245 | 0.943 | 1 | – | 0 | 0.000 | 0.268 | 1.062 |
| 0.372 | 0.916 | 0.418 | 1.110 | 0.375 | 0.927 | 0.412 | 1.083 | 0 | 0.000 | 1 | – | 0.384 | 0.962 |
| 0.250 | 1.049 | 0.235 | 0.968 | 0.247 | 1.032 | 0.242 | 1.007 | 0.252 | 1.061 | 0.235 | 0.969 | 1 | – |
| 0.386 | 0.966 | 0.407 | 1.054 | 0.383 | 0.954 | 0.404 | 1.041 | 0.387 | 0.966 | 0.405 | 1.042 | 0 | 0.000 |
| 0.418 | 1.088 | 0.378 | 0.922 | 0.413 | 1.064 | 0.388 | 0.960 | 0.422 | 1.106 | 0.380 | 0.928 | 0.412 | 1.062 |
| 0.241 | 0.942 | 0.265 | 1.073 | 0.243 | 0.956 | 0.262 | 1.056 | 0.236 | 0.918 | 0.265 | 1.069 | 0.245 | 0.965 |
| 0.238 | 1.087 | 0.208 | 0.915 | 0.234 | 1.064 | 0.214 | 0.950 | 0.243 | 1.116 | 0.210 | 0.925 | 0.233 | 1.059 |
| 0.411 | 0.952 | 0.438 | 1.062 | 0.412 | 0.956 | 0.432 | 1.039 | 0.406 | 0.934 | 0.435 | 1.052 | 0.415 | 0.969 |
| 0.234 | 1.098 | 0.196 | 0.875 | 0.231 | 1.081 | 0.201 | 0.902 | 0.241 | 1.142 | 0.201 | 0.904 | 0.228 | 1.061 |
| 0.306 | 0.972 | 0.323 | 1.049 | 0.306 | 0.970 | 0.325 | 1.061 | 0.307 | 0.978 | 0.320 | 1.037 | 0.313 | 1.006 |
| 0.190 | 1.057 | 0.170 | 0.926 | 0.190 | 1.060 | 0.175 | 0.959 | 0.192 | 1.074 | 0.172 | 0.938 | 0.186 | 1.034 |
| 0.507 | 0.950 | 0.539 | 1.080 | 0.507 | 0.949 | 0.539 | 1.079 | 0.505 | 0.944 | 0.535 | 1.064 | 0.517 | 0.989 |
| 0.097 | 0.998 | 0.095 | 0.984 | 0.097 | 1.006 | 0.094 | 0.973 | 0.095 | 0.984 | 0.098 | 1.010 | 0.096 | 0.991 |
| 0.621 | 0.965 | 0.642 | 1.058 | 0.618 | 0.955 | 0.637 | 1.033 | 0.618 | 0.955 | 0.638 | 1.039 | 0.624 | 0.980 |
| 0.413 | 0.923 | 0.421 | 0.953 | 0.438 | 1.020 | 0.427 | 0.975 | 0.426 | 0.971 | 0.427 | 0.977 | 0.405 | 0.892 |
| 0.587 | 1.083 | 0.579 | 1.049 | 0.562 | 0.981 | 0.573 | 1.026 | 0.574 | 1.029 | 0.573 | 1.024 | 0.595 | 1.121 |
| 0.295 | 1.624 | 0.135 | 0.603 | 0.263 | 1.382 | 0.175 | 0.824 | 0.310 | 1.744 | 0.131 | 0.584 | 0.269 | 1.422 |
| 0.158 | 0.941 | 0.141 | 0.822 | 0.195 | 1.212 | 0.188 | 1.162 | 0.184 | 1.130 | 0.166 | 0.998 | 0.180 | 1.102 |
| 0.245 | 0.968 | 0.234 | 0.911 | 0.247 | 0.978 | 0.170 | 0.608 | 0.241 | 0.947 | 0.221 | 0.844 | 0.203 | 0.758 |
| 0.135 | 0.763 | 0.178 | 1.060 | 0.141 | 0.803 | 0.156 | 0.900 | 0.137 | 0.774 | 0.214 | 1.327 | 0.153 | 0.881 |
| 0.091 | 0.880 | 0.135 | 1.382 | 0.098 | 0.960 | 0.140 | 1.443 | 0.048 | 0.441 | 0.126 | 1.280 | 0.092 | 0.890 |
| 0.049 | 0.682 | 0.116 | 1.732 | 0.036 | 0.492 | 0.115 | 1.721 | 0.058 | 0.811 | 0.083 | 1.202 | 0.068 | 0.963 |
| 0.026 | 0.746 | 0.061 | 1.779 | 0.020 | 0.558 | 0.056 | 1.621 | 0.022 | 0.621 | 0.059 | 1.713 | 0.036 | 1.028 |

| Do not want to reside | Do not want to reside_odds | Warm | Warm_odds | Aloof | Aloof_odds | Fascinating | Fascinating_odds | Not fascinating | Not fascinating_odds | Want to play | Want to play_odds | Want to examine deliberately | Want to examine deliberately_odds |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0.216 | 1.009 | 0.212 | 0.982 | 0.216 | 1.006 | 0.214 | 0.996 | 0.214 | 0.999 | 0.213 | 0.988 | 0.220 | 1.034 |
| 0.174 | 1.002 | 0.164 | 0.934 | 0.183 | 1.062 | 0.160 | 0.904 | 0.179 | 1.035 | 0.162 | 0.915 | 0.174 | 1.002 |
| 0.101 | 0.986 | 0.096 | 0.924 | 0.109 | 1.072 | 0.096 | 0.925 | 0.105 | 1.025 | 0.093 | 0.903 | 0.102 | 0.995 |
| 0.400 | 1.018 | 0.373 | 0.908 | 0.417 | 1.092 | 0.367 | 0.885 | 0.408 | 1.051 | 0.362 | 0.866 | 0.396 | 1.002 |
| 0.091 | 1.030 | 0.089 | 1.000 | 0.088 | 0.989 | 0.090 | 1.011 | 0.089 | 1.003 | 0.089 | 1.002 | 0.092 | 1.046 |
| 0.334 | 0.974 | 0.349 | 1.045 | 0.328 | 0.952 | 0.353 | 1.061 | 0.333 | 0.970 | 0.356 | 1.075 | 0.337 | 0.990 |
| 0.292 | 1.001 | 0.288 | 0.980 | 0.300 | 1.042 | 0.284 | 0.961 | 0.296 | 1.021 | 0.278 | 0.935 | 0.290 | 0.991 |
| 0.252 | 0.985 | 0.266 | 1.063 | 0.245 | 0.950 | 0.270 | 1.082 | 0.248 | 0.967 | 0.271 | 1.090 | 0.253 | 0.989 |
| 0.390 | 1.042 | 0.364 | 0.931 | 0.398 | 1.075 | 0.353 | 0.889 | 0.394 | 1.057 | 0.345 | 0.859 | 0.387 | 1.029 |
| 0.172 | 0.976 | 0.182 | 1.046 | 0.170 | 0.961 | 0.183 | 1.054 | 0.172 | 0.976 | 0.184 | 1.062 | 0.171 | 0.974 |
| 0.497 | 1.030 | 0.482 | 0.969 | 0.497 | 1.032 | 0.476 | 0.948 | 0.497 | 1.028 | 0.470 | 0.927 | 0.495 | 1.021 |
| 0.125 | 1.012 | 0.122 | 0.981 | 0.125 | 1.009 | 0.120 | 0.966 | 0.125 | 1.012 | 0.123 | 0.995 | 0.126 | 1.020 |
| 0.560 | 0.994 | 0.562 | 1.002 | 0.561 | 0.997 | 0.564 | 1.009 | 0.561 | 0.997 | 0.563 | 1.006 | 0.560 | 0.994 |
| 0.364 | 0.975 | 0.380 | 1.045 | 0.361 | 0.961 | 0.384 | 1.063 | 0.363 | 0.971 | 0.385 | 1.065 | 0.369 | 0.994 |
| 0.234 | 1.015 | 0.234 | 1.013 | 0.231 | 1.001 | 0.231 | 0.999 | 0.232 | 1.007 | 0.227 | 0.976 | 0.233 | 1.011 |
| 0.250 | 0.956 | 0.275 | 1.085 | 0.244 | 0.925 | 0.278 | 1.102 | 0.249 | 0.949 | 0.287 | 1.154 | 0.255 | 0.983 |
| 0.441 | 1.037 | 0.415 | 0.935 | 0.448 | 1.070 | 0.413 | 0.926 | 0.441 | 1.040 | 0.402 | 0.884 | 0.437 | 1.024 |
| 0.234 | 0.975 | 0.256 | 1.097 | 0.226 | 0.935 | 0.257 | 1.104 | 0.231 | 0.960 | 0.256 | 1.097 | 0.235 | 0.981 |
| 0.443 | 1.018 | 0.415 | 0.910 | 0.457 | 1.080 | 0.411 | 0.893 | 0.449 | 1.046 | 0.403 | 0.866 | 0.441 | 1.013 |
| 0.433 | 0.962 | 0.465 | 1.095 | 0.423 | 0.923 | 0.471 | 1.122 | 0.430 | 0.950 | 0.476 | 1.141 | 0.434 | 0.966 |
| 0.244 | 1.042 | 0.225 | 0.937 | 0.249 | 1.070 | 0.221 | 0.915 | 0.244 | 1.045 | 0.213 | 0.873 | 0.244 | 1.043 |
| 0.277 | 0.961 | 0.295 | 1.051 | 0.275 | 0.950 | 0.299 | 1.069 | 0.278 | 0.964 | 0.301 | 1.083 | 0.279 | 0.970 |
| 0.184 | 1.029 | 0.176 | 0.972 | 0.188 | 1.054 | 0.173 | 0.956 | 0.184 | 1.027 | 0.167 | 0.913 | 0.187 | 1.051 |
| 0.251 | 0.972 | 0.272 | 1.083 | 0.240 | 0.918 | 0.279 | 1.121 | 0.246 | 0.948 | 0.284 | 1.152 | 0.253 | 0.980 |
| 0.403 | 1.042 | 0.375 | 0.927 | 0.412 | 1.086 | 0.370 | 0.907 | 0.404 | 1.049 | 0.362 | 0.879 | 0.403 | 1.042 |
| 0 | 0.000 | 0.250 | 1.051 | 0.235 | 0.968 | 0.252 | 1.062 | 0.237 | 0.977 | 0.253 | 1.068 | 0.242 | 1.004 |
| 1 | – | 0.388 | 0.974 | 0.403 | 1.036 | 0.386 | 0.964 | 0.401 | 1.026 | 0.381 | 0.945 | 0.401 | 1.026 |
| 0.391 | 0.972 | 1 | – | 0 | 0.000 | 0.421 | 1.099 | 0.387 | 0.957 | 0.430 | 1.141 | 0.394 | 0.984 |
| 0.257 | 1.027 | 0 | 0.000 | 1 | – | 0.238 | 0.927 | 0.259 | 1.039 | 0.229 | 0.881 | 0.255 | 1.017 |
| 0.218 | 0.971 | 0.237 | 1.078 | 0.212 | 0.935 | 1 | – | 0 | 0.000 | 0.248 | 1.146 | 0.219 | 0.978 |
| 0.430 | 1.027 | 0.411 | 0.953 | 0.435 | 1.048 | 0 | 0.000 | 1 | – | 0.400 | 0.909 | 0.427 | 1.016 |
| 0.210 | 0.956 | 0.235 | 1.105 | 0.199 | 0.893 | 0.240 | 1.135 | 0.207 | 0.937 | 1 | – | 0 | 0.000 |
| 0.317 | 1.023 | 0.309 | 0.984 | 0.315 | 1.015 | 0.307 | 0.976 | 0.315 | 1.012 | 0 | 0.000 | 1 | – |
| 0.176 | 0.966 | 0.189 | 1.053 | 0.173 | 0.947 | 0.191 | 1.068 | 0.176 | 0.965 | 0.195 | 1.093 | 0.178 | 0.978 |
| 0.528 | 1.035 | 0.510 | 0.963 | 0.531 | 1.047 | 0.506 | 0.947 | 0.527 | 1.030 | 0.500 | 0.923 | 0.530 | 1.043 |
| 0.096 | 0.996 | 0.097 | 0.999 | 0.097 | 1.007 | 0.096 | 0.990 | 0.097 | 1.004 | 0.097 | 1.002 | 0.095 | 0.984 |
| 0.636 | 1.032 | 0.622 | 0.972 | 0.636 | 1.031 | 0.618 | 0.954 | 0.635 | 1.027 | 0.615 | 0.941 | 0.634 | 1.020 |
| 0.449 | 1.068 | 0.383 | 0.814 | 0.490 | 1.257 | 0.408 | 0.902 | 0.451 | 1.076 | 0.359 | 0.734 | 0.417 | 0.936 |
| 0.551 | 0.937 | 0.617 | 1.228 | 0.510 | 0.796 | 0.592 | 1.109 | 0.549 | 0.929 | 0.641 | 1.363 | 0.583 | 1.069 |
| 0.191 | 0.912 | 0.283 | 1.529 | 0.149 | 0.680 | 0.305 | 1.699 | 0.169 | 0.785 | 0.304 | 1.691 | 0.186 | 0.885 |
| 0.145 | 0.848 | 0.168 | 1.012 | 0.156 | 0.929 | 0.178 | 1.083 | 0.150 | 0.882 | 0.189 | 1.167 | 0.183 | 1.121 |
| 0.228 | 0.878 | 0.229 | 0.886 | 0.234 | 0.908 | 0.243 | 0.955 | 0.247 | 0.976 | 0.268 | 1.088 | 0.210 | 0.791 |
| 0.199 | 1.215 | 0.155 | 0.896 | 0.189 | 1.134 | 0.134 | 0.758 | 0.190 | 1.143 | 0.162 | 0.942 | 0.179 | 1.064 |
| 0.098 | 0.959 | 0.088 | 0.849 | 0.137 | 1.405 | 0.060 | 0.562 | 0.121 | 1.215 | 0.033 | 0.305 | 0.101 | 0.987 |
| 0.096 | 1.402 | 0.058 | 0.821 | 0.078 | 1.121 | 0.046 | 0.635 | 0.084 | 1.212 | 0.031 | 0.417 | 0.098 | 1.441 |
| 0.044 | 1.266 | 0.019 | 0.518 | 0.057 | 1.661 | 0.035 | 0.985 | 0.040 | 1.151 | 0.014 | 0.386 | 0.044 | 1.255 |

| Lively | Lively_odds | Calm | Calm_odds | Atmosphere of urban | Atmosphere of urban_odds | Atmosphere of rural area | Atmosphere of rural area_odds | Gender Male | Male_odds | Female | Female_odds | Age 10th | 10th_odds |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0.211 | 0.978 | 0.219 | 1.025 | 0.211 | 0.978 | 0.215 | 1.001 | 0.181 | 0.807 | 0.240 | 1.159 | 0.180 | 0.802 |
| 0.169 | 0.965 | 0.176 | 1.017 | 0.179 | 1.038 | 0.175 | 1.005 | 0.195 | 1.150 | 0.158 | 0.890 | 0.070 | 0.357 |
| 0.101 | 0.982 | 0.103 | 1.010 | 0.103 | 1.003 | 0.102 | 0.994 | 0.132 | 1.328 | 0.080 | 0.764 | 0.044 | 0.402 |
| 0.379 | 0.933 | 0.401 | 1.024 | 0.404 | 1.034 | 0.399 | 1.016 | 0.449 | 1.247 | 0.355 | 0.839 | 0.214 | 0.416 |
| 0.085 | 0.957 | 0.091 | 1.033 | 0.084 | 0.945 | 0.090 | 1.018 | 0.058 | 0.637 | 0.112 | 1.294 | 0.085 | 0.954 |
| 0.349 | 1.044 | 0.335 | 0.982 | 0.335 | 0.980 | 0.336 | 0.986 | 0.306 | 0.857 | 0.365 | 1.120 | 0.403 | 1.314 |
| 0.288 | 0.979 | 0.293 | 1.004 | 0.298 | 1.030 | 0.293 | 1.006 | 0.300 | 1.038 | 0.286 | 0.972 | 0.278 | 0.932 |
| 0.265 | 1.053 | 0.250 | 0.977 | 0.254 | 0.998 | 0.252 | 0.985 | 0.250 | 0.976 | 0.258 | 1.019 | 0.363 | 1.664 |
| 0.361 | 0.920 | 0.390 | 1.043 | 0.382 | 1.006 | 0.388 | 1.032 | 0.388 | 1.033 | 0.375 | 0.975 | 0.253 | 0.552 |
| 0.182 | 1.044 | 0.171 | 0.972 | 0.176 | 1.006 | 0.173 | 0.987 | 0.169 | 0.954 | 0.180 | 1.035 | 0.233 | 1.431 |
| 0.478 | 0.954 | 0.496 | 1.026 | 0.487 | 0.991 | 0.495 | 1.020 | 0.488 | 0.992 | 0.491 | 1.006 | 0.441 | 0.821 |
| 0.121 | 0.975 | 0.126 | 1.018 | 0.124 | 1.003 | 0.125 | 1.008 | 0.109 | 0.864 | 0.135 | 1.107 | 0.085 | 0.657 |
| 0.561 | 1.000 | 0.560 | 0.995 | 0.563 | 1.005 | 0.561 | 0.998 | 0.543 | 0.930 | 0.575 | 1.058 | 0.569 | 1.030 |
| 0.380 | 1.042 | 0.367 | 0.985 | 0.364 | 0.973 | 0.367 | 0.986 | 0.335 | 0.857 | 0.397 | 1.121 | 0.447 | 1.379 |
| 0.230 | 0.994 | 0.233 | 1.009 | 0.228 | 0.984 | 0.233 | 1.009 | 0.233 | 1.011 | 0.230 | 0.992 | 0.266 | 1.202 |
| 0.275 | 1.090 | 0.252 | 0.966 | 0.260 | 1.006 | 0.253 | 0.971 | 0.229 | 0.853 | 0.281 | 1.121 | 0.350 | 1.547 |
| 0.416 | 0.937 | 0.440 | 1.034 | 0.431 | 0.998 | 0.437 | 1.020 | 0.474 | 1.188 | 0.399 | 0.875 | 0.342 | 0.685 |
| 0.251 | 1.072 | 0.233 | 0.969 | 0.239 | 1.001 | 0.235 | 0.984 | 0.217 | 0.883 | 0.255 | 1.094 | 0.399 | 2.120 |
| 0.422 | 0.936 | 0.446 | 1.030 | 0.439 | 1.002 | 0.442 | 1.016 | 0.483 | 1.196 | 0.404 | 0.870 | 0.251 | 0.429 |
| 0.464 | 1.087 | 0.432 | 0.956 | 0.442 | 0.997 | 0.437 | 0.976 | 0.423 | 0.922 | 0.458 | 1.064 | 0.637 | 2.212 |
| 0.222 | 0.921 | 0.245 | 1.049 | 0.233 | 0.982 | 0.241 | 1.027 | 0.230 | 0.965 | 0.241 | 1.027 | 0.155 | 0.593 |
| 0.299 | 1.069 | 0.278 | 0.965 | 0.287 | 1.008 | 0.280 | 0.976 | 0.288 | 1.016 | 0.283 | 0.988 | 0.365 | 1.443 |
| 0.174 | 0.960 | 0.186 | 1.045 | 0.175 | 0.970 | 0.182 | 1.015 | 0.177 | 0.982 | 0.182 | 1.013 | 0.154 | 0.828 |
| 0.272 | 1.082 | 0.249 | 0.963 | 0.253 | 0.980 | 0.252 | 0.977 | 0.252 | 0.978 | 0.260 | 1.017 | 0.388 | 1.838 |
| 0.373 | 0.919 | 0.404 | 1.050 | 0.396 | 1.015 | 0.398 | 1.023 | 0.387 | 0.978 | 0.397 | 1.017 | 0.251 | 0.517 |
| 0.248 | 1.038 | 0.240 | 0.992 | 0.239 | 0.988 | 0.239 | 0.990 | 0.226 | 0.917 | 0.253 | 1.066 | 0.315 | 1.451 |
| 0.383 | 0.954 | 0.401 | 1.028 | 0.393 | 0.995 | 0.399 | 1.019 | 0.409 | 1.063 | 0.383 | 0.954 | 0.367 | 0.888 |
| 0.415 | 1.075 | 0.390 | 0.969 | 0.397 | 0.997 | 0.394 | 0.983 | 0.352 | 0.823 | 0.433 | 1.154 | 0.549 | 1.840 |
| 0.241 | 0.941 | 0.257 | 1.029 | 0.254 | 1.011 | 0.255 | 1.015 | 0.285 | 1.183 | 0.227 | 0.871 | 0.183 | 0.666 |
| 0.236 | 1.073 | 0.217 | 0.965 | 0.221 | 0.989 | 0.219 | 0.977 | 0.210 | 0.926 | 0.233 | 1.058 | 0.332 | 1.727 |
| 0.411 | 0.950 | 0.429 | 1.025 | 0.425 | 1.008 | 0.427 | 1.017 | 0.441 | 1.075 | 0.410 | 0.946 | 0.348 | 0.727 |
| 0.234 | 1.097 | 0.209 | 0.950 | 0.218 | 1.004 | 0.213 | 0.971 | 0.180 | 0.792 | 0.246 | 1.173 | 0.322 | 1.710 |
| 0.307 | 0.973 | 0.319 | 1.030 | 0.308 | 0.979 | 0.315 | 1.011 | 0.301 | 0.946 | 0.321 | 1.042 | 0.283 | 0.869 |
| 1 | – | 0 | 0.000 | 0.181 | 0.997 | 0.178 | 0.979 | 0.175 | 0.961 | 0.186 | 1.030 | 0.234 | 1.385 |
| 0 | 0.000 | 1 | – | 0.515 | 0.983 | 0.524 | 1.018 | 0.516 | 0.985 | 0.522 | 1.011 | 0.444 | 0.737 |
| 0.096 | 0.997 | 0.096 | 0.991 | 1 | – | 0 | 0.000 | 0.098 | 1.015 | 0.096 | 0.988 | 0.094 | 0.973 |
| 0.618 | 0.956 | 0.634 | 1.024 | 0 | 0.000 | 1 | – | 0.629 | 1.001 | 0.629 | 0.999 | 0.585 | 0.831 |
| 0.419 | 0.944 | 0.430 | 0.987 | 0.439 | 1.025 | 0.433 | 1.001 | 1 | – | 0 | 0.000 | 0.433 | 1.000 |
| 0.581 | 1.059 | 0.570 | 1.013 | 0.561 | 0.976 | 0.567 | 0.999 | 0 | 0.000 | 1 | – | 0.567 | 1.000 |
| 0.266 | 1.401 | 0.175 | 0.822 | 0.200 | 0.969 | 0.191 | 0.914 | 0.205 | 1.000 | 0.205 | 1.000 | 1 | – |
| 0.197 | 1.232 | 0.171 | 1.037 | 0.147 | 0.865 | 0.150 | 0.882 | 0.166 | 1.000 | 0.166 | 1.000 | 0 | 0.000 |
| 0.252 | 1.003 | 0.224 | 0.859 | 0.251 | 1.000 | 0.248 | 0.981 | 0.251 | 1.000 | 0.251 | 1.000 | 0 | 0.000 |
| 0.135 | 0.763 | 0.182 | 1.085 | 0.219 | 1.365 | 0.184 | 1.097 | 0.170 | 1.000 | 0.170 | 1.000 | 0 | 0.000 |
| 0.088 | 0.855 | 0.110 | 1.087 | 0.117 | 1.173 | 0.106 | 1.050 | 0.102 | 1.000 | 0.102 | 1.000 | 0 | 0.000 |
| 0.045 | 0.626 | 0.094 | 1.371 | 0.034 | 0.471 | 0.085 | 1.234 | 0.070 | 1.000 | 0.070 | 1.000 | 0 | 0.000 |
| 0.017 | 0.465 | 0.044 | 1.275 | 0.031 | 0.885 | 0.037 | 1.045 | 0.035 | 1.000 | 0.035 | 1.000 | 0 | 0.000 |

| 20th | 20st_odds | 30th | 30st_odds | 40th | 40st_odds | 50th | 50st_odds | 60th | 60st_odds | More than70 | More than70_odds |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0.262 | 1.301 | 0.196 | 0.891 | 0.212 | 0.983 | 0.170 | 0.751 | 0.240 | 1.156 | 0.419 | 2.640 |
| 0.229 | 1.411 | 0.182 | 1.057 | 0.230 | 1.417 | 0.240 | 1.502 | 0.127 | 0.688 | 0.095 | 0.496 |
| 0.158 | 1.642 | 0.117 | 1.155 | 0.084 | 0.799 | 0.148 | 1.516 | 0.077 | 0.728 | 0.095 | 0.915 |
| 0.402 | 1.025 | 0.437 | 1.183 | 0.473 | 1.372 | 0.531 | 1.729 | 0.371 | 0.900 | 0.419 | 1.101 |
| 0.066 | 0.726 | 0.092 | 1.041 | 0.075 | 0.827 | 0.050 | 0.545 | 0.168 | 2.068 | 0.216 | 2.833 |
| 0.362 | 1.102 | 0.354 | 1.069 | 0.259 | 0.681 | 0.297 | 0.824 | 0.305 | 0.856 | 0.333 | 0.973 |
| 0.238 | 0.760 | 0.279 | 0.937 | 0.298 | 1.032 | 0.460 | 2.065 | 0.233 | 0.736 | 0.323 | 1.158 |
| 0.258 | 1.016 | 0.219 | 0.818 | 0.238 | 0.914 | 0.197 | 0.719 | 0.204 | 0.748 | 0.220 | 0.825 |
| 0.305 | 0.714 | 0.374 | 0.972 | 0.443 | 1.294 | 0.546 | 1.954 | 0.556 | 2.041 | 0.399 | 1.079 |
| 0.156 | 0.872 | 0.180 | 1.033 | 0.144 | 0.792 | 0.186 | 1.074 | 0.120 | 0.642 | 0.124 | 0.664 |
| 0.419 | 0.753 | 0.482 | 0.971 | 0.516 | 1.110 | 0.552 | 1.282 | 0.651 | 1.944 | 0.533 | 1.188 |
| 0.141 | 1.158 | 0.125 | 1.009 | 0.154 | 1.293 | 0.109 | 0.861 | 0.142 | 1.170 | 0.124 | 0.999 |
| 0.537 | 0.906 | 0.578 | 1.070 | 0.546 | 0.940 | 0.574 | 1.054 | 0.509 | 0.810 | 0.656 | 1.492 |
| 0.390 | 1.090 | 0.372 | 1.008 | 0.260 | 0.597 | 0.354 | 0.934 | 0.349 | 0.913 | 0.430 | 1.282 |
| 0.190 | 0.782 | 0.202 | 0.843 | 0.215 | 0.912 | 0.272 | 1.240 | 0.327 | 1.617 | 0.199 | 0.828 |
| 0.350 | 1.545 | 0.222 | 0.818 | 0.237 | 0.892 | 0.194 | 0.691 | 0.116 | 0.378 | 0.124 | 0.405 |
| 0.398 | 0.871 | 0.416 | 0.936 | 0.473 | 1.182 | 0.500 | 1.317 | 0.556 | 1.651 | 0.581 | 1.824 |
| 0.198 | 0.790 | 0.181 | 0.708 | 0.187 | 0.734 | 0.272 | 1.191 | 0.178 | 0.693 | 0.172 | 0.663 |
| 0.469 | 1.132 | 0.490 | 1.230 | 0.452 | 1.058 | 0.549 | 1.558 | 0.513 | 1.349 | 0.485 | 1.205 |
| 0.421 | 0.915 | 0.432 | 0.957 | 0.352 | 0.684 | 0.394 | 0.820 | 0.309 | 0.563 | 0.333 | 0.629 |
| 0.200 | 0.809 | 0.220 | 0.913 | 0.248 | 1.065 | 0.314 | 1.482 | 0.389 | 2.060 | 0.409 | 2.238 |
| 0.334 | 1.257 | 0.280 | 0.977 | 0.237 | 0.778 | 0.275 | 0.950 | 0.145 | 0.427 | 0.162 | 0.483 |
| 0.203 | 1.165 | 0.121 | 0.630 | 0.165 | 0.899 | 0.248 | 1.506 | 0.295 | 1.905 | 0.285 | 1.821 |
| 0.284 | 1.149 | 0.246 | 0.946 | 0.207 | 0.755 | 0.120 | 0.395 | 0.211 | 0.775 | 0.162 | 0.558 |
| 0.392 | 0.997 | 0.345 | 0.814 | 0.494 | 1.508 | 0.489 | 1.477 | 0.465 | 1.346 | 0.656 | 2.954 |
| 0.261 | 1.113 | 0.195 | 0.760 | 0.217 | 0.871 | 0.217 | 0.873 | 0.233 | 0.955 | 0.247 | 1.035 |
| 0.343 | 0.802 | 0.357 | 0.853 | 0.463 | 1.320 | 0.380 | 0.940 | 0.538 | 1.787 | 0.495 | 1.503 |
| 0.402 | 1.017 | 0.363 | 0.862 | 0.363 | 0.862 | 0.343 | 0.790 | 0.331 | 0.749 | 0.210 | 0.401 |
| 0.237 | 0.921 | 0.234 | 0.908 | 0.279 | 1.151 | 0.340 | 1.530 | 0.280 | 1.155 | 0.409 | 2.055 |
| 0.238 | 1.090 | 0.216 | 0.956 | 0.176 | 0.745 | 0.131 | 0.526 | 0.145 | 0.592 | 0.220 | 0.981 |
| 0.381 | 0.838 | 0.416 | 0.969 | 0.472 | 1.220 | 0.503 | 1.380 | 0.505 | 1.394 | 0.485 | 1.282 |
| 0.247 | 1.180 | 0.232 | 1.084 | 0.207 | 0.939 | 0.071 | 0.276 | 0.095 | 0.376 | 0.086 | 0.338 |
| 0.343 | 1.151 | 0.261 | 0.776 | 0.329 | 1.078 | 0.309 | 0.984 | 0.436 | 1.705 | 0.388 | 1.398 |
| 0.215 | 1.237 | 0.181 | 1.002 | 0.144 | 0.761 | 0.157 | 0.843 | 0.116 | 0.595 | 0.086 | 0.425 |
| 0.535 | 1.065 | 0.463 | 0.796 | 0.556 | 1.157 | 0.560 | 1.176 | 0.695 | 2.102 | 0.656 | 1.766 |
| 0.086 | 0.874 | 0.097 | 1.000 | 0.124 | 1.326 | 0.112 | 1.172 | 0.047 | 0.463 | 0.086 | 0.878 |
| 0.566 | 0.769 | 0.620 | 0.963 | 0.679 | 1.247 | 0.657 | 1.130 | 0.764 | 1.906 | 0.656 | 1.127 |
| 0.433 | 1.000 | 0.433 | 1.000 | 0.433 | 1.000 | 0.433 | 1.000 | 0.433 | 1.000 | 0.433 | 1.000 |
| 0.567 | 1.000 | 0.567 | 1.000 | 0.567 | 1.000 | 0.567 | 1.000 | 0.567 | 1.000 | 0.567 | 1.000 |
| 0 | 0.000 | 0 | 0.000 | 0 | 0.000 | 0 | 0.000 | 0 | 0.000 | 0 | 0.000 |
| 1 | – | 0 | 0.000 | 0 | 0.000 | 0 | 0.000 | 0 | 0.000 | 0 | 0.000 |
| 0 | 0.000 | 1 | – | 0 | 0.000 | 0 | 0.000 | 0 | 0.000 | 0 | 0.000 |
| 0 | 0.000 | 0 | 0.000 | 1 | – | 0 | 0.000 | 0 | 0.000 | 0 | 0.000 |
| 0 | 0.000 | 0 | 0.000 | 0 | 0.000 | 1 | – | 0 | 0.000 | 0 | 0.000 |
| 0 | 0.000 | 0 | 0.000 | 0 | 0.000 | 0 | 0.000 | 1 | – | 0 | 0.000 |
| 0 | 0.000 | 0 | 0.000 | 0 | 0.000 | 0 | 0.000 | 0 | 0.000 | 1 | – |

# Alerts Clustering for Intrusion Detection Systems: Overview and Machine Learning Perspectives

Wajdi Alhakami

Department of Computer Sciences, College of Computers and Information Technology,
Taif University, Taif, Saudi Arabia, KSA

*Abstract*—The tremendous amount of the security alerts due to the high-speed alert generation of high-speed networks make the management of intrusion detection computationally expensive. Evidently, the high-level rate of wrong alerts disproves the Intrusion Detection Systems (IDS) performances and decrease its capability to prevent cyber-attacks which lead to tedious alert analysis task. Thus, it is important to develop new tools to understand intrusion data and to represent them in a compact forms using, for example, an alert clustering process. This hot topic of research is studied here and an understandable taxonomy followed by a deep survey of main published works related to intrusion alert management is presented in this paper. The second part of this work exposes different useful steps for designing a unified IDS system on the basis of machine learning techniques which are considered one of the most powerful tools for solving certain problems related to alert management and outlier detection.

*Keywords*—*Intrusion detection systems; alert clustering; taxonomy; survey; machine learning*

## I. INTRODUCTION

Intrusion Detection Systems (IDSs) are widely deployed into servers for data security purposes. However, these systems produce a lot of false positive alerts making the task of security analysts difficult such as taking suitable actions for them. This problem has received considerable attention from researchers given that manual intrusion alerts management is extremely fastidious and computationally expensive. Consequently, it will not be easy to manage alerts and to take appropriate actions for them. Thereby, the automation of the process of alert management is a necessity. It turns out that this problem becomes more difficult with the context of high-speed networks [1], [2]. Indeed, new severe issues related especially to the scalability, the real time constraints, and efficiency create a major challenge to the success of IDS [3], [4]. In particular, the big quantity of the security alerts leads to an expensive intrusion detection process. Evidently, the high-level rate of wrong alerts disproves the IDS performances and decrease its capability to prevent cyber-attacks that lead to tedious alert analysis task. In particular, alert clustering and outlier detection are crucial problems for taking suitable actions and for security threats understanding [5], [1], [2]. The content of this paper is a taxonomy and a survey related to the intrusion alert management. In fact, an understandable taxonomy, especially for beginner researchers, is presented. It is related to intrusion detection system with special emphasize on intrusion alert management problem. An informative description is presented for intrusion detection systems, for both misuse and anomaly-based detection, and for low- and high-level alert management (i.e. alert ranking, normalization, clustering, correlation, etc.)

[6], [7], [8]. The second part of this paper is dedicated to expose how to design a possible IDS framework using machine learning techniques which are considered one of the most powerful tools for solving certain problems related to alert management and outlier detection. The rest of this paper is organized as follows. In the next section, a taxonomy related to the current area of research is presented. Section 3 is devoted to survey main existing works in the literature. Section 4 describes how can machine learning techniques be involved and considered as an interesting alternative to deal with such problem. Finally, we conclude the paper in the last section.

## II. TAXONOMY

This section presents a comprehensible taxonomy, especially for beginner researchers, for intrusion detection systems (IDS) and intrusion alert management problems. In particular, informative descriptions are given for intrusion detection techniques, alert management, alert clustering/classification and alert correlation.

*1) Intrusion Detection Systems (IDSs):* Intrusion detection systems are widely deployed into both hosts and networks to protect assets. To ensure security, most of developed IDSs apply mainly the so-called misuse-based (named also signature-based) or anomaly-based IDS techniques [9]. The key purpose of these techniques is to help the security administrator to fully recognize what the IDS is doing. Fig. 1 illustrate main intrusion detection techniques.

- **Misuse (Signature)-based detection**

Misuse-based intrusion detection techniques have been proven to detect effectively attacks without generating a great number of false alarms. For this reason, they are broadly approved in the most commercial systems. Such technique can be used to detect known attacks, so, we have to create signatures (patterns) for known attacks and store them into databases as a priori information. However, this kind of approach cannot detect unknown attacks and therefore they must be frequently updated with signatures of new attacks.

- **Anomaly-based detection**

Anomaly-based approaches (named behaviour-based) are used to detect unknown (novel or irregular) and known attacks on the basis of their profiles or statistical models. These models employ labeled data to model and train anomaly detection as a classification problem. This kind of approach try to find behavior (attack behavior) deviating from normal one. In general, techniques driven from pattern recognition (parametric

Fig. 1. Classification of intrusion detection approaches.

and non-parametric statistical models, neural networks, rule-based algorithms, hidden Markov model, etc.) are applied to identify normal data and abnormal (anomalous) ones. For example, the user can be notified about the existing of unusual behavior if the network activity deviates from its normal state. If normal behavior is well trained and well modeled, unusual behavior can be labelled and identified as intrusive. Thus, normal behavior should be well defined otherwise, if it is not, a lot of false alerts will be generated. Thus, the modeling problem, which is a difficult task, is very important for the design of anomaly-based approach. Finally, when compared to misuse-based approaches, anomaly-based ones are more efficient and faster, but, they generate a lot of false positive rate.

*2) Intrusion detection alert management:* Alert management function address two different processing: low-level processing and high-level processing (Fig. 2). Each level is able to accomplish specific objectives such as alert pre-processing, alert correlation, alert fusion, etc. Through management function, alerts are processed in order to help the security administrator to understand what the IDSs are addressing.

- **Low-level alert management**

The main purpose of low-level alert management is to facilitate further high-processing. It is required to achieve low processing such as updating alert attributes to a standardized format, scoring and prioritizing alerts, translating all alert's attributes into numerical values, etc. For instance, ranking alerts is required to mark and select only significant alerts for further investigation. In the literature, only few works have been proposed to deal with the importance of low level alert management for the overall system evaluation [3], [4]. Authors showed the importance of this step by defining various metrics to score and prioritize alerts such as applicability, sensor status, severity, alert relationship, etc.

- **Alert Normalization**

Alert attributes vary from one IDS to another. Thus, alerts from diverse sensors are encoded with different formats. For this reason, it is important to unify all information to make easy further processing. The normalization step will convert all diversified alert attributes into an appropriate unified representation such as the well known standard "Intrusion Detection Message Exchange Format" (IDMEF). Such standard is able to provide a flexibility for any possible extension since it is based on XML language. Each alert can be represented by the IDMEF data model as shown in Fig. 3.

- **Alert scoring/ranking/prioritizing**

Ranking and prioritizing alerts is a useful tool to evaluate the relative importance of such alerts. It helps in reduction the amount of incoming alerts by quickly discard irrelevant alerts. Some developed procedures [3], [4], [10] employed mainly the following properties: the integrity, the secrecy, and the availability metrics. In general, these properties are calculated easily since they are stored in the IDS database. As output, each alert is assigned with high or medium or low score. To achieve an appropriate scoring task as proposed in [3], [4], the pseudo-code for alert scoring in Fig. 4 can be used.

- **High-level Alert Management**

Automated and intelligent high-level alert management is a potential task helping the administrator to analyze properly alerts, and to save his time and effort. Alert management can be defined as a generalization function related to a variety of specific operations like alert classification (or clustering), alert fusion (or aggregation), and alert correlation. These operations are considered crucial since they provide an abstraction of a set of alerts. In the literature several approaches were proposed to solve the problem of alert management from different angles and some of them are presented in Section 3.

- **Alert Aggregation/Fusion/Merging**

The role of the aggregation function is to group alerts having same common characteristics such as the source IP, the target IP, the type of the attack, etc. The fusion/merging function attempts to combine a group of alerts into one hyper-alert. The latter should be representative of each belongs to the same cluster (or component). Some relevant methods [6], [7], [8] were developed in this context to help the analyst to take rapidly, through one global meta-alert, appropriate action against the seriousness of the attack.

Fig. 2.   Classification of Intrusion Alert Management Techniques.



Fig. 3.   The IDMEF data model.



Fig. 4.   Pseudo-code for alert scoring.

- **Alert Clustering/Classification**

This function attempts to cluster or classify alerts that match the same attack occurrence and share common features like source/target IP address or port number. Clustering output leads to reduce the number of alerts. Several criteria, as presented in [6], can be defined to achieve this task as a "similarity relation" to connect alerts to specific cluster.

- **Alert correlation**

Alert correlation is used to discover any relationship between alerts in response to launched attacks in order to achieve a final goal. Some criteria [3] are defined to achieve this and provide different ways to study the relationship between the attacks.

## III. SURVEY ON INTRUSION DETECTION ALERT MANAGEMENT

It is noteworthy that several promising approaches have been developed, in the past decades, to solve the problem of intrusion alert management [11], [12]. Intrusion alert management methods can be categorized into four main categories: predefined attack scenarios-based approaches, similarity-based approaches, prerequisites and consequences-based approaches, and hybrid approaches. The main common objective of these approaches is to categorize alerts and to reduce false positive ones.

### A. Predefined Attack Scenarios

This category takes into account only the known attack scenarios in order cluster alerts. These attack scenarios are learned in general from different datasets [13], [7], [14]. In order to construct the detected attack and to correlate alerts, each sequence of alerts must be compared with a known attack. The main advantage of this category is its ability to discover the causal relationship between attacks. However, the main drawback of this category is it is limited to known attacks and not unknown ones which is not very helpful for discovering and detecting new ones. Several solutions are suggested, like in [13], where attack scenarios through chronicle language is proposed. Another work is proposed in [7] that uses explicit rules to solve such problem. In [14] attack scenarios are constructed by comparing probability measures. This probability is defined as a metric and is usually calculated through a training data.

### B. Similarity-based Approaches

The second category of approaches involves the definition of a similarity metric between alert's attributes (e.g. source/target IP address, port number, etc.) to classify alerts [15], [16], [17], [18], [6]. Such metric may discover the relationships between different alerts. The obtained score, after calculating the similarity metric, decides if these alerts will be correlated or no. Although this category can be considered as effective for many cases; however, it cannot find the main causal relationship between alerts. For example, in [16], authors grouped alerts into common cluster using an clustering algorithm. A probabilistic-based distance method for alert clustering was also implemented in [15]. The suggested probabilistic model consists of a unified mathematical framework

with appropriate similarity functions for each alert feature. As a result, alerts are grouped if the similarity measures are closely matched.

### C. Prerequisites and Consequence-based Approaches

A third category has the role to match prerequisites with the consequences based on the dependencies between alerts [19], [20], [21], [22]. Two or more attacks are correlated if any of the prerequisites of the later attack match any of the consequences of the early one. With this principle, causal relationships between attacks can be successfully identified, and it is will possible to build a new attack scenarios by connecting each attack into a sequence of causal relations. The main advantage of this kind of approach is its simplicity to find out the casual relationship between alerts, but the process of discovering individual attacks is computationally expensive. A typical work is proposed in [21] which is based on logical predicates to construct prerequisites and consequences model of attacks.

### D. Hybrid Approaches

Hybrid methods are proposed to overcome limitation of applying only single algorithm and to solve the alert management problems with several techniques at the same time [23], [24], [25], [26], [11]. Some interesting papers [12], [11] demonstrate that hybrid approaches provide better flexibility. For instance, a hybrid fuzzy-based anomaly IDS utilizing hidden Markov model (HMM) detector and a normal database detector to minimize false alert rate was developed in [24]. Another decision support system (DSS) for online network behavior monitoring is proposed in [23]. The developed classification model involves three phases: alert preprocessing, model constructing and rule refining. In [25], an effective algorithm is implemented for filtering false alert in network-based IDSs. The proposed filter involves three main components which are based on statistical properties of the alerts. In [8], [11], a collaborative architecture for multiple IDSs to detect real-time network intrusions is also developed. More advanced works are proposed in the literature such as the one published in [27] that enables alert aggregation and minimizing false alerts using an anomaly detector technique. In [26], authors proposed an interesting framework based on structured patterns technique for aggregating input alerts in real-time.

## IV. DATASETS AND EVALUATING METRICS

### A. Datasets

For evaluation purposes, several challenging datasets are provided for researcher working in this field and several of these datasets are publicly available to be used. In the following, some of well known datasets are presented.

- **ISCX dataset [28]**

The ISCX (Information Security Centre of Excellence) is one of the widely used dataset. Records are defined by simulation and based on eleven features.

- **TUIDS dataset [29]**

The TUIDS dataset [29] was prepared by the University of Tezpur, in which several attack scenarios are performed. The used representative features are labeled into normal or attack.

- **KDDCup'1999 dataset [30]**

The KDDCup'1999 is the most important and used dataset for IDS performance evaluation. It is generated through several simulations and contains more than 4 million records. Each records is defined on the basis of 41 features either as normal or abnormal attacks.

- **CICIDS'2017 dataset [31]**

The CICIDS'2017 dataset is created for Cybersecurity and it contains both attack and normal scenarios as for the case of ISCX dataset.

- **Kyoto 2006+ dataset [32]**

This dataset is also a challenging benchmark used for real traffic data analysis. It involves 24 features.

### B. Evaluating Metrics

When evaluating IDSs, several factors should be considered such as the cost, the ease-of-use, the speed, the memory/CPU, the effectiveness, and scalability. The performance is usually evaluated and expressed using the following metrics: True positives (TP), True negatives (TN), False positives (FP), False negatives (FN), sensitivity (or True positive rate: TPR), specificity (or True negative rate), and precision. These metrics are defined as follows:

- True Positive Rate (TPR) = $TP/(FN + TP)$.
- False Negative Rate (FNR) = $TN/(FP + TN)$.
- False Positive Rate (FPR) = $FP/(FP + TN)$.
- Accuracy = $(TN + TP)/(TP + FP + TN + FN)$.

### V. Machine Learning (ML) Perspectives

Machine learning-based approaches (Bayesian approaches, Neural Networks, Statistical mixture models, SVM, Hidden Markov model, genetic algorithms, etc.) [33], [34], [35], [36], [37], [38], [39], [1], [2] have been proposed as a powerful techniques to solve several issues related to IDS, alert classification and intrusion detection problems. In particular, they are considered as effective tools for complex data modeling able to represent alerts in a compact form, to filter and to reduce the huge quantity of false alerts and to identify abnormal activities. Moreover, they offer high flexibility to train classifiers and to identify attacks based on a well predefined or extracted specific features. Their use, which is based on the using of a prior and newly acquired information, has proven to be of great importance in this growing area in order to improve the performance of IDS. In the literature, numerous machine learning-based algorithms were implemented for alert classification/clustering [35], [39], [2]. In particular, support vector machines (SVM) is widely employed since it is able to filter efficiently false alert and also it is considered by an important number of researchers in the context of intrusion alert management [40], [41], [42]. Indeed, an SVM-based network intrusion detector is implemented in [40] and its performance is well studied in

[41]. A system for alert and attacks grouping is also developed by [43]. In the subsequent work, the expectation maximization (EM) algorithm is investigated to combine resulted groups into one single attack. Probabilistic models are also investigated online alert aggregation [44]. The later work used the so-called maximum likelihood method to estimate the statistical model's parameters. An anomaly-based algorithm that uses a discriminative machine learning model is implemented to detect intrusions attempts [45]. In fact, input intrusions can be modeled as outliers via a principled probabilistic approach. Moreover, finite mixtures models are mainly used to detect both previously seen (known) and unknown attacks. The same authors proposed another interesting classifier in [2] for online intrusion detection.

### VI. Machine Learning (ML) based Intrusion Alert Clustering

According to the literature review, many works show that machine learning approaches can be very useful for intrusion alert clustering and outlier detection [44], [43], [46]. A lot of works share some common steps which are described in the following sub-sections. Thus, the main objective of this section is to present for the interested reader how can to design a unified ML-based solution that includes several steps. In particular, a case study will be described throughout next sections.

### A. ML-based Framework Design

As shown in Fig. 5, a possible generic solution based on machine learning concepts for intrusion detection and management, alert clustering/classification and outlier detection is presented. Useful technical details related to the implementation of this solution are also provided.

- A preprocessing step: This step is necessary in order to unify and rank all alert information. Moreover, at this level, it is suitable to translate all alert's attributes into numerical values because some of them are in the form of non-numerical values such as SourceIPaddress, DestinationIPaddress, ServiceProtocol and AlertType. These attributes must be mapped into a numerical value.

- A feature extraction/selection step: This step is often used to simplify and accelerate further processing, it would be better to select only significant alert features. Determining an optimal feature set while preserving high accuracy is a challenging problem. To deal with this problem, several algorithms were developed in the literature. Most of them study the relationships between alerts.

- A clustering of normal/abnormal alerts step: Tthis is the main important step and the challenge question is how to develop an accurate intrusion detection model for both alert clustering and abnormal alert detection (outlier detection)? Many supervised and unsupervised machine learning algorithms have been applied to solve this issue. A good choice of a machine learning (ML) technique helps in obtaining effective clustering results.

Fig. 5. General ML-based framework for alert clustering.

### B. Problem Modeling

According to the literature review, many works suppose that attacks may be considered as a random processes generating alerts [44], [43], [46]. For instance, one of the most interesting techniques that can be applied is mixture models associated with maximum likelihood principle. More specifically, the method "expectation-maximization: EM" [44] can be a reasonable choice. This method has the advantage to avoid the restrictions imposed by other algorithms and can aggregate efficiently similar alerts. On the other hand, it would be more interesting if someone considers an effective strategy called "a bootstrap sampling strategy" within the EM algorithm in order to improve the overall process and to speed up the processing time of aggregating similar-alerts from the output of IDSs. As a result, an optimal representative set of input alerts will be determined according to some well defined criteria thanks to the strategy of bootstrap sampling.

*1) Alert features selection:* In order to study the relationships between alerts, it is indispensable to analyze their attributes (features). Among all of them, only few features contribute mainly to this relationship. Hence, identifying main features is a crucial step for further processing. To address this problem, several algorithms were developed, and many of them consider a lot alert features in the process which is very expensive. Furthermore, the correlation procedure cannot make use a big number of attributes given the restrictions imposed by the high-speed networks environment. Thus, determine an optimal feature set while preserving high accuracy is a challenging problem for alert management process. To meet this challenge, it seems adequate to follow a procedure that involves for instance two machine learning-based algorithms: principal component analysis (PCA) and a multi-class support vector machine. Why PCA ? since its components are orthogonal to each other and this characteristic has proven to be a useful statistical technique for dimension reduction and multivariate analysis [47] and guarantees a robust convergence and speedup training as confirmed in [48]. SVM is considers as a robust technique especially when dealing with big data. SVM is scalable and has high performance when compared to existing methods such as artificial neural networks (ANN). Now how can these two techniques be used ? First, PCA can be used to select an optimal subset of most relevant attributes.

Then, a multi-class support vector machine can be applied to classify alerts into meaningful clusters based on the selected features. If the selected features are not sufficient and the clustering step fails, then, additional attributes are required to increase the clustering precision and accuracy. This process will be repeated until finding a good compromise between a high performance and a small number of alert features. The following algorithm can be used for alert's attributes election.

```
program- Alert's Attributes Selection
  begin
      Run the PCA-algorithm;
      Rank Alert's Attributes
      Determine initial subset of attr.
      Fa :=initial attributes ;
      Fa := {F1, F2,....,Fm}; (m < 41)
      Fr := {All Attributes} - Fa ;
      NumberSelectedAttributes :=  m;
  Repeat
      {Classify dataset using SVM in N
      classes:(Normal, DoS, U2R, R2L, Probe)}
      {Compute the classification accuracy}
      IF (Accuracy < \epsilon) then
          {Select best Attr. with best rank.}
          bestAttribute := fb;
          Fa :=  Fa + fb;
          NumberSelectedAttributes ++ ;
      Else
          {Return final selected Attr. }
          return Fa ;
      End IF
  Until convergence (accuracy is achieved)
end.
```

*2) Dimensionality reduction:* To speed up the step of alert clustering, it would be interesting to reduce the data dimension by taking, for example, into account a preprocessing step of data sampling. Among of the motivating techniques, the "bootstrap-sampling" [49] can be examined. Bootstrap is a data resampling method which was introduced as a tool for estimating the sample distribution of statistics. It is applied successfully in many pattern classification problems. The key idea of the Bootstrap is to generate new samples (random

samples) to replace original data. The process of determining a random sample is repeated many times until finding an empirical distribution of the statistic. The process of sampling has to reduce the complexity of clustering algorithms (e.g the EM algorithm). From a technical point of view, the initial dataset will be replaced by only a small samples which are closely representative to the initial dataset. Then, the clustering algorithm will estimate statistics on one of these samples. With this manner, statistics can be calculated easily and in "real time". Details concerns the process of the bootstrap sampling combined with the EM algorithm are given as follows:

Let's consider that initial data contains n alerts noted by $A = (A_1, ..., A_n)$. The key problem is how to find randomly a representative Bootstrap sample from the initial dataset denoted $A^* = (A_1^*, ..., A_n^*)$. This problem is solved in the following way:

1) Step 1: From the initial sample $A = (A_1, ..., A_n)$, create an empirical probability distribution $F$ which consists in placing the probability of $1/n$ for each alert (all alerts have the same probability).

2) Step 2: Given the empirical distribution function, F, (original data set), generate a new random sample of size n with replacement: this is called the "bootstrap resample".

3) Step 3: Calculate the model statistics through the EM algorithm for this resample $(\theta^*)$ instead of the original sample$(\theta)$.

4) Step 4: Repeat steps 2 and 3 $B$-times ($B$ is the number of bootstrap samples) in order to generate B resamples and to obtain an approximation of the distribution. The size of B depends on the tests to be runned on the data.

Now, to estimate the optimal sample size, one can take advantage of some criteria [50]. These criteria can be applied easily in the context of intrusion alert clustering. The appropriate size of the strapped sample is determined as follow:

- K : represents the number of alerts having different attributes each others,

- $\pi_i$ is the a priori probability of a particular alert.

- $\epsilon$ is a fixed small value,

The probability $P_i$ that a particular alert "i" which is from the sample is given by: $P_i = 1 - (1 - \pi_i)^n$. If the condition $(n\pi_i > 4)$, the probability $P_i$ can be approximated as: $P_i = 1 - e^{-n\pi_i}$. According to this condition, if we have many similar alerts then at least one of them should exist in the sample. It is equivalent to maximize the joint probability: $P_n = \prod_{i=1}^{K}(1 - e^{-n\pi_i})$. This problem is equivalent to minimize the derived of the logarithm the joint probability. Now, denote by $n_0$ the optimal bootstrap sample size. The value of $n_0$ is determined as:

$$\begin{cases} Size(n_0) = \sum_{i=1}^{K} \frac{\pi_i e^{-n_0 \pi_i}}{1 - \pi_i e^{-n_0 \pi_i}} < \epsilon \\ \\ n_0 > 4K \end{cases}$$

*3) Intrusion alert clustering:* At this stage, the challenge question is how to design an accurate intrusion detection model for both alert clustering and abnormal alert detection (outlier detection)? Many supervised and unsupervised machine learning algorithms have been applied to solve this issue. In particular, using an enhanced version of the EM algorithm which is combined with the bootstrap sampling making it an attractive solution. The EM is one of the most frequently used technique for estimating the probability density functions (PDF) in both univariate and multivariate cases. It is used especially to model a set of feature vectors by a mixture of statistical distributions. These distributions are then used to model the observation vectors. There are some researchers who have tried to apply EM-algorithm for alert clustering such as in [51], [52], [38], [53]. From a technical point of view, the distribution of the generated alerts can be approximated according to multivariate probability distribution given that an attack instance is considered as a random process. Let's consider an alert $A$ consists of d attributes. For example, for the case of Gaussian distribution, the density function is defined as: $f(A_d; \theta_k) = \frac{1}{\sqrt{2\pi|\Sigma|_k}} e^{-\frac{1}{2}(a_d - \mu_k)^T \Sigma_k^{-1}(a_d - \mu_k)}$ Where:

- $\theta_k = (\mu_k, \Sigma_k))$.

- $a_d$ represents the $d^{th}$ feature of the alert a.

- $\mu_k$ and $\Sigma_k$ are respectively the mean and the covariance matrix.

The posteriori probability to be calculated for each alert corresponds to the labeled classes that should be determined. The aim is to assign the best label $L_n$ to each alert $A$ where$L_n \in \{c_1, ..., c_k\}$, $c_n$ are the classes of the mixture model and k is the number of classes. The output of the problem are: model's parameters associated to each class which are $\theta_k = (\mu_k, \Sigma_k)$; and a posteriori probability $\gamma_k$ for each alert $A$. The algorithm is based on two main steps: E (expectation) and M (maximization). In E-step, incoming alerts are assigned to classes which leads to a compact partition $A$ with $K$ classes. The second step is the M step, where an optimal values of the parameters of the developed model is determined. The main steps of the EM algorithm are illustrated in the following pseudo-code.

```
begin
    1.  Initialization-step:
```

$$\pi_k := \frac{1}{K} \tag{1}$$

```
    2.  Expectation-step:
```

$$\gamma_{nk}^{(q)} := P(c_k/A_n, \theta_k^{(q)}) := \frac{\pi_k^{(q)} f_k(A_n; \theta_k^{(q)})}{\sum_{l=1}^{K} \pi_l^{(q)} f_l(A_n; \theta_l^{(q)})} \tag{2}$$

```
    3.  Maximization-step:
```

$$\begin{cases} \pi_k^q := \frac{\sum_{n=1}^{N} \gamma_{nk}}{N} \\ \mu_k^q := \frac{\sum_{n=1}^{N} \gamma_{nk} A_n}{\sum_{n=1}^{N} \gamma_{nk}} \\ \sigma_k^q := \frac{\sum_{n=1}^{N} \gamma_{nk}(A_n - \mu_k^{(q)})^2}{\sum_{n=1}^{N} \gamma_{nk}} \end{cases} \quad (3)$$

N: Number of alerts.

$$\begin{cases} \sum_{i=1}^{K} \pi_i = 1 \\ \pi_i \geq 0 \end{cases} \quad (4)$$

```
K :number of components (classes)
4.  Repeat step 2 and 3 until
convergence based on this criterion:
```

$$|\gamma_{nk}^{(q+1)} - \gamma_{nk}^{(q)}| < \epsilon \quad (5)$$

### c) *Online Classification Process*

It is possible to extend the previous algorithm to an online fashion in order to allow alert classification in real-time. The online classification is a new version over the classical batch version where the parameters are re-calculated each time a new alert is coming. The main reason why an online algorithm is desirable is to avoid huge computational and memory savings. Thus, it does not require the whole data set to be available at each iteration. For online clustering, the mixture model parameter estimates from each previous iteration are used to initialize the next iteration. If a new alert is observed, it can be associated with an existing cluster or to a new created one. This process is performed on the basis of the most likely component using the obtained measures from the E-step of the EM-algorithm. If the alert cannot be assigned to an existing existing component (i.e. there is no similarity with previous alerts), it is assumed as a new alert instance and therefore a new component will be created for it. A possible pseudo-code for online alert classification can be summarized as follows:

```
program for Online Alert Classification

 begin
  While (new alert "a" is received) do
  - find most likely component for "a"
```

$$k^* := argmax_k(\gamma_k(\theta_k)) \quad (6)$$

```
  - add the new alert to the component
```

$$C_{old-k} := C_{k^*}$$
$$C_{k^*} := C_{k^*} \cup \{a\}$$

```
  - update the new statistics:
```

$$\begin{cases} \pi_{k^*} \\ \mu_{k^*} \\ \sigma_{k^*} \end{cases} \quad (7)$$

```
      if |θ_k − θ_k*| ≥ Threshold

        - discard previous changes:
```

$$C_{k^*} := C_{old-k}$$

```
        - create new component for new alert
        - update number of components
      else
        - accept changes
```

$$C_{k^*} := C_{k^*} \cup \{a\}$$

Finally, redundant alerts in each cluster can be fused into a so-called "Meta-Alert". Redundant alerts are supposed have equal attribute values. Meta-Alerts are needed for the security expert reports and may be investigated further in order to detect more complex attack scenarios. A typical algorithm for meta-alert generating is given in the following:

```
program for Meta-alert generating

 begin
    MergeAlerts := 0;
    Repeat for each class Ck
      Repeat For each alert Ai in Ck
        IF (Attr(Ai) == Attr(Ai+1))
            Delete  alert  Ai;
            Meta-Alert := Ai+1;
            MergeAlerts := MergeAlerts++;
      End IF
    End Repeat
  End Repeat
 end.
```

*4) Outlier detection:* Outliers are anomalies' observations that do not conform to the normal behavioral of the dataset and deviate a lot from the other observations since their values are very different from the data values. Given that, some statistical parameters such as the mean and the standard deviation are sensitive to outlier detection; it is recommended to apply for example more robust well known distance measures such as the so-called "Mahalanobis distance" to filter false positive alerts (outliers). Mahalanobis measure is a multidimensional version of a z-score which is based on clustering between variables and depends on estimated parameters of the multivariate distribution. It measures the distance of any alert to the center alert (multidimensional mean of the alert-class), given the covariance (multidimensional variance). The Mahalanobis distance is scale-invariant (not dependent on the scale of measurements) and takes into account the correlations of the data set. These properties are not retained by the classical Euclidean distance. Formally, the Mahalanobis distance between a particular multivariate alert vector "a" and the mean value $\mu$ is defined as:

$$MDist(a, \mu) = \sqrt{(a - \mu)^T \sum{}^{-1}(a - \mu)} \quad (8)$$

Where $\Sigma$ is the covariance matrix of the "normal" data. All data alerts which are located far away from the mean alert (center of the data) are considered as outliers. In other word, multivariate outliers $a_i$ have large distance value $MDist$. Formally, the Mahalanobis distance follows the chi-square distribution with p degrees of freedom $\chi_p^2$. So, an alert is considered as outlier if we have $MDist(a, \mu) > \chi(0.975)$.

```
program for Alert-Outliers Detection
  {Input : different alert classes.
   Output: alerts as outliers. }
 begin
    - Compute mean value of each class;
    - Compute the covariance matrix
     Repeat for each class Ci,
        For each alert aj in Ci,
          - Calculate Mahalanobis distance
            between ai and mean(MDist).
            If( MDist > chi(0.975) ) Then
              - Set aj as Outlier-Alert;
              - Delete aj from class Ci;
            End IF
        End For
     End Repeat
 end.
```

## VII. Conclusion and Discussion

The main important problem related to the developed intrusion detection systems (IDSs) in the literature is that they produce a lot of false positive alerts for the same attack. Therefore, it becomes too difficult to distinguish between normal and abnormal alerts, to classify them accurately and to take correct actions for them. A lot of promising researches have been developed but many of them have apparent limitations. It is noted that despite more than twenty years' efforts on the field of intrusion detection systems, a lot of issues still not yet solved. For example, some developed methods for IDSs are not qualified to recognize all kind of intrusions (anomalies), they cannot ensure that all alerts are true positives, and they fail to identify main malicious activities. In this study, a deep review of some well known clustering methods is presented. Then, a particular focus is dedicated for machine learning techniques which are considered as attractive alternatives to address main issues related to IDS systems. Moreover, this paper presents useful technical details for designing and implementing a unified framework by taking part some effective machine-learning algorithms. Such framework can be helpful for interested readers in this field of research.

## References

[1] H. Sallay, A. Ammar, M. B. Saad, and S. Bourouis, "A real time adaptive intrusion detection alert classifier for high speed networks," in *2013 IEEE 12th International Symposium on Network Computing and Applications, Cambridge, MA, USA, August 22-24, 2013*, 2013, pp. 73–80.

[2] H. Sallay and S. Bourouis, "Intrusion detection alert management for high-speed networks: current researches and applications," *Security and Communication Networks*, vol. 8, no. 18, pp. 4362–4372, 2015.

[3] K. Alsubhi, I. Aib, and R. Boutaba, "Fuzmet: a fuzzy-logic based alert prioritization engine for intrusion detection systems," *Int. J. Netw. Manag.*, vol. 22, no. 4, pp. 263–284, 2012.

[4] F. Valeur, *Real-Time Intrusion Detection Alert Correlation*. University of California, Santa Barbara: Phd thesis, 2006.

[5] W. Alhakami, A. ALharbi, S. Bourouis, R. Alroobaea, and N. Bouguila, "Network anomaly intrusion detection using a nonparametric bayesian approach and feature selection," *IEEE Access*, vol. 7, pp. 52 181–52 190, 2019.

[6] F. Cuppens, "Managing alerts in a multi-intrusion detection environment," in *Proceedings of the 17th Annual Computer Security Applications Conference*, ser. ACSAC '01, 2001, pp. 22–31.

[7] H. Debar and A. Wespi, "Aggregation and correlation of intrusion-detection alerts," in *Recent Advances in Intrusion Detection*, 2001, pp. 85–103.

[8] J. Yu, Y. R. Reddy, S. Selliah, S. Reddy, V. Bharadwaj, and S. Kankana-halli, "Trinetr: An architecture for collaborative intrusion detection and knowledge-based alert evaluation," *Advanced Engineering Informatics*, vol. 19, no. 2, pp. 93 – 101, 2005.

[9] M. E. Whitman and H. J. Mattord, *Principles of Information Security*, 3rd ed. Boston, MA, United States: Course Technology Press, 2007.

[10] F. Valeur, G. Vigna, C. Kruegel, and R. A. Kemmerer, "A comprehensive approach to intrusion detection alert correlation," *IEEE Trans. Dependable Secur. Comput.*, vol. 1, no. 3, pp. 146–169, 2004.

[11] C. V. Zhou, C. Leckie, and S. Karunasekera, "A survey of coordinated attacks and collaborative intrusion detection," *Computers and Security*, vol. 29, no. 1, 2010.

[12] C.-F. Tsai, Y.-F. Hsu, C.-Y. Lin, and W.-Y. Lin, "Intrusion detection by machine learning: A review," *Expert Systems with Applications*, vol. 36, no. 10, pp. 11 994 – 12 000, 2009.

[13] B. Morin and H. Debar, "Correlation of intrusion symptoms: an application of chronicles," in *In Proceedings of the 6th International Conference on Recent Advances in Intrusion Detection (RAID'03)*, 2003, pp. 94–112.

[14] O. Dain and R. K. Cunningham, "Fusing a heterogeneous alert stream into scenarios," in *In Proceedings of the 2001 ACM workshop on Data Mining for Security Applications*, 2001, pp. 1–13.

[15] A. Valdes and K. Skinner, "Probabilistic alert correlation," in *Recent Advances in Intrusion Detection*, 2001, pp. 54–68.

[16] K. Julisch, "Clustering intrusion detection alarms to support root cause analysis," *ACM Trans. Inf. Syst. Secur.*, vol. 6, no. 4, pp. 443–471, 2003.

[17] S. Staniford, J. A. Hoagland, and J. M. McAlerney, "Practical automated detection of stealthy portscans," *Journal of Computer Security*, vol. 10, no. 1/2, pp. 105–136, 2002.

[18] O. M. Dain and R. K. Cunningham, "Building scenarios from a heterogeneous alert stream," in *IEEE Workshop on Information Assurance and Security*, 2001, pp. 231–235.

[19] F. Cuppens and A. Miège, "Alert correlation in a cooperative intrusion detection framework," in *Proceedings of the 2002 IEEE Symposium on Security and Privacy*, 2002, pp. 202–2015.

[20] P. Ning, Y. Cui, and D. S. Reeves, "Constructing attack scenarios through correlation of intrusion alerts," in *Proceedings of the 9th ACM conference on Computer and communications security*, 2002, pp. 245–254.

[21] P. Ning, Y. Cui, D. S. Reeves, and D. Xu, "Techniques and tools for analyzing intrusion alerts," *ACM Trans. Inf. Syst. Secur.*, vol. 7, no. 2, pp. 274–318, 2004.

[22] S. J. Templeton and K. Levitt, "A requires/provides model for computer attacks," in *Proceedings of the 2000 workshop on New security paradigms*, New York, NY, USA, 2000, pp. 31–38.

[23] N.-Y. Jan, S.-C. Lin, S.-S. Tseng, and N. P. Lin, "A decision support system for constructing an alert classification model," *Expert Systems with Applications*, vol. 36, no. 8, pp. 11 145 – 11 155, 2009.

[24] X. D. Hoang, J. Hu, and P. Bertok, "A program-based anomaly intrusion detection scheme using multiple detection engines and fuzzy inference," *Journal of Network and Computer Applications*, vol. 32, no. 6, pp. 1219 – 1228, 2009.

[25] G. P. Spathoulas and S. K. Katsikas, "Reducing false positives in intrusion detection systems," *Computers and Security*, vol. 29, no. 1, pp. 35 – 44, 2010.

[26] R. Sadoddin and A. A. Ghorbani, "An incremental frequent structure mining framework for real-time alert correlation," *Computers and Security*, vol. 28, pp. 153 – 173, 2009.

[27] F. Maggi, M. Matteucci, and S. Zanero, "Reducing false positives in anomaly detectors through fuzzy alert aggregation," *Information Fusion*, vol. 10, no. 4, pp. 300 – 311, 2009.

[28] A. Shiravi, H. Shiravi, M. Tavallaee, and A. A. Ghorbani, "Toward developing a systematic approach to generate benchmark datasets for intrusion detection," *Computers & Security*, vol. 31, no. 3, pp. 357–374, 2012.

[29] P. Gogoi, M. H. Bhuyan, D. K. Bhattacharyya, and J. K. Kalita, "Packet and flow based network intrusion dataset," in *Contemporary Computing - 5th International Conference, IC3 2012, Noida, India, August 6-8, 2012. Proceedings*, 2012, pp. 322–334.

[30] M. Tavallaee, E. Bagheri, W. Lu, and A. A. Ghorbani, "A detailed analysis of the KDD CUP 99 data set," in *2009 IEEE Symposium on Computational Intelligence for Security and Defense Applications, CISDA 2009, Ottawa, Canada, July 8-10, 2009*, 2009, pp. 1–6.

[31] I. Sharafaldin, A. H. Lashkari, and A. A. Ghorbani, "Toward generating a new intrusion detection dataset and intrusion traffic characterization," in *Proceedings of the 4th International Conference on Information Systems Security and Privacy, ICISSP 2018, Funchal, Madeira - Portugal, January 22-24, 2018.*, 2018, pp. 108–116.

[32] J. Song, H. Takakura, Y. Okabe, M. Eto, D. Inoue, and K. Nakao, "Statistical analysis of honeypot data and building of kyoto 2006+ dataset for NIDS evaluation," in *Proceedings of the First Workshop on Building Analysis Datasets and Gathering Experience Returns for Security, BADGERS@EuroSys 2011, Salzburg, Austria, April 10, 2011*, 2011, pp. 29–36.

[33] A. Patcha and J.-M. Park, "An overview of anomaly detection techniques: Existing solutions and latest technological trends," *Comput. Netw.*, vol. 51, no. 12, pp. 3448–3470, 2007.

[34] F. Najar, S. Bourouis, N. Bouguila, and S. Belghith, "A fixed-point estimation algorithm for learning the multivariate ggmm: application to human action recognition," in *2018 IEEE Canadian Conference on Electrical & Computer Engineering (CCECE)*, 2018, pp. 1–4.

[35] S. Zanero and S. M. Savaresi, "Unsupervised learning techniques for an intrusion detection system," in *Proceedings of the 2004 ACM symposium on Applied computing*, 2004, pp. 412–419.

[36] S. Bourouis, A. Zaguia, N. Bouguila, and R. Alroobaea, "Deriving probabilistic SVM kernels from flexible statistical mixture models and its application to retinal images classification," *IEEE Access*, vol. 7, pp. 1107–1117, 2019.

[37] T. Pietraszek and A. Tanner, "Data mining and machine learning-towards reducing false positives in intrusion detection," *Inf. Secur. Tech. Rep.*, vol. 10, no. 3, pp. 169–183, 2005.

[38] F. Najar, S. Bourouis, N. Bouguila, and S. Belghith, "Unsupervised learning of finite full covariance multivariate generalized gaussian mixture models for human activity recognition," *Multimedia Tools and Applications*, pp. 1–23, 2019.

[39] P. Laskov, P. Dussel, C. Schafer, and K. Rieck, "Learning intrusion detection: supervised or unsupervised," in *IMAGE ANALYSIS AND PROCESSING, PROC. OF 13TH ICIAP CONFERENCE.*, 2005, pp. 50–57.

[40] S.-J. Horng, M.-Y. Su, Y.-H. Chen, T.-W. Kao, R.-J. Chen, J.-L. Lai, and C. D. Perkasa, "A novel intrusion detection system based on hierarchical clustering and support vector machines," *Expert Syst. Appl.*, vol. 38, no. 1, pp. 306–313, 2011.

[41] D. Fisch, A. Hofmann, and B. Sick, "On the versatility of radial basis function neural networks: A case study in the field of intrusion detection," *Information Sciences*, vol. 180, no. 12, pp. 2421–2439, 2010.

[42] L. Khan, M. Awad, and B. Thuraisingham, "A new intrusion detection system using support vector machines and hierarchical clustering," *The VLDB Journal*, vol. 16, no. 4, pp. 507–521, 2007.

[43] R. Smith, N. Japkowicz, M. Dondo, and P. Mason, "Using unsupervised learning for network alert correlation," in *21st conference on Advances in artificial intelligence*, ser. Canadian AI'08, 2008, pp. 308–319.

[44] A. Hofmann and B. Sick, "Online intrusion alert aggregation with generative data stream modeling," *IEEE Transactions on Dependable and Secure Computing*, vol. 8, no. 2, pp. 282–294, 2011.

[45] H. Sallay, S. Bourouis, and N. Bouguila, "Web service intrusion detection using a probabilistic framework," in *Progress in Systems Engineering*, vol. 1089, 2015, pp. 161–166.

[46] T. Shon and J. Moon, "A hybrid machine learning approach to network anomaly detection," *Inf. Sci.*, vol. 177, no. 18, pp. 3799–3821, 2007.

[47] I. Jolliffe, Ed., *Principal Component Analysis*, ser. 3rd ed. Springer-Verlag, New York, 2002.

[48] E. Oja, "Neural networks, principal components, and subspaces," *Int. J. Neural Syst.*, vol. 1, no. 1, pp. 61–68, 1989.

[49] B. Efron, "Better bootstrap confidence intervals," *Journal of the American Statistical Association*, vol. 82, no. 397, pp. 171–185, 1987.

[50] C. Banga and F. Ghorbel, "Optimal bootstrap sampling for fast image segmentation: application to retina image," in *IEEE international conference on Acoustics, speech, and signal processing*, 1993, pp. 638–641.

[51] A. Hofmann and B. Sick, "Online intrusion alert aggregation with generative data stream modeling," *IEEE Trans. Dependable Sec. Comput.*, vol. 8, no. 2, pp. 282–294, 2011.

[52] S. Bourouis, Y. Laalaoui, and N. Bouguila, "Bayesian frameworks for traffic scenes monitoring via view-based 3d cars models recognition," *Multimedia Tools and Applications*, pp. 1–21, 2019.

[53] I. Channoufi, S. Bourouis, N. Bouguila, and K. Hamrouni, "Image and video denoising by combining unsupervised bounded generalized gaussian mixture modeling and spatial information," *Multimedia Tools Appl.*, vol. 77, no. 19, pp. 25 591–25 606, 2018.

# Feature based Algorithmic Analysis on American Sign Language Dataset

Umair Muneer Butt[1], Basharat Husnain[2], Usman Ahmed[3], Arslan Tariq[4],
Iqra Tariq[5], Muhammad Aadil Butt[6], Dr. Muhammad Sultan Zia[7]
Department of Computer Science
University of Lahore
Chenab Campus, Pakistan

*Abstract*—**Physical disability is one of the factor in human beings, which cannot be ignored. A person who can't listen by nature is called deaf person. For the representation of their knowledge, a special language is adopted called 'Sign-Language'. American Sign Language (ASL) is one of the most popular sign language that is used for learning process in deaf persons. For the representation of their knowledge by deaf persons, a special language is adopted 'Sign-Language'. American Sign Language contains a set of digital images of hands in different shapes or hand gestures. In this paper, we present feature based algorithmic analysis to prepare a significant model for recognition of hand gestures of American Sign Language. To make a machine intelligent, this model can be used to learn efficiently. For effective machine learning, we generate a list of useful features from digital images of hand gestures. For feature extraction, we use Matlab 2018a. For training and testing, we use weka-3-9-3 and Rapid Miner 9 1.0. Both application tools are used to build an effective data modeling. Rapid Miner outperforms with 99.9% accuracy in auto model.**

*Keywords*—**Hand gesture recognition; pre-processing; weka; rapid miner; HOG; LBP; auto model**

## I. Introduction

Sign language provides a big aid and convenience in human life [1] and used especially by deaf persons and by other people to add weight in conversation. Visual representation by hands, delivers a meaningful message to others [2]. Sign Language consists in three forms: one is called facial expression, second is hand gestures and third is called body postures [1], [2]. In our daily life, we mostly use our body postures and facial expression to deliver meaningful information to others. The goal of the communication is achieved when the senders message fully interpreted by the receiver with full of emotions. Hand gestures and facial expressions play an important role in the learning process of deaf persons. Sign language is greatly influenced by hand gestures recognition. Hand gesture plays vital role in understanding sign language [1]. It can be taken from live camera in the form of moving hand gestures or in the form of still images [3]. In our research we will consider only still images of hand gestures.

The persons who can't listen by birth are called deaf persons. Deaf persons can't listen any voice through their ears.Teaching them verbally is not effective way of communication. There is a need of special language for their learning purpose. That is called a "Sign Language". Sign Language is used to understand the conveyed message from others. American Sign Language has 24 different hand postures. Each



Fig. 1. Hand Gestures[21]

posture shows a unique ASL letter. The following Fig. 1 shows American Sign Language alphabets. Sign Language field is very vast. The study of Hand gestures are always being a very tough to learn. A machine could not be recognized gestures until or unless the machine is professionally trained.

The above data set of sign American Sign Language is taken form a well-known website "kaggle" [21]. The data set did not contain the letter "J" and "Z". It is because visually similarity of these two signs with others.

According to a rough counting in a research [3], almost there are 500,000 to 2,000,000 deaf people's uses the sign language for communication with one another. The counting figure may be different from other proved research, but everyone would be agreed that the sign language is at the third most wanted and most used all over the world [3].

We can build a model to recognized hand gestures using

Fig. 2. Skin detection.

different techniques. In the past, developers used a finger technique, in which a user uses a finger mouse to capture fingers [4], using skin colour detection from any useful algorithm [3], gloves technique that was used neural net [5], feature extraction technique using Scale-Invariant Feature Transform(SIFT) algorithm [6]. All of these techniques are very tough to implement. From the above source [3] talked, a skin detection algorithm is used to detect the skin colour.

A special environment is created for skin detection with sufficient lightning conditions. There are some constraints that needs to be satisfied. First, the background color must be different from skin color. Second, algorithms fails to perform well under different backgrounds and colored clothes. For skin detection, the user should be there in a specially created environment in which a sufficient light was required. In case of less light or different background and cloth colours, the skin detection algorithm did not work properly and did not detect skin properly. Fig. 2 shows the detail. Just like skin detection algorithms, the neural net and SIFT techniques are also difficult to implement. A neural net algorithm takes a lot of time to process digital images.

In this paper, we present simple but efficient technique for ASL recognition. We provide a comprehensive analysis on different techniques with feature extraction and different algorithms. we use tools like (Weka and Rapid Miner) and achieved 99% accuracy on test data. Methods used in these techniques, experimental results and assumptions are described in coming Sections II, III, IV and V accordingly.

## II. LITERATURE REVIEW

Focus of our research work is sign language recognition using hand gestures. It is important to understand gestures

so that true semantics of communication can be grasped. According to a rough counting in a research [3], there are more than 500,000 to 2,000,000 deaf people's using the sign language for communication with one another. The counting figure may be different from other researcher's research, but everyone would be agreed that the sign language is at the third most wanted and most used all over the world [3]. There are different sign languages such as American Sign Language (ASL), Indian Sign Language (ISL), Arabic Sign Language (ArSL), Tamil Sign Language (TSL), Koran Sign Language (KSL), Japanese sign Language (JSL) and many more [1]. In our research work, we are focusing on American Sign Language (ASL).

Gesture recognition was used in 1993 for the first time [3]. Later for recognizing dynamic gestures, Dynamic Time Wrapping (DTW) technique was used [4]. Hidden Markov Model (HMM) was also used for recognizing sign language's shape [5], [6]. They used HMM efficiently and accuracy of sign language recognition reached to 94%. Later, it was found that accuracy dropped to 47.6%, when system was by a person other than those images were used for training. If both person's images are used for training then accuracy level increased [7]. Major limitation of HMM was its context dependency. HMM was used with 3D data to classify 53 ASL and attained accuracy of 89.91% [8].

Image acquisition and Pre-processing is the backbone of gesture recognition. In the Past, image is acquired using Leap Motion Controller (LMC), Kinect and vision based approaches [19]. LMC can acquire signals 200 frames per sec [18]. It has been widely used for hand gesture recognition tasks [20].

In the past, researchers used many methods for recognizing hand gestures. Some used a finger technique, in which a user uses a finger mouse to capture fingers [9], using skin colour detection from any useful algorithm [10], [11], gloves technique that used neural nets [12], [13], feature extraction technique using "SIFT" algorithm [14], [15]. Viola- Jones method was used for detecting skin, skin colour was used to detect hand. After hand detection features are extracted using SIFT and Support Vector Machine(SVM) is used for classification purpose [15].

Skin detection techniques are more sensitive in this process, the user should be there in a specially examined environment, which requires a specific intensity light. There are various other constraint like: The background colour must be different form skin colour, light should be constant and background and clothes should be simple.Various algorithms fails to perform in skin detection, if these condition cannot be fulfilled [9]. Fig. 2 shows the detail. Just like skin detection algorithms, the neural net and SIFT techniques are also not efficient both in accuracy and time. A neural net algorithm takes relatively more time than other techniques to process digital images [4]. So, Neural net is not suitable for real time skin detection [16]. K Nearest Neighbour algorithm was used with PCA and achieve 96% accuracy [17].

## III. PROPOSED METHODOLOGY

The proposed method is based on combination of His- togram of Oriented Features (HOG), Local Binary Patterns (LBP) and statistical features. It focuses on the algorithmic

(a) Some area of Images contains unnecessory contents after reduction

(b) Accurate image after reduction

Fig. 3. Accurate and inaccurate image reduction [21].



Fig. 4. Data Set with labels

| Label in CSV File | Assigning Letters | Label in CSV File | Assigning Letters | Label in CSV File | Assigning Letters |
|---|---|---|---|---|---|
| 0 | A | 8 | I | 17 | R |
| 1 | B | 9 | K | 18 | S |
| 2 | C | 11 | L | 19 | T |
| 3 | D | 12 | M | 20 | U |
| 4 | E | 13 | N | 21 | V |
| 5 | F | 14 | O | 22 | W |
| 6 | G | 15 | P | 23 | X |
| 7 | H | 16 | Q | 24 | Y |

Fig. 5. Labels in CSV File

analysis of different tools and techniques with respect to time and accuracy. Sign language recognition can be divided into four major steps [16].

1) Image Acquisition
2) Image Segmentation
3) Feature Extraction
4) Hand Gesture Recognition

### A. Image Acquisition and Pre-Processing

Data set for proposed work are taken from a well-known source "Kaggle". Other Data sources are also visits but we could not found enough data for hand gestures in digital images or Comma Separated Value (CSV) file format. At Kaggle, we found two data sets for hand gestures one is in the form of set of colour images and other is in the form of CSV file as shown in Fig. 5. Colour images has 9 folders and each folder has 241 colour images with an excel file. Excel file contains image name and the images dimensions from x1, y1, x2 and y2. After Pre-processing (reduction of images according to x1, y1, x2 and y2 given points in excel file) in MATLAB 2018a, we found images are not good as they still contains some unnecessary contents. In these images we found some images did not have proper cutting contents. According to the given dimensions in excel file, some hand gestures was cute and they did not express the accurate meaning of sign language. Fig. 3 shows the cute area of hand gestures.

Data set in the form of CSV files with the following name "signministtrain" and "signministtest" are checked. The training data set contains 27,455 digital images record and test data set contains 7,172 digital images records. These files contain pixel values of a grey scale digital image in the form 785 columns and the last column contains the class of each image. First 784 columns have pixel values ofeach image with dimension 28x28. Have a look of these CSV file in the following Fig. 4. ' First and very important task is to

separate each file into its original graphical form, from its pixel values. For this purpose we use set of instruction in MATLAB to convert each pixel into an image. After analysing the "signministtrain" file, we have the following labels or classes for image dataset. In CSV file 5 each record has a label in numeric format, which means that each numeric digit is represent a sign language letter.

The following algorithm takes each row form csv or excel file from very first record to end of the file and reshape the each row vector from 1x784 columns to 28x28 columns vector. 28x28 column vector stores in an array and convert it into a graphical image file and store it on the given location.

1) Read CSV file and convert into an excel file format.
2) Resize each 1x784 column to 28x28 vector column by reading each record in excel file.
3) Generate digital images by reading each 28x28 vector.
4) Store the file at given location in digital form.

After executing this algorithm we got 27,455 training and 7,172 test images for hand gesture data set. After pre-processing the following Fig. 6 shows data set is generated in digital images. After analysing the "signministtrain" file, we have the following labels or classes for image dataset.

### B. Feature Extraction

We use simple feature extraction technique in this paper to make is simple to simplest. MATLAB 2018a is used for feature extraction techniques. HOG (Histogram of Oriented Gradient) and LBP (Local Binary Pattern) are the important feature extraction techniques using in MATLAB.

Fig. 6.    After pre-processing, hand gestures data set.

TABLE I.    HOG AND LBP FEATURES

| HOG Features | | | LBP Features |
|---|---|---|---|
| Cell Size | Block Size | Num of Bins | Cell Size |
| [8 8] | [2 2] | 9 | [28 28] |
| [8 8] | [2 2] | 9 | [28 28] |
| [8 8] | [2 2] | 9 | [28 28] |
| [8 8] | [2 2] | 9 | [28 28] |

TABLE II.    STATISTICAL FEATURES

| Statistical Features | | | |
|---|---|---|---|
| A.M | STD | Variance | Skewness |
| Yes | No | No | No |
| Yes | Yes | No | No |
| Yes | Yes | Yes | No |
| No | No | Yes | No |

*3) Statistical Feature Measurements:* Based on above two techniques HOG and LBP, we use some additional statistical techniques for better feature extraction as shown in Table II. We use Mean, Standard deviation, variance and skewness for additional feature extraction techniques. The algorithm reads the files directory which contains training data set images. Get each file one by one and extract the HOG, LBP and other statistical feature measurements. Store this features into a CSV file on the specified location. The following algorithm is used for feature extraction in MATLAB 2018a.

1) Read stored images one by one from specified directory.
2) Generate HOG, LBP and other statistical feature measurements.
3) Set labels against each feature.
4) Store features into a CSV file in specified directory.



Fig. 7.    HOG cell size and features.

## IV.    EXPERIMENTS AND RESULTS

Weka-3-9-3 and Rapid Miner 9.1.0 are used for experiments of training and test data models as shown in Fig. 8.

### A. Weka

Weka has a large collection of algorithms for creating effective models in machine learning techniques. Weka provides important facilities in regard of data preparation and classification. Weka also has regression algorithms and clustering algorithms for unsupervised learning. Decision trees and random forest algorithm are also include in weka for supervised learning. We use Navie Bayes, Lazy IBK and Random Forest algorithms.

### B. Rapid Miner

Rapid Miner is the most latest software used for machine learning, data mining, deep learning and text mining. Rapid Miner introduced in 2006 and it has wonderful GUI and provides a lot of options to build a model for machine learning [22]. The algorithms KNN , Neural Net, Generalized Leaner Model, Deep Learning, Naïve Bayes, Random Forest and Decision Trees are used in Rapid Miner 9.1.0 for effective machine learning model.

*1) HOG (Histogram of Oriented Gradient):* Histogram of Oriented Gradient is a very useful technique to extract features in digital image processing. Histogram of Oriented Gradient has a variety of parameter to extract efficient features. Here, we discuss some of these parameters.

*a) HOG Cell Size:* HOG cell size is 2-element vector. It specified the number of pixel in digital image. Cell size [2 2] show that HOG get the features of 2-by-2 pixels. Cell size may vary on different values, it may be 4-by-4, 8-by-8, 16-by16 or 31-by-32 as shown in Table I.In Larger images, we set the cell size large. As we increase the cell size, number of features are also increased. In our feature extraction technique, we use cell size [2, 2], [4, 4] and finally [8, 8] for better experiments. Experiments are shown in experiments and results section.Default value of cell is [8 8]. Fig. 7 shown the number of features on cell size [2 2], [4 4] and [8 8].

*b) Num Bins:* Num Bins describes the number of features value contained in a cell. Default value of Num Bins is 9, which means total 9 number of features will collected from a cell. The following command is used for HOG feature extraction: extractHOGFeature (image path, 'Cell Size' [8 8], 'Block Size' [2 2], Num bins=9).

*2) LBP (Local Binary Pattern):* To encode the texture information, LBP (Local Binary Pattern) technique is used vastly. For LBP feature extraction, the following command is used: No-of-Features = extractLBPFeatures (Image Path).

TABLE III.     EXPERIMENT1 RESULTS

| HOG Features | | | LBP Features |
|---|---|---|---|
| Cell Size | Block Size | Num Bins | Cell Size |
| [4 4] | [2 2] | 9 | [28 28] |
| [8 8] | [2 2] | 9 | [28 28] |
| [8 8] | [2 2] | 9 | No |
| Extracted Features | Results | | |
| No of Features | Lazy.IBK | Random Forest | |
| 1355 | 95.24% | 94.11% | |
| 204 | 96.75% | 94.95% | |
| 145 | 96.50% | 94.24% | |

TABLE IV.     EXPERIMENT2 RESULTS

| Extracted Features | Results | |
|---|---|---|
| No. of Features | Lazy.IBK | Random Forest |
| 232 | 96.83% | 95.60% |
| 260 | 96.40% | 95.20% |
| 288 | 96.19% | 94.30% |
| 232 | 96.87% | 95.39% |

## C. Experiment 1

The following Table III shows the detail of experiment 1. Different results are showing on different parameter settings of cell size as shown in Fig. 9. Lazy.IBK and Random forest gives the accuracy of 96.75% and 94.95% on cell size [8,8] in weka.

## D. Experiment2

In our 2nd experiment we use some statistical measurement for extracting effective features. Table IV shows the detail of experiment 2. Lazy IBK gives the highest accuracy 96.87% on combination of HOG, LBP, Arithmetic mean and Variance. Random Forest gives the highest accuracy at HOG, LBP and Arithmetic mean.

## E. Experiment3

In 3rd experiment we use only HOG features with LBP and increase the HOG cell size from [2,2] or [4,4] to [8,8] with Numbins 12. We got much better accuracy using Lazy.IBK up to 97.37%. Whereas Random Forest did not achieve much better accuracy. At NumBins 15, Lazy.IBK achieve better accuracy up to 97.70% whereas Random forest achieve accuracy 95.70%. Table V shows the detail of experiment.

## F. Experiment4

The final experiment in weka we use NumBins 25 with cell size [8,8] and Lazy.IBK achieved 98.24% accuracy. Random Forest did not achieve much better accuracy than previously achieved 95.70%. Table VI shows the detail of experiment 4 and different graphs represents the data set accuracy of confusion matrix generated by Lazy.IBK. The below Fig. 9

TABLE V.     EXPERIMENT3 RESULTS

| HOG Features | | | LBP Features |
|---|---|---|---|
| Cell Size | Block Size | Num Bins | No. of Features |
| [8 8] | [2 2] | 12 | [28 28] |
| [8 8] | [2 2] | 15 | [14 14] |
| [8 8] | [2 2] | 15 | No |
| Extracted Features | Results | | |
| No. of Features | Lazy.IBK | Random Forest | |
| 348 | 97.37% | 95.70% | |
| 477 | 97.21% | 95.94% | |
| 241 | 97.70% | 95.70% | |

TABLE VI.     EXPERIMENT4 RESULTS

| HOG Features | | | LBP Features |
|---|---|---|---|
| Cell Size | Block Size | Num Bins | No. of Features |
| [8 8] | [2 2] | 25 | No |
| [8 8] | [2 2] | 23 | No |
| [8 8] | [2 2] | 27 | No |
| [8 8] | [2 2] | 29 | No |
| Extracted Features | Results | | |
| No. of Features | Lazy.IBK | Random Forest | |
| 401 | 98.24% | 95.64% | |
| 369 | 98.14% | 95.62% | |
| 433 | 98.14% | 95.58% | |
| 465 | 98.22% | 95.12% | |



Fig. 8.   Training and Testing Data Graph

shows the comparisons of both algorithms on different number of features. From the graph it shows that Lazy.IBK gives much better accuracy than Random Forest algorithm.

## G. Experiment5

After lot of experiments, we use a different tool "Rapid Miner 9.1.0". Rapid Miner is very sophisticated tool used for data mining. Algorithms (K-NN, Neural Net, Generalized Linear Model, Deep Learning, Naive Bayes, Random Forest and Decision Tree) are used to train a model and test our data set. Table VII shows the test results. Unfortunately, we could not achieve much better results than Weka. We got highest result 98.03 using K-NN.

Fig. 10, 11 shows the graph of accuracy and time taken of Rapid Miner algorithms. KNN shows top accuracy and time



Fig. 9.   LazyIBK and Random Forest Result Comparison

TABLE VII.     RESULT COMPARISON OF DIFFERENT CLASSIFIERS ON RAPID MINER

| Algorithm | Accuracy (%) | Time |
|---|---|---|
| KNN | 98.03 | 0:07:41 |
| NN | 97.56 | 0:41:50 |
| GLM | 96.83 | 1:07:43 |
| GL | 96.72 | 5:45:13 |
| NB | 74.57 | 0:03:56 |
| RD | 40.63 | 6:19:52 |
| DT | 11.34 | 4:50:03 |

Fig. 10.    Algorithms accuracy comparison.



Fig. 11.    Algorithms time comparison.

in the graph.

### H. Experiment6

After lot of experiments 8 in Rapid Miner, we decided to use auto model facility in Rapid Miner to build a model for training data. After building model, we achieved impressive results. Using HOG features with Cell [8 8] and NumBin 25, we got the success to build a model with 100% accuracy result using 'Generalized Linear Model' in auto model, in Rapid Miner. Using 'Deep Learning' we achieve 99.9% results whereas "Naïve Bayes" achieved 89.5% our results. Fig. 12 shows the detail of our achieved results on given test data set.

## V.    Conclusion and Result Analysis

From the above experiments, numbers of well-known algorithms are test on training and testing data provided by kaggle. Their results are clearly giving a message that Rapid Miner using auto model gives 100% accuracy. Whereas building up model using Lazy.IBK in Weka 3-9-3 gives 98.24% accuracy. Naive Bayes and Decision Tree did not achieve much better results and we did not add their results in this paper. On the other hand in Rapid Miner we use the following algorithms: KNN (K-Nearest Neighbour), Neural Net, Generalized Leaner



Fig. 12.    Results on Rapid Miner using Auto Model Algorithm.

Model, Deep Learning, Naïve Bayes, Random Forest and Decision Trees. Auto model is also used in Rapid Miner with following algorithms: "Nave Bayes", "Generalized Leaner Model" and "Deep Learning" in Rapid miner. In Rapid Miner 9.1.0, using different algorithms, we achieved highest accuracy 98.03% from K-NN (K Nearest Neighbor) algorithm.

In Rapid Miner using auto model, "Generalized Linear Model" produced 100% results whereas "Deep Learning" also produced 99.9%results and "Naive Bayes" achieved 89.5% results.

Rapid Miner performs extra ordinary performance on test data. Rapid Miner achieved 100% results as compare to Weka tool.

### References

[1] M. J. Cheok, Z. Omar, and M. H. Jaward, "A review of hand gesture and sign language recognition techniques," Int. J. Mach. Learn. Cybern., vol. 10, no. 1, pp. 131–153, 2019.

[2] A. Akoum and N. Al Mawla, "Hand Gesture Recognition Approach for ASL Language Using Hand Extraction Algorithm," J. Softw. Eng. Appl., vol. 08, no. 08, pp. 419–430, 2015.

[3] J. (Junta D. de la S. E. de N. (SEN)) Sancho Rieger et al., "Vision Based Gesture Recognition Using Neural Networks Approaches: A Review," Int. J. Hum. Comput. Interact., vol. 32, no. 1, pp. 480–94, 2008.

[4] H. Zhong, J. Shi, and M. Visontai, "Detecting Unusual Events in Video," IEEE Conf. Comput. Vis. Pattern Recognit., vol. 2, no. June, p. II–819, 2004.

[5] R. Yang and S. Sarkar, "Gesture recognition using hidden Markov models from fragmented observations," Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit., vol. 1, pp. 766–773, 2006.

[6]   T. Starner, J. Weaver, and A. Pentland, "Real-time American sign language recognition using desk and wearable computer based video," IEEE Trans. Pattern Anal. Mach. Intell., vol. 20, no. 12, pp. 1371–1375, 1998.

[7]   K. Grobel and M. Assan, "Isolated sign language recognition using hidden Markov models," pp. 162–167, 2002.

[8]   C. Vogler and D. Metaxas, "ASL recognition based on a coupling between HMMs and 3D motion analysis," pp. 363–369, 2002.

[9]   B. K. Ko and H. S. Yang, "Finger mouse and gesture recognition system as a new human computer interface," Comput. Graph., vol. 21, no. 5, pp. 555–561, Sep. 1997.

[10]  T. Khan and H. Pathan, "Hand Gesture Recognition based on Digital Image Processing using MATLAB," Int. J. Sci. Eng. Res., vol. 6, no. 9, pp. 338–346, 2015.

[11]  K. M. Lim, A. W. C. Tan, and S. C. Tan, "A feature covariance matrix with serial particle filter for isolated sign language recognition," Expert Syst. Appl., vol. 54, pp. 208–218, 2016.

[12]  M. Maraqa and R. Abu-Zaiter, "Recognition of Arabic Sign Language (ArSL) using recurrent neural networks," 1st Int. Conf. Appl. Digit. Inf. Web Technol. ICADIWT 2008, pp. 478–481, 2008.

[13]  V. Adithya, P. R. Vinod, and U. Gopalakrishnan, "Artificial neural network based method for Indian sign language recognition," 2013 IEEE Conf. Inf. Commun. Technol. ICT 2013, no. Ict, pp. 1080–1085, 2013.

[14]  N. M, "Alphabet Recognition of American Sign Language : A Hand Gesture Recognition Approach Using Sift Algorithm," Int. J. Artif. Intell. Appl., vol. 4, no. 1, pp. 105–115, 2013.

[15]  M. HafizurRahman and J. Afrin, "Hand Gesture Recognition using Multiclass Support Vector Machine," Int. J. Comput. Appl., vol. 74, no. 1, pp. 39–43, 2013.

[16]  S. Suharjito, M. C. Ariesta, F. Wiryana, and G. P. Kusuma, "A Survey of Hand Gesture Recognition Methods in Sign Language Recognition," Pertanika J. Sicence Technol., vol. 26, no. 4, pp. 1659–1675, 2018.

[17]  C. F. F. Costa Filho, R. S. De Souza, J. R. Dos Santos, B. L. Dos Santos, and M. G. Fernandes Costa, "A fully automatic method for recognizing hand configurations of Brazilian sign language," Rev. Bras. Eng. Biomed., vol. 33, no. 1, pp. 78–89, 2017.

[18]  M. Mohandes, S. Aliyu, and M. Deriche, "Arabic sign language recognition using the leap motion controller," IEEE Int. Symp. Ind. Electron., no. June, pp. 960–965, 2014.

[19]  H. Bhavsar and J. Trivedi, "Review on Feature Extraction methods of Image based Sign Language Recognition system," Indian J. Comput. Sci. Eng., vol. 8, no. 3, pp. 249–259, 2017.

[20]  L. Quesada, G. López, and L. A. Guerrero, "Sign Language Recognition Using Leap Motion," pp. 277–288, 2015.

[21]  "Sign Language MNIST — Kaggle." [Online]. Available: https://www.kaggle.com/datamunge/sign-language-mnist. [Accessed: 26-Mar-2019].

[22]  Hofmann, Markus and K. Ralf,RapidMiner: Data mining use cases and business analytics applications, CRC Press, 2013

# A Comprehensive Survey on the Performance Analysis of Underwater Wireless Sensor Networks (UWSN) Routing Protocols

Tariq Mahmood[1], Faheem Akhtar[2], Khali ur Rehman[4],
Saqib Ali[5], Fawaz Mahiuob Mokbal[6]
Faculty of Information Technology,
Beijing University of Technology,
Beijing, China 100124

Sher Daudpota[3]
Department of Computer Science
Sukkur IBA University
Sukkur 65200, Pakistan

*Abstract*—**The probe of innovative technologies is a furious issue of the day for the improvement of underwater wireless sensor network devices. The undersea is a remarkable and mystical region which is still unexplored and inaccessible on earth. Interest has been increasing in monitoring the medium of underwater for oceanographic data collection, surveillance application, offshore exploration, disaster prevention, commercial, scientific investigation, attack avoidance, and other military purposes. In underwater milieus, the sensor networks face a dangerous situation due to intrinsic water nature. However, significant challenges in this concern are high power consumption of acoustic modem, high propagation latency in data transmission, and dynamic topology of nodes due to wave movements. Routing protocols working in UWSN has low stability period due to increased data flooding which causes nodes to expire quickly due to unnecessary data forwarding and high energy consumption. The quick energy consumption of nodes originates large coverage holes in the core network. To keep sensor nodes functional in an underwater network, dedicated protocols are needed for routing that maintain the path connectivity. The path connectivity consumes more energy, high route updated cost with a high end to end delay for the retransmission of packets. So, in this paper, we are providing a comprehensive survey of different routing protocols employed in UWSN. The UWSN routing protocols are studied and evaluated related to the network environment and quality measures such as the end to end delay, dynamic network topology, energy consumption and packet delivery ratio. The merits and demerits of each routing protocol are also highlighted.**

*Keywords*—*Underwater Wireless Sensor Networks (UWSN); routing protocols; end-to-end delay; energy consumptions*

## I. INTRODUCTION

Wireless sensor networks (WSNs) is an emerging technology of the day which is used to structure large number of separated tiny embedded sensor nodes used in monitoring and sensing of data from the aqueous environment [1], [2]. In present years, wireless sensor network applications span in the different fields used in weather monitoring, pollution monitoring, military, health, home, and commercial companies [3].
Wireless sensor networks comprises of structured and unstructured networks for sensing of huge amount of aqueous data [4]. The deployment of structure or unstructured networks is usually depend on the environments that needs to be moni-

tored. Mobile Social Networks is a modern distributed buffer storage approach used for data exchange and communication between the mobile users to enhance the network performance concerned to content delivery ratio, e2e delay and throughput [5].

In underwater wireless sensor network; routing is different from the terrestrial wireless sensor networks due to limited bandwidth, energy, node mobility, and end to end delay in the data packet transmission [6]. The energy efficiency is an important prerequisite to their reliable operation and resource management. The routing protocol techniques play an crucial role in energy efficiency that supports network quality of service [7]. Since every application has different quality factors and challenges, there is need to diversify routing protocol having ability to fulfill the application requirements [8].

Fig. 1 illustrate about the architecture of underwater wireless sensor network [9], [10]. In underwater many static acoustic sensor node are distributed over seabed and mobile sensor nodes that move freely with water current. Sensor nodes nous data from environments and detect the movement of submarine that works as autonomous underwater vehicle. In this figure, submarine acquire data from acoustic sensor nodes, aggregate the collected data and store data temporarily. The submarine, static node and mobile sensor node communicate each other and their relevant cluster head through acoustic signals. The submarine forward the aggregated data to surface sink base station. Surface Sink Base station froward the received data to man-controlled computer using radio signal through satellite communication.

TDOA localization (Time difference of arrival technology) is a significant technology instigated in sensor node to ascertain the source location in the real time or time critical applications to evade from interruption arise due to multipath channel [11]. The autonomous underwater vehicle (UUVs, AUVs, Submarine) move in water in fixed pattern to communicate with sensor nodes through short range and high rate data link. AUVs near the base station negotiate to the sink node to forward the collected data. The efficient underwater communication poses significant problems due to intrinsic absorption nature of water. Currently, there has been a growing awareness in monitoring underwater media, analysis of water quality, water pollution monitoring [12], [13] such as biological, chemical

Fig. 1. Underwater Sensor Network Architecture.

and nuclear, micro-organism or fish movement tracking, pressure movements, temperature movement, disaster prevention, underwater oil and gas pipeline corrosion detection, military and home land security application [14], [15]. The acoustic signal also faces many problems due to high error rate, low available bandwidth, node mobility, less propagation speed, a wide class of security threats and malicious attacks and high end to end delay [16], [17]. In recent years the wireless sensor network applications span in the different fields that are used in weather monitoring, water pollution monitoring such as biological, chemical, nuclear, micro-organism tracking, and disaster prevention [18]. In UWSN, link failure is burning issue due to node mobility. Node mobility will create the holes in the network causing increase in end to end delay in data transimission. Existing routing protocols have been analyzed which investigate the optimization performance of network services, node mobility, end to end delay and energy draining of sensors nodes. Large number of routing protocol have been developed working in UWSN [19], [20]. These protocols evaluate the performance efficiency in respect to end to end delay, node mobility, network throughput, and energy consumption. In order to design the efficient UWSN routing protocol researcher faces many challenges which attenuated to the medium are node mobility, end-to-end propagation delay and energy saving [21].

This rest of the paper is organized as follows, Section 2 presents overview and working of UWSN routing protocol. Section 3 presents the comparison of different routing protocol working in underwater wireless sensor network. Section 4 presents the detail of comparative results. Finally, the paper is concluded and future research directions and issues are pointed out in Section 5.

### A. Major Challenges in Design of UWSN

Factors affecting on the propagation of underwater acoustic signals have become the designing challenges for UWSN. Following are the factors that effect on propagation of underwater acoustic signals: [22]

- Bandwidth: The bandwidth available is extremely limited due to water absorption.

- Propagation delay: In underwater network, the propagation delay is five times higher than RF (terrestrial channel). The RF speed is 3x108 ms-1 whereas speed of acoustic signal is about 1.5x103 ms-1. The low speed of sound causes multi-path propagation to stretch over time delay. It effects real time application of UWSN.

- Shadow zones: Due to the underwater extreme characteristics like density and temperature, high bit error rates and temporary losses of connectivity occur.

- Energy: Limited battery power is difficult to recharge.

- Attenuation: Due to decreased amplitude and intensity of a signal.

- The devices for underwater sensor networks are more expensive and have limited availability in market.

- Noise from machinery, shipping and movement of the fish or animals are concerned in UASN.

### B. Differences between UWSN and TWSN

The underwater wireless sensor network is different from the terrestrial wireless sensor network due to the unique characteristics of the water. Following are the major differences shown UWSN and TWSN [23], [24]

- Communication Method: WSN uses radio signals whereas UWSN uses acoustic signals.

- Cost: Terrestrial sensors are inexpensive while underwater sensors are expensive due to transceivers complexity and protection.

- Power: UWSN needs more power than the TWSN because acoustic signals cover long distance and more complex signal techniques are implemented whereas RF needs less power.

- Memory: UWSN requires more memory for caching the data because the connection of acoustic signal can be disabled by shadow zones (less than 100 meters area) while this issue is not treated in terrestrial WSN.

- Difference in Deployment: The sensors are deploying densely in terrestrial sensor application like in tracking system while it will be costly in UWSN to deploy densely. It is not easy to deploy sensors in UWSN in densely.

- Performance: The performance of TWSN is better than UWSN.

- Mobility: Sensor nodes in underwater are mobile whereas in terrestrial network they are fixed

### C. Routing Protocols

The routing is the basic task of network layer used to determine the route from source to destination. The network layer is the administrator that tells how the messages are routed within the networks. In Underwater WSN; routing is different from the terrestrial WSN due limited amount of bandwidth, node mobility for ocean current and end to end delay in data packet transmission [25]. Therefore, in order to hold the network together, there is a need to develop the routing strategies. The design of routing protocol for UWSN is

concerned with saving energy and node mobility in the long-term non-time critical applications. The researchers have made numerous efforts to develop efficient routing protocol while considering the unique characteristics of underwater network [26].

There are mainly three categories namely that determine the path:

*1) Proactive Routing Protocols (Table Driven):* The core function of this protocol is to maintain the routing table containing all routing information to find routes from node to node [25], [27], [28]. This protocol reduces message latency brought by routing discovery. The proactive routing protocol first generates a signal on predefined route to establish the route. All nodes update route information in their routing table. The protocol establishes the route because every time topology is modified due to link failure and node failure in underwater wireless sensor networks. In UWSN, memory and energy are main reasons to avoid the proactive routing protocol in UWSN.

*2) Reactive Routing Protocols:* In reactive routing protocol, node start the route finding process when a route is needed to destination. Once route is established, it is maintained by routing table and is remained in routing table until it is needed. This protocol is more suitable for dynamic environments. This protocol is usually used by source initiated by flooding method [25], [27], [28] . This results an increase in message latency unsuitable for UWSNs.

*3) Geographic Routing Protocols:* Source to destination path is established in geographic protocol by controlling location information[25], [27], [28], [29]. In this scenario, source node selects next forwarder node that is based on the location information of neighbour node. In the underwater environments, it is challenging to attain an accurate location information due to the node movement in water current.

This paper is organized as follows. Section 2 presents overview and working for different routing protocols working in Underwater Wireless sensor network. Section 3 presents evaluations and analysis of different routing protocols based on information of different metric regarding network performance, architecture and design. Section 4 presents the conclusion and point out the future open research issues.

## II. BACKGROUND

This section describes the functioning of different routing protocol operational in Underwater Wireless Sensor Network. Metric like node mobility, end to end delay, deployment, routing approach, energy consumption and packet delivery ratio are used as depiction of routing protocol. According to our review and explored search material we come to know three major categories of these protocols namely Geography based, flat based/Multipath based and clustering/hierarchical based routing protocols.

### A. Geography based Routing Protocol

Location information is required to determine the distance between two selective nodes so that end to end delay and energy consumption can be calculated approximately. This category is subdivided into following specific routing protocols; these are VBF, HH-VBF and AURP. Let us discuss the each of them in the following subsections.

TABLE I. VBF DATA PACKET [8]

| Features subset | Source Position | Target Position | Forwarder Position | Range Field | Radius Field |
|---|---|---|---|---|---|



Fig. 2. VBF Routing Protocol with a Single Virtual Path.

*1) Vector-Based Forwarding Routing Protocol:* VBF is geographically based (location-based) routing protocol, proposed for Underwater WSN [30], [31]. The VBF routing protocol handle node mobility and maintains the routing path at medium level. The VBF gives high data delivery rates, energy efficiency and robust. There is no node state information is required because all nodes are involve in packets forwarding. Data packets is forwarded with redundant path from sender to destination. Sink deals with node failure and packet loss problem.

**Routing**: Table I illustrates about the all field of data packet routed in VBF routing protocol. Each data packets carries routing information of source position, target position, forwarder position. The Range field is used to handle the node mobility and Radius field which contain pre-defined threshold width that is used to decide as a forwarder by sensor nodes if they are nearby the routing pipe.

Fig. 2 explains, how VBF routing protocol build a virtual path in network and clearly show a virtual path build for nodes A, B, and C. In VBF, a "routing pipe or virtual pipe or vector" is established between sender node to target node that embeds its own position, sink position, and its position as a rely node in the packet and broadcasts this packet.. All packets are forwarded through this vector pipe from source to sink node. Only node that close to the vector pipe have ability to forward the data packets from source node to destinated sink. Routing pipe have not only limiting the network traffic significantly but also can easily control the dynamic topology. Every node keeps the information of its location. All the nodes estimate their position by determining the forwarder distance and Angle of Arrival (AoA) of receiving packet. If a node is near to vector pipe with respect to pre-defined threshold distance value, node update its own position in FP field of data packet and forward otherwise discard. VBF is working in two modes: Sink initiated Query:It is location depended query. Sink initiate query for interested area by issuing the INTEREST query packets that contains the source location, target location in the sink coordination system. Source initiated Query:It is location independent query. When

source is interested to send data after setting its location in the source coordination system, it broadcast $DATA_READY$ packets to all nodes.

Desirable Factor:The node competence to forward packets to another nodes are calculated by desirable factor.

**Merits**:

- Achieve robustness against the node failure.
- In dense area, end to end delay is minimum.
- Energy efficient..

**Demerits**:

- In VBF node mobility is not handled efficiently.
- Small data delivery in sparse area because some time, few node lie in routing pipe that are responsible for packet forwarding.
- VBF is very sensitive for pre-define threshold radius because routing performance can affect by threshold radius. Some nodes send data packets again and again from source to sink which drain their energy and increase end to end delay.
- Multiple node are involved in routing act as relay node.
- Due to 3-handshake nature, VBF produce communication overhead which cannot consider the link quality during this operation.



Fig. 3. HH-VBF, Each Node Separate Virtual Path.

*2) HH-VBF: Hop by - Hop Vector - Based Forwarding Routing Protocol:* Hop-by-Hop Vector Based Forwarding Routing Protocol (HH-VBF) is working similar to VBF routing protocol.

Fig. 3 shows each forwarder node builds its separate virtual pipe instead of one virtual pipe between source and destination node in HH-VBR [32], [33], while in VBF, a single "vector" is built between sender to destination in the entire network.

**Routing**: In HH-VBF, intermediate node make forwarding decision based on its present location. Forwarding process is similar to VBF, when a node receive a packet, it keeps packet for short time. This awaiting time is proportional to desirable factor that decide which node is suitable as forwarder node. It is calculated by measuring the distance from the routing vector, angle and transmission range between the nodes. On the expiry of awaiting time, node forward that packets which has smallest desirableness factor.



Figure 1. Underwater Acoustic Sensor Network with multiple AUVs.

Fig. 4. AUVs working in Underwater Acoustic Network.

Like the VBF, in HH-VBF the desirableness factor of each node depends upon:

- Distance of source node to forwarder node.
- Distance of source node to the routing pipe.
- Angle form at forwarder between vector, from forwarder to sink and forwarder to source node.

After recieving a packet, node hold the packet for the time interval $T_{adaptation}$.

**Merits**:

- It gives the better delivery ratio in sparse area.
- Handle node mobility efficiently.
- HH-VBF produce the signaling overhearing due to hop by hop nature.
- Ability to identify delivery route information even nodes are scattered in networks.

**Demerits**:

- By increasing node density will result to increase end to end delay.
- Node density will result to increase energy consumption.
- Node density will result to low data delivery ratio.
- Handle node mobility at medium level due to hop by hop routing pipe.
- Low energy efficiency.
- Difficult to recharge the battery.

*B. AURP: AURP Routing Protocol for UWSN*

AURP Routing Protocol propose to attain the high data rates and low energy consumption in underwater wireless sensor networks [34]. AURP routing protocol use heterogeneous acoustic channels for communication and handle the movement of numbers of AUVs (autonomous underwater vehicles). AUV are used as a relay nodes.

**Routing**: Fig. 4 explains the network architecture of AURP that comprises of U-Sensor nodes, Gateways, AUVs, and surface node/mother ship. The U-Sensor nodes send the aggregated sense data directly to Gateway by using mid range acoustic channel or by multihop fashion. Gateway forward aggregated data to sink directly or through AUVs when it passes near the Gateway through high data rates channel.

Fig. 5. Selection of Optimal forwarded Node in DBR Protocol.

The AUVs deliver received data to sink using short range high data rates channel. Sink sent aggregated data to surface station through fiber optic cables. The AUVs also use low data rates interface for long distance to sent urgent data. AUVs movement is controlled by surface station by sending control signal. The Gateway and sink broadcast their interest toward the sensor node to receive data periodically. AUVs routing protocol is working in four modes:

*Sink Node*: A sink node broadcast messages every time which is used by U-Sensor to find the next node. Sink node negotiate with AUVs to establish the link for data transmission. Sink received sense data form the U-Sensor directly or other U-Sensor relaying data.

*Sensor Nodes*:U-Sensor nodes periodically forward collected data to sink. U-Sensor determine the next node by maintaining the timer.

*Gateway Nodes*: Gateway node flood messages periodically that is recieved by U-Sensors nodes. Gateway receive data from U-Sensors and stores in its queue until AUVs negotiate with gatewsay. AUVs:Multiple autonomous underwater vehicle collect data from gateway and send to sink.

**Merits**:

- Achieve high data rates.
- Low energy consumption.
- Control the mobility of multiple AUVs

Location free routing protocol which need node's related other information like pressure, depth, dynamic address in greedy fashion routing instead of fully geographical node information. Following are the flat based routing protocols in UWSN.

### C. Flat Based / Depth Based Routing Protocol

*1) DBR: Depth-Based Routing :* DBR is location based routing protocol use greedy algorithm [35] to forward the packet from source to sink node. Certain protocol required the full dimensional location information of sensor node in underwater wireless sensor network. In DBR routing protocol each node use depth base information. The sensor node equipped with inexpensive depth hardware to calculate the depth pressure locally.

**Routing**: Fig. 5 shows n1, n2, n3 are neighbor sensor nodes of sensor node S and circle line represents node S

transmission range. Node S broadcasts a packet to n1, n2, n3 neighboring nodes in transmission rang. Node n1 and n2 is qualified forwarding node whereas n1 is selected as optimal forwarded due low depth and node n3 is below of S node, so it discard the packet. Fig. 6 show, multiple stationery data sink are deployed at sea surface in underwater. Sensor nodes are randomly deployed at different depths which sense data from environment and sends data packet containing depth information to its optimal neighboring node through multi-hop fashion (greedy algorithms). On receiving data packets, the neighboring node calculates it depth via pressure hardware and compares it depth "dc" with the packet containing depth "dp". If (dc¡dp) node depth is less then packet containing depth, then packets are forwarded to the next node otherwise packets discard. Data packet is deliver by hop to hop manner to surface sink and base station.

TABLE II. DBR DATA PACKET [25]

| Source ID | Packet Sequence number | Depth | Data |
|-----------|------------------------|-------|------|

In Table II, Data packet of DBR contains source ID, Packet seq #, depth information and original data. "Source ID" is source node identifier, "Packet sequence" unique sequence # assign to data packet by the source node, "Depth" is the depth information that is used to update node by node when the data packet is forwarded.

In DBR, there are two factor involve in collisions and



Fig. 6. DBR node deployment and network working architecture.

redundant transmission of packet. First, multiple path forward data from each nodes by flooding approach and second is, every node send the same data packet many time / repeatedly. In order to control the same packet delivery, the DBR uses priority queue to decrease number of forwarded nodes and packet history buffer to handle the packet retransmission. To prevent from collisions, redundant packet transmission and high overhead, the DBR use holding time for each received packet based on packet embedded depth dp and node own depth dc. Different node have different holding time. Each node wait for holding time when it is expired, it forward the data packet.

**Merits**:

- It can easily handle the mobility of nodes and also handle the traffic at multi-sink on ocean surface.
- It achieves high data rates in sparse area.
- No need for full geographical information.

**Demerits**:

- High energy consumption.
- High end to end delay in dense medium.
- Communication problems in sparse area due to greedy algorithms.
- The data packet is forwarded in broadcast fashion so number of duplicate packets are forwarded which decrease the network performance.
- In sparse area, if depth position of two nodes are same, the network performance is reduced because continuously finding the suitable forwarder.
- Due to sparse and dense area, the complexity is also increased which lead not only to consume more energy and packet losses but also in-efficient use of memory.
- Link failure and hole arise due to draining of energy on the top layer that effect on the network performance.



Fig. 7.   H2-DAB Hop-ID Assigning Process.

*2) Hop-by-Hop Dynamic Addressing Based Routing Protocol: H2-DAB:* $H^2$DAB "Hop-by-Hop Dynamic Addressing Based Routing" is the first greedy and dynamic address based routing protocol in UWSN [36]. The H2-DAB does not require the extra hardware and location information like other greedy protocol in UWSN. The basic purpose of H2-DAB is to solve continuous nodes movement problems. The movement problems is solved by allocating dynamic address to sensor nodes so that the float sensor nodes get the new dynamic address according to their position in different depth level.

**Routing**: Fig. 7 elaborates the ID assigning process in H2-DAB network that consists of multiple static sink located on ocean surface which collect data from sensor fixed with bottom in greedy fashion. The ocean is divided into different depth level (8 to 10 layer) and nodes are deployed randomly move horizontally and vertically. At initial "99", is default Hop ID assigned to every floating node, after receiving Hello packet from sink the node will updated their Hop ID.
H2-DAB is comprising on two different phases 1) assign the dynamically address to static surface nodes, mobile nodes and sensor nodes, 2) sent data by using these addresses.
Table III illustrate about the all field of data packet routed in H2-DAB routing protocolDuring first phase, the dynamic address is allocated to each node by hello message. One static Hop ID is allocated to anchor bottom sensor nodes while the

moving node and surface sink are assigned with two type of Hop ID.

TABLE III.         DBR DATA PACKET [25]

| Sender Hop ID | Next Node ID | Packet Seq # | Destination ID | Data |
|---|---|---|---|---|

During the second phase, data packets are send toward the sink node. To select the forwarder, source node send a inquiry request message to each neighbor node within its communication range. The less Hop ID range node is selected as forwarder, so data is forwarded toward the sink through greedy fashion. Due to mobility of nodes the Hop ID is updated after a time of interval.

**Merits**:

- Gives high data delivery ratio.
- Handle the node mobility without updating the routing table.
- Reduce the congestion of nodes that near to the surface sink.
- It works without extra hardware and maintain any information in routing table.

**Demerits**:

- It is difficult to deploy the mobile nodes at different layer as compared to random deployment.
- Nodes near to sink are working frequently causing drain large amount of energy.
- Communication problem in sparse area due to greedy algorithms is not consider.
- Dynamic addressing phase is completed in short period of time which decreases the network performance.
- Link quality is not considered by single hop that will result in high packet loss and reduce the reliability of network.

*3) Energy-efficient depth-based routing protocol for underwater wireless sensor networks:* EEDBR is an energy efficient localization free routing protocol working in Underwater WSN. The EEDBR [37] is sender based routing approach in which sender node opt the forwarder node depending on node residual energy and depth information from its neighbor nodes. It does not consider the link failure which is important parameter in ocean environment. EEDBR, void critical problem in greedy approach can't handle. Energy draining in dense network is very high due the unnecessary data forwarding. The low-depth nodes expire earlier due to the increased load of data forwarding, causing a less number of available neighbors for the remaining nodes.

**Merits**:

- EEDBR handle the node mobility with water current.
- EEDBR handle the rapid energy consumption of the node near to sink.
- EEDBR is greedy routing approach that uses only depth information and residual energy. There is no need of full location information about the sensor node for routing.

Fig. 8. MPR Basic Procedure.



Fig. 9. (a) Single Path, (b) Partial Multipath, and (c) Multipath.

**Demerits**:

- It does not consider the link failure which is important parameter in ocean environment.
- EEDBR, void critical problem in greedy approach can't handle.
- Energy draining in dense network is very high due the unnecessary data forwarding.
- The low-depth nodes expire earlier due to the increased load of data forwarding, causing a less number of available neighbors for the remaining nodes.

*4) Towards Delay-Sensitive Routing in Underwater Wireless Sensor Networks:* Delay Sensitive DBR [38] is location free routing protocol working in underwater wireless sensor network that is formulated specially for delay sensitive application. Delay Sensitive DBR is enhanced form of DBR in which routing is carry out depend on hold time and depth information of the node. All ordinary sensor nodes forward the sensed data within their transmission scope. A neighbor node that is located in low depth area from the source node calculates their holding time to receive packet. DSDBR working as a greedy algorithm in which data packet is forwarded from source node toward base station through multi-hop fashion. Each qualified neighbor node calculates the forwarding value for the received data packet that is helpful to compute the holding time.

$$H_t = \frac{(\alpha - (TL)_i q_i/\mu)H\_T_{max}}{v_{AC}(TL\_min)} \qquad (1)$$

Where $H_t$ is holding time of received data packet that calculates by each node during which a node stay data packet in its buffer. TL is received packet transmission loss that measure in dB. q is the speed in m/s of received packet. $\mu$ and $\alpha$ is constant depend on network scope. $H\_T_{max}$ is maximum holding time of received packet express in sec. $v_{AC}$ is acoustic signal speed denoted in m/sec. $T\_min$ is the minimum transmission loss between two sensor node that is express in dB. An optimal value of holding time is used by sensor node to limiting the redundant packet transmission. The low depth nodes will not forward the packet on overhearing the received packet. Therefore the DSDBR result the minimize the end to end delay by using of holding time and weight function. In stability period of network, there is trade of among throughput and end to end delay.

**Merits**:

- Networks try to eradicate the distance transmission by the selection of optimal forwarder node depending on holding time and received packet transmission loss.
- DSDBR exhibit a reduced amount of end to end delay by compromising on low stability period and lesser throughput.

**Demerits**:

- The constant depth threshold causes the selection of the same nodes as data forwarders again and again; resulting the quick energy consumption of these nodes.
- It faces tradeoff between minimize the end to end delay and increase in consumption of total energy.
- It show low network stability period that decrease the network performance.

### D. MPR: Multipath Routing

The MPR [39] Routing Protocol construct a routing path form source to destination node. During construction of routing path, multipath is utilized between the sender and receiver node which contains a series of Multi-subpath. Multipaths is a subpath from source to it two-hop neighbor node via a relay node in the neighbourhood of both source and destination nodes.

**Routing**: Fig. 8 shows the construction of routing path, multipath is utilized between the sender and receiver node. Multipaths is route of two-hop neighbor node via a relay node in the neighbourhood of both source and destination nodes. Nodes anchor at sea bottom forward sense data toward the sink at ocean surface. These data packets divided into time slot by source node based on bandwidth. Two hop transmissions are used to send data packets. The destination node receives many packets from different relay nodes, so the MPR prevents from collision at receiving node. The packet arrival is different with different multisubpath.

Fig. 9 (a). There is a single path for data flow between two nodes. The single path waste time during data transmission because it is not fault tolerant. Fig. 9 (b and c), There are multiple path for data transmission between two nodes which have load balancing advantage. In multiple path load balance is attained that decrease the packets drop ratio. The high robustness is also achieved in multipath. Battery lifetime is

Fig. 10.   DUCS, Node deployment and working Network.



Fig. 11.   DUCS Time Line for Node Selection.

improved in multipath routing protocol.

The MPR complete their operation into three phases:

- In first phase, the sender establish routing path. The sender node is required to keep the information of two hop relay from its nearest neighbor node and send to the next hop.
- In second round, the intermediate node "I" is selected by using these information which is collected from propagation delay.
- In third phase, source node check each node to avoid from collision.

**Merits**:

- It has higher throughput in dense area while in sparse area that has low throughput.
- It uses multiple paths and therefore has more overhead.
- It has low end to end propagation delay due to multiple path available.
- It has high packet delivery ratio.

**Demerits**:

- It deals high energy consumption.
- It uses wire for connection.
- Redundant node create the backup route.
- Redundancy creates the contention among nodes.

*E. ARP: Adaptive Routing Protocol*

Adaptive Routing Protocol [40] is used in underwater wireless sensor networks which perform adaptive routing based on application requirements and messages nature. Due to node movement and sparse deployment, UWSN is divided into layers. To achieve the routing performance requirements, protocol get the message redundancy and resource rearrangements. To control packet adaptively, various type of messages have different requirements. The aim of adaptive protocol is to attain trade-off between delivery ratio, energy consumption and propagation delay. The Adaptive routing protocol also provide different services based on data priority level.

**Routing:** Underwater which float freely in two Dimensional plane controlled by the buoyancy nodes. Sink node is placed in center of water surface. All sensor nodes have knowledge about their position due to the localization algorithms.

Sensor nodes use two types of packet, HELLO packet and data packet. HELLO packet contains information exchange with the neighbor nodes. HELLO packet also contain neighbor discovery information. The data packet is sited in payload. Each sensor nodes perform three types of action, 1) Discovery of neighbor; 2) Calculate priority; 3) Decision for routing.

In neighbor discovery mode, each node broadcast HELLO packet to other nodes periodically. The piggy-back ACKs approach is used by each node to broadcast HELLO packet.In second step, packet priority is calculated by vector information which contains, packet age, emergency level, nodes density and battery status. In last step, routing decision is divided into four level. Each level is corresponding to a routing state.

**Merits**:

- Achieve different set of services for different types of data packets based on priority.
- It achieves high delivery ratio.
- It is reliable and also efficient for its bandwidth and energy.

**Demerits**:

- Performance is not good in UWSN.
- Trade-off among energy consumption and delivery ratio.
- Node Mobility is not handled.
- It is not able to analyze the network performance.

### III. CLUSTERING BASED / HIERARCHICAL ROUTING PROTOCOL

In hierarchical routing Protocols, nodes are arranged into clusters where a node having low energy can be used to sense the data from its surrounding and forward the sensed data to respective cluster head while a node having high energy can be elected as a cluster head to aggregate the data received from the sensor node and throw it to sink [41]. By this not only the reduction of energy consumption,but it also achieves the equalization of traffic load and scalability. In this scheme following routing protocols are used along their overview, merits and demerits.

**Routing**: Fig. 10 explains the sensor nodes are structured into cluster where one sensor node is selected as a cluster head for

TABLE IV.      COMPARISON OF UWSN ROUTING PROTOCOLS

| Protocol | Routing Approach | Deployment | End-2-End Delay | Mobility | Energy Efficiency | Delivery Ratio | Localization needed | Rate |
|---|---|---|---|---|---|---|---|---|
| VBF [8,9,10, 17] | Geography base (Flooding) | Dense | Low | Low | Medium | Low | Yes | two packet per second Packet size: 76 byte |
| HH-VBF [9, 10, 17] | Geography base (Flooding) | Dense & Sparse | High | Medium | Low | Medium | Yes | one packet per 10 second Packet size: 50 byte |
| AURP [16] | Geography base (Flooding) | Dense & Sparse | Medium | Low | Medium | High | Yes | 48 kbps |
| DBR [10, 12] | Flat base/depth (Flooding) | Dense & Sparse | High | Medium | Low | High | Partially | one packet per second Packet size: 50 byte |
| $H^2$-DAB [13, 17] | Flat/Addresss base (Flooding) | Dense & | Medium | Medium | Medium | High | No | one packet per second Packet size: 50 byte |
| MPR [15] | Flat base Multipath | Dense | Medium | Low | Medium | Medium | Yes | 10 kbps |
| ADOPTIVE [18] | Flat/Priority base | Dense & Sparse | Medium | Low | Medium | High | Yes | one packet per second (speed 0-5 ms) |
| DUCS [14] | Clustering (Distributed) | Dense | High | Medium | Medium | Medium | No | 6.6 kbps (speed 1.5 ms) |
| MCCP [14] | Clustering (Distributed) | Dense & Sparse | High | Medium | High | Low | Yes | 2 to 3 m/sec (3-5 km/h) |
| EEDBR [23] | Depth (Flooding) | Dense & Sparse | Medium | High | Low | Medium | Partially | 64 bytes every 15 seconds |
| DSDBR [22] | Depth (Flooding) | Dense & Sparse | Low | Low | Low | High | No | 10 bytes (speed 2-3 knots) |

all other nodes. All cluster member nodes forward data packets to their respective cluster head. The cluster head is used for aggregation process on sensed data.Aggregated data is sent to sink through multihop routing. The cluster head performed aggregation function on data and send aggregated data to sink through multihop routing. The cluster head is responsible for inter-cluster and intra-cluster communication. The cluster head randomly selected in order void the draining of battery.

Fig. 11 explains the DUCS routing protocol operation is divided into two rounds, setup round and steady-state round. In steady state, a network is distributed in cluster and data frame is formed with unique id. In steady state rounds, sensor node send numerous frames to each cluster head according to their schedule that composite of series of data massage.

**Merits**:

- It achieves high data delivery ratio.
- It has lower network overhead and increase the network throughput.
- It is energy efficient.

**Demerits**:

- The node mobility affect on cluster structure and reduce cluster life.
- A cluster head only forward the aggregated data to another cluster head.
- Due to water current, the cluster head move away from each other. So they cannot communicate with each other, if any other cluster is not laid down between them.

*1) MCCP: Minimum Cost Clustering Routing Protocol:* In terrestrial wireless networks, LEACH protocol is proposed in which cluster are made with optimal number of cluster head using distributed approach. The cluster head have constant knowledge about node distribution. Due to water intrinsic nature and non uniformly deployment of nodes in underwater, LEACH is unsuitable in underwater sensor network. The cluster head is made without assuming the constant node distribution. The cluster head movement is also not considered therefore traffic load in different area is un-balanced. Both factors are based on centric approach of cluster head. The MCCP protocol [42] is recommended to handle these problem and increase energy proficiency and prolong network life.

**Routing**:MCCP exploit cluster base methodology in which the cluster is formed by using the three parameters.

1). Total energy needed to send data form cluster member to cluster head; 2) remaining energy of cluster members and its cluster head; 3) position of cluster head and its relative members. MCCA (minimum cost clustering algorithms) use centralized methodology to choice the cluster head. The MCCP is advanced of MCCA, in which distributed approach is used to select the cluster head.

By this approach, all sensor nodes are candidates for the cluster members as well as cluster head. Neighbor sets and non neighbor sets are constructed by each candidate in order to form the cluster head. Particular cluster head calculates its average cost and broadcast to all members associated with cluster head ID within two hop range. Every cluster member nodes compare their costs with receiving cluster head cost. Cluster member become cluster head, if its cost is minimum than the cluster head cost. It INVITE a message to all other nodes in the cluster to become its members otherwise join message is send to specific cluster head.

Finally a TDMA schedule is defined by the nominated cluster and forwarded to members of the cluster.

**Merits**:

- It helps in traffic balancing due to formation of more cluster head.
- Number of Cluster members are depending on the sink location and cluster head. The cluster near to sink has less quantity of cluster member.
- It is capable for load traffic balancing by re-clustering node periodically.

**Demerits**:

- It does not support in multi-hop fashion.
- Re-clustering is completed in many days or months.
- Due to mobility, different nodes leave and enter the cluster can reduce the cluster efficiency.

## IV. Comparison of UWSN Routing Protocols

This section describes a comparison in the form of table for different routing protocol functioning in Underwater Wireless Sensor Network. Metric like node mobility, end to end delay, deployment, routing approach, energy consumption and packet delivery ratio are used for analysis of these routing protocols.

In Table IV, routing methodology carries about networks performance enhancement in UWSNs, however it may also faces a variety of intrinsic challenges challenges such as node high mobility due to water flow, 3D environment, high path loss, high rout update cost, low bandwidth, and high propagation delay. Today, there is no any routing protocol can knob all of these challenges. This paper highlighted the basic issues of acoustic communications, data routing and difference between the TWSNs and UWSNs. A collection of such challenges of UWSNs and comprehensive analysis on routing protocols of Underwater Wireless Sensor Networks are briefly discussed as shown in Table IV.
In accordance with the statistics prerequisite for ascertaining the optimistic progress area toward base station, we distributed routing protocols in three categories. A lot number of well-known routing protocols were studied and their merit and demerit are explained with respect to categories. Additionally, pros and cons of each discussed routing protocol is presented which may critically evaluate functionality of each discussed protocol. Furthermore, these protocols were compared to each other based on their features and their simulation conditions. As per literature studied, we can presume the following assumptions as motivating factors which is essential to publish a survey that mainly focuses on those protocols performing issues & challenges on routing / channel assignment at network layer.

## V. Conclusion and Future Work

Routing is an underlying issue of any network, especially in the underwater wireless sensor network. Routing protocols are used to find out different routes that a packet should track over a topology. The design of efficient routing protocol is the fundamental and critical issue of a network layer. The presented research highlights the node mobility and end to end delay issues of UWSN routing protocol at the network layer. The primary motivation of this research is to investigate the performance of different UWSN routing protocols considering different scenarios which are in high demand these days. The researchers made numerous efforts to develop an efficient routing protocol while recognising the unique characteristics of an underwater network. In future work, it is highly needed to create such a routing strategies which can be used to hold the network together i-e underwater and land network, the saving of energy, end to end delay, and node mobility.

## References

[1] G. Han, J. Jiang, C. Zhang, T. Q. Duong, M. Guizani, and G. K. Karagiannidis, "A survey on mobile anchor node assisted localization in wireless sensor networks." *IEEE Communications Surveys and Tutorials*, vol. 18, no. 3, pp. 2220–2243, 2016.

[2] I. F. Akyildiz, D. Pompili, and T. Melodia, "Underwater acoustic sensor networks: research challenges," *Ad hoc networks*, vol. 3, no. 3, pp. 257–279, 2005.

[3] G. Xu, W. Shen, and X. Wang, "Applications of wireless sensor networks in marine environment monitoring: A survey," *Sensors*, vol. 14, no. 9, pp. 16 932–16 954, 2014.

[4] J. Yick, B. Mukherjee, and D. Ghosal, "Wireless sensor network survey," *Computer networks*, vol. 52, no. 12, pp. 2292–2330, 2008.

[5] R. Akhtar, S. Leng, I. Memon, M. Ali, and L. Zhang, "Architecture of hybrid mobile social networks for efficient content delivery," *Wireless Personal Communications*, vol. 80, no. 1, pp. 85–96, 2015.

[6] U. Prathap, P. D. Shenoy, K. Venugopal, and L. Patnaik, "Wireless sensor networks applications and routing protocols: survey and research challenges," in *Cloud and Services Computing (ISCOS), 2012 International Symposium on*. IEEE, 2012, pp. 49–56.

[7] K. Song, B. Ji, and C. Li, "Resource allocation for relay-aided underwater acoustic sensor networks with energy harvesting," *Physical Communication*, vol. 33, pp. 241–248, 2019.

[8] A. Yalcuk and S. Postalcioglu, "Evaluation of pool water quality of trout farms by fuzzy logic: monitoring of pool water quality for trout farms," *International Journal of Environmental Science and Technology*, vol. 12, no. 5, pp. 1503–1514, 2015.

[9] N.-S. N. Ismail, L. A. Hussein, and S. H. Ariffin, "Analyzing the performance of acoustic channel in underwater wireless sensor network (uwsn)," in *Mathematical/Analytical Modelling and Computer Simulation (AMS), 2010 Fourth Asia International Conference on*. IEEE, 2010, pp. 550–555.

[10] F. Akhtar, J. Li, M. Azeem, S. Chen, H. Pan, Q. Wang, and J.-J. Yang, "Effective large for gestational age prediction using machine learning techniques with monitoring biochemical indicators," *The Journal of Supercomputing*, online 2019. [Online]. Available: https://doi.org/10.1007/s11227-018-02738-w

[11] I. Memon, D. Jamro, F. A. Mangi, M. Basit, and M. Memon, "Source localization wireless sensor network using time difference of arrivals (tdoa)," *International Journal of Scientific and Engineering Research*, vol. 4, no. 7, pp. 1046–1054, 2013.

[12] A. Khan and L. Jenkins, "Undersea wireless sensor network for ocean pollution prevention," in *Communication Systems Software and Middleware and Workshops, 2008. COMSWARE 2008. 3rd International Conference on*. IEEE, 2008, pp. 2–8.

[13] A. Watt, M. R. Phillips, C.-A. Campbell, I. Wells, and S. Hole, "Wireless sensor networks for monitoring underwater sediment transport," *Science of The Total Environment*, 2019.

[14] H. Saeed, S. Ali, S. Rashid, S. Qaisar, and E. Felemban, "Reliable monitoring of oil and gas pipelines using wireless sensor network (wsn) remong," in *System of Systems Engineering (SOSE), 2014 9th International Conference on*. IEEE, 2014, pp. 230–235.

[15] F. Akhtar, j. Li, Y. Pei, and M. Azeem, "A semi-supervised technique for lga prognosis," *Proceedings of The International Workshop on Future Technology FUTECH 2019*, pp. 36–37, 2018.

[16] A. Davis and H. Chang, "Underwater wireless sensor networks," in *Proceedings of IEEE Oceans*, 2012, pp. 1–5.

[17] G. Yang, L. Dai, G. Si, S. Wang, and S. Wang, "Challenges and security issues in underwater wireless sensor networks," *Procedia Computer Science*, vol. 147, pp. 210–216, 2019.

[18] E. Felemban, F. K. Shaikh, U. M. Qureshi, A. A. Sheikh, and S. B. Qaisar, "Underwater sensor network applications: A comprehensive survey," *International Journal of Distributed Sensor Networks*, vol. 11, no. 11, p. 896832, 2015.

[19] J. Yifeng and S. Lin, "Nir: Uwsn routing protocol based on node neighbor information," in *Future Information Technology and Management Engineering (FITME), 2010 International Conference on*, vol. 2. IEEE, 2010, pp. 219–222.

[20] A. Celik, N. Saeed, B. Shihada, T. Y. Al-Naffouri, and M.-S. Alouini, "End-to-end performance analysis of underwater optical wireless relaying and routing techniques under location uncertainty," *arXiv preprint arXiv:1901.09357*, 2019.

[21] A. Nayyar, V. Puri, and D.-N. Le, "Comprehensive analysis of routing protocols surrounding underwater sensor networks (uwsns)," in *Data Management, Analytics and Innovation*. Springer, 2019, pp. 435–450.

[22] D. B. Kilfoyle and A. B. Baggeroer, "The state of the art in underwater acoustic telemetry," *IEEE Journal of oceanic engineering*, vol. 25, no. 1, pp. 4–27, 2000.

[23] M. T. Kheirabadi and M. M. Mohamad, "Greedy routing in underwater acoustic sensor networks: a survey," *International Journal of Distributed Sensor Networks*, vol. 9, no. 7, p. 701834, 2013.

[24] M. Ayaz, I. Baig, A. Abdullah, and I. Faye, "A survey on routing techniques in underwater wireless sensor networks," *Journal of Network and Computer Applications*, vol. 34, no. 6, pp. 1908–1927, 2011.

[25] A. Wahid and K. Dongkyun, "Analyzing routing protocols for underwater wireless sensor networks," *International Journal of Communication Networks and Information Security (IJCNIS)*, vol. 2, no. 3, 2010.

[26] P. Carroll, S. Zhou, K. Mahmood, H. Zhou, X. Xu, and J.-H. Cui, "On-demand asynchronous localization for underwater sensor networks," in *Oceans, 2012*. IEEE, 2012, pp. 1–4.

[27] G. G. Xie and J. H. Gibson, "A network layer protocol for uans to address propagation delay induced performance limitations," in *OCEANS, 2001. MTS/IEEE Conference and Exhibition*, vol. 4. IEEE, 2001, pp. 2087–2094.

[28] S. Iyer and D. V. Rao, "Genetic algorithm based optimization technique for underwater sensor network positioning and deployment," in *Underwater Technology (UT), 2015 IEEE*. IEEE, 2015, pp. 1–6.

[29] I. Ullah, J. Chen, X. Su, C. Esposito, and C. Choi, "Localization and detection of targets in underwater wireless sensor using distance and angle based algorithms," *IEEE Access*, vol. 7, pp. 45 693–45 704, 2019.

[30] P. Xie, J.-H. Cui, and L. Lao, "Vbf: vector-based forwarding protocol for underwater sensor networks," in *International conference on research in networking*. Springer, 2006, pp. 1216–1221.

[31] Y. Bayrakdar, N. Meratnia, and A. Kantarci, "A comparative view of routing protocols for underwater wireless sensor networks," in *OCEANS, 2011 IEEE-Spain*. IEEE, 2011, pp. 1–5.

[32] H. Bhambri and A. Swaroop, "Underwater sensor network: Architectures, challenges and applications," in *Computing for Sustainable Global Development (INDIACom), 2014 International Conference on*. IEEE, 2014, pp. 915–920.

[33] N. Nicolaou, A. See, P. Xie, J.-H. Cui, and D. Maggiorini, "Improving the robustness of location-based routing for underwater sensor networks," in *Oceans 2007-Europe*. IEEE, 2007, pp. 1–6.

[34] S. Yoon, A. K. Azad, H. Oh, and S. Kim, "Aurp: An auv-aided underwater routing protocol for underwater acoustic sensor networks," *Sensors*, vol. 12, no. 2, pp. 1827–1845, 2012.

[35] H. Yan, Z. J. Shi, and J.-H. Cui, "Dbr: depth-based routing for underwater sensor networks," in *International conference on research in networking*. Springer, 2008, pp. 72–86.

[36] M. Ayaz and A. Abdullah, "Hop-by-hop dynamic addressing based (h2-dab) routing protocol for underwater wireless sensor networks," in *2009 international conference on information and multimedia technology*. IEEE, 2009, pp. 436–441.

[37] A. Wahid, S. Lee, H.-J. Jeong, and D. Kim, "Eedbr: Energy-efficient depth-based routing protocol for underwater wireless sensor networks," in *Advanced Computer Science and Information Technology*. Springer, 2011, pp. 223–234.

[38] N. Javaid, M. R. Jafri, S. Ahmed, M. Jamil, Z. A. Khan, U. Qasim, and S. S. Al-Saleh, "Delay-sensitive routing schemes for underwater acoustic sensor networks," *International Journal of Distributed Sensor Networks*, vol. 11, no. 3, p. 532676, 2015.

[39] Y.-S. Chen, T.-Y. Juang, Y.-W. Lin, and I.-C. Tsai, "A low propagation delay multi-path routing protocol for underwater sensor networks," *Internet Technology Journal*, vol. 11, no. 2, pp. 153–165, 2010.

[40] Z. Guo, G. Colombi, B. Wang, J.-H. Cui, D. Maggiorini, G. P. Rossi *et al.*, "Adaptive routing in underwater delay/disruption tolerant sensor networks," *Wireless on Demand Network Systems and Services*, pp. 31–39, 2008.

[41] M. C. Domingo and R. Prior, "A distributed clustering scheme for underwater wireless sensor networks," in *Personal, Indoor and Mobile Radio Communications, 2007. PIMRC 2007. IEEE 18th International Symposium on*. IEEE, 2007, pp. 1–5.

[42] P. Wang, C. Li, and J. Zheng, "Distributed minimum-cost clustering protocol for underwater sensor networks (uwsns)," in *Communications, 2007. ICC'07. IEEE International Conference on*. IEEE, 2007, pp. 3510–3515.

# A Novel Framework for Drug Synergy Prediction using Differential Evolution based Multinomial Random Forest

Jaspreet Kaur[1], Dilbag Singh[2], Manjit Kaur[3]

Computer Science & Engineering,
Apex Institute of Technology
Chandigarh University
Gharuan, Punjab, India

*Abstract*—An efficient prediction of drug synergy plays a significant role in the medical domain. Examination of different drug-drug interaction can be achieved by considering the drug synergy score. With an rapid increase in cancer disease, it becomes difficult for doctors to predict significant amount of drug synergy. Because each cancer patient's infection level varies. Therefore, less or more amount of drug may harm these patients. Machine learning techniques are extensively used to estimate drug synergy score. However, machine learning based drug synergy prediction approaches suffer from the parameter tuning problem. To overcome this issue, in this paper, an efficient Differential evolution based multinomial random forest (DERF) is designed and implemented. Extensive experiments by considering the existing and the proposed DERF based machine learning models. The comparative analysis of DERF reveals that it outperforms existing techniques in terms of coefficient of determination, root mean squared error and accuracy.

*Keywords—Machine learning; random forest; drug synergy*

## I. INTRODUCTION

With changing lifestyle, more and destructive diseases are occurring because of poor dietary habits, absence of physical exercises, alcohol utilization, etc. Most common harmful diseases nowadays are cancer, obesity, heart disease, stroke and type II diabetes. Utilization of single drug gives just single supplement; however experiencing big illnesses likewise influences entire body [1]. Therefore, combination of different drugs is required for giving appropriate treatment and legitimate supplements to the body. Mix of various drugs is essentially relies on two things i.e., drug oriented or disease oriented. If mixture is drug oriented then it will concentrate on combination of drugs without knowing disease [1]. However, if it is disease oriented, then many different combinations of drugs are possible. In this way, making pair of various drugs and blends isn't easy. Forecast of drug synergy score is an ill caused issue [2]. For this purpose, various machine learning techniques are also implemented, many techniques are also compared in terms of different parameters like accuracy [2]. Prediction of drug synergy is very important to prevent many harmful diseases. It assumes an effective job in medical domain for preventing particular cancer agents. Machine learning methods has a capacity to limit the synergy estimation errors and thus, used ensemble based differential evolution to optimize the SVM regression technique. Developed a Synergistic field-aware factorization machine SyFFM which uses pharmacological data

to inspect and forecast different combinations of drugs [3]. While making different drug combinations, it is very important to contemplate various sources of information like chemical, biological, pharmacological and network knowledge [4]. To minimize the side effects of different combinations of drugs, it is very important to contemplate natural belongings and network knowledge of drugs which helps to find efficient drug combinations [4]. Combining various drugs also leads to drug toxicity and analyzed drug combinations over in vitro normal cell lines. Also, produced various combinations whose effect on normal cell lines is less [5]. Multiple drug combination treatment is much more effective than single drug. Therefore, various Ayurveda complex combinations are explored with their benefits [6]. Also, features made of DNN which are of high-level are more powerful than carefully assembled features for predicting the cell penetrability of fundamentally assorted synthetic mixes in Caco-2 cell lines [7]. Current vaccine adjuvants featuring and the benefit of immune drug synergy to adjuvant and immunization model is made. The attention was on new advances which are studied and applied adjuvant on immune synergies and immunization improvement [8]. Taiji has been exhibited, software with a high-performance for quick and precise estimation of drug synergy dependent on the winning algorithm [9]. Also, synergy prediction work process can carry compound prioritization in huge scale medicate screenings, and synergy stratification work process can choose where the viability of medications definitely known for inciting synergy is higher [10].

A computational framework biological device is built, DrugComboExplorer, to distinguish pathways of driver signaling and anticipate combinations of synergistic drug by incorporating the learning inserted in huge measures of accessible omics and pharmacogenomics information [11]. Drug synergy, of numerous types, can be anticipated with high degrees of exactness with significant clinical potential [12]. This leap understanding of joint systems of action will take into account the plan of balanced combinatorial therapeutics on a vast scale, across various cancer types. A synergy score dependent on the contrast between the drug and the single drug dose response- bends. The CSS-based synergy score can distinguish genuine synergistic and antagonistic drug combinations. An exploratory computational pipeline has been depicted, named as target addition scoring (TAS), that numerically converts the profiles of drug response to target fixation marks, and

in this way gives a ranking of potential therapeutic targets as indicated by their utilitarian significance in a specific cancer sample [13]. A complete survey of the different drug-re purposing techniques concentrating on the computational methodology [14]. Mechanism synergy estimation used well-distinguished information of biology to anticipate interaction of drugs dependent on medication target interactions [15].

**Contributions:** Following are our main contributions in this paper:-

1) Initially, various machine learning techniques are analyzed to estimate drug synergy score.
2) Thereafter, to overcome the issue or parameters tuning issue, in this paper, an efficient Differential evolution based multinomial random forest (DERF) is designed and implemented.
3) Comparative analysis are also drawn between the existing and the proposed machine learning models in terms of coefficient of determination, root mean squared error and accuracy.

The rest of paper is summarized as follows: Section II discusses Related work. Mathematical preliminaries are given in Section III. Section IV provides description of proposed technique. Section V demonstrates Performance analysis. Conclusion remarks are discussed in Section VI.

## II. RELATED WORK

In this section, Related work has been demonstrated which is as follows: For the analysis of drug combinations Sałat et al., 2015, proposed an universal support vector regression (SVR)-based technique that significantly increases the isobolographic investigation [16]. Until drugs are endorsed, many side effects are not perceivable in clinical preliminaries . Therefore, Zhang et al., 2016, formulated drugs that are endorsed, drug and terms of side effect– symptom relationship as a recommend-er framework, and change the issue of estimating side effects into a recommend-er task [17]. Weinstein et al., 2017, Substantial number of drug combinations are given, to organize the best treatments with computational strategies are critical and presented different methodologies to predict synergistic drug interactions [15]. Wang et al., 2017, proposed another computational technique from drug molecular structure for estimating DTIs and protein succession by utilizing the stacked auto-encoder of deep learning which can sufficiently separates the crude information data. Conventional machine learning methods to deal with sensitivity of drug estimation accept that training information and test information must be in a similar feature space and have the equivalent basic distribution [18]. Two methods of exchange learning are presented by Turki et al., 2017, which combine the auxilliary information from the related undertaking with the training information of the objective task to enhance the forecast execution of the target task on the test information [19]. Tsigelny et al., 2018, examined a few machine learning techniques that have been effectively executed in a few instances of combination drug therapy from hypertension, HIV, irresistible sicknesses to cancer [20].
Hemalatha et al., 2018, considered a survey of methodologies that have been presented to handle sensitivity of drug estimation issue particularly as for the customized Cancer treatment

[21]. Yuan et al., 2018, aimed to build up an applicable new SVM model by consolidating the majority of the highlights of the atomic property-based descriptors and fingerprints to improve the precision for the BBB permeability prediction [22]. Olier et al., 2018, examined the learning of quantitative structure activity relationships (QSARs) as a contextual analysis of meta-learning [23]. This area of application is the most noteworthy societal significance, because it is main step to develop new medicines. Levine et al., 2018, examined how slacked linear regression can be utilized to distinguish the physiologic impacts of drugs from information in the electronic health record (EHR) [24]. Tom et al., 2018, concentrated on new advances and applied immune synergies to adjuvant and vaccine improvement. Su et al., 2019, presented Deep-Resp-Forest that has exhibited the promising utilization of deep learning and deep forest approach on the drug reaction estimation tasks [25]. As per Li et al., 2019, Spectroscopy of near-infrared joined with chemometrics was utilized to analyze the fundamental dynamic components including caffeic acid, chlorogenic acid, ursodesoxycholic acid, luteoloside, chenodeoxycholic acid and baicalin in the Tanreqing injection [26]. Mofrad et al., 2019, created and approved a clinically pertinent decision tree in the finding of Alzheimer's disease (AD) for the utilization of cerebrospinal liquid biomarkers [27]. Strategy is proposed by Bashar et al., 2019, to assess the Heart rate (HR) from wearable gadgets using random forest algorithm [28]. Ogunleye et al., 2019, given another information partitioning rule utilizing the mean of the information sections to develop the tree until the child nodes are little in size [29]. Then, connection is made between the local regression and leave nodes to improve the goals of the node outputs. Randomization is presented at tree development and creation of forest. Zhao et al., 2019, created and approved a estimation model by information extracted of eGFR from a territorial health framework [30]. Feng et al., 2019, proposed a strategy for support vector machine dependent on the dragonfly algorithm (DA-SVM) in a offshore oil field to estimate the short-team load of the microgrid [31]. As per Kestenbaum et al., 2019, Regression is a numerical method used to evaluate the relationship between two or more study variables [32]. Norman et al., 2019, Reported flu immunization inclusion in kids with medical comorbidities remains inadequate. Through the investigation and examination of the information of the department of respiration that is combined with the information of the quality of air measurement, meteorological measurement, and time measurement, Jin et al., 2019, considered its related features and sets up a multidimensional estimation features model dependent on a BP neural system [33]. An attempt has been made by Rajalakshmy et al., 2019 to extract a couple of important time area features from sEMG signals. [34]. The examination by Hazelden et al., 2019, identified a few factors which might be helpful to recognize patients as high risk for hospitalization and the following stages will be to decide and consider the role of the drug specialist in preventing hospitalization of these patients [35]. However, machine learning based drug synergy prediction approaches suffer from the parameter tuning problem. Also, these techniques suffer from poor computational speed.

## III. Mathematical Preliminaries

In this section, Multinomial random forest (MRF) has been presented. Also, Multinomial random forest is compared with random forest.

### A. Random Forest for Regression

There are three aspects in which random forest and Multinomial random forest (MRF) is different. In place of bootstrap technique, partitioning procedure of training set has been used. In attribute set choice, randomness has been introduced at each internal node, and for selecting splitting point , an impurity-based multinomial distribution has been utilized. By doing these, we guarantee that each attribute and each conceivable splitting point get an opportunity for the selection. The key is to limit the negative effect of the randomness, which is required for consistency, on prediction execution.

Suppose $\mathcal{Z}_n$ indicates a data set consisting of $n$ *r.r.d.* examinations. Every examination has the form $(P, Q)$, in which $P \in R^Z$ denotes the $Z$-dimensional attributes and $Q \in \{1, \cdots, N\}$ is the correlating mark of the examination.

*1) Training Data Set Partition:* To make a tree, the training set is isolated arbitrarily into two non-overlapping subsets. Different roles have been played by the two subsets. To make the structure of a tree, one subset has been used; the samples in this subset are called as structure points. Once a tree is made, on the premises of the other subset, the marks of its leaves will be re-directed; the samples in this subset are called as estimation points. In this procedure, the structure points are utilized just to change the shape of the tree by deciding the splitting point in each interior node, and for the final prediction, the estimation points have been utilized. For guaranteeing consistency of the tree, the training set's partitioning and the detachment of their roles are essential.

To assemble another tree, the training set is re-divided arbitrarily and freely. The sizes of the two subsets are fixed. The proportion of the two sections is parameterized by partition rate $= |\text{Structure points}|/|\text{Estimation points}|$.

*2) Tree Construction:* When we compare the proposed Multinomial random forest (MRF) with the original Random forest (RF) then in place of the bootstrap technique, partitioning of the training set has been done. In a rational way, some randomness has been introduced while selecting candidate features and splitting point. There are different ways to select the candidate features, but each feature has to be selected with positive likelihood at every split. The case has been discussed with the use of Bernoulli distribution for more appropriate comparison.

The initial change in Multinomial random forest (MRF) is to randomly select candidate features. Draw randomly from a Bernoulli distribution $C(t)$ especially for every interior node. Randomly select it as a feature set if it is 1; else move with the native process for feature set selection (*e.g.* choose randomly $O(Z)$ features, where $O(Z) = \sqrt{Z}$ or $O(Z) = \log Z$).

In the problem of classification , decrease in the impurity is caused by splitting point $v$ is represented by:

$$V(w) = L(\mathcal{Z}^K) - \frac{|\mathcal{Z}^{K_l}|}{|\mathcal{Z}^K|} L(\mathcal{Z}^{K_l}) - \frac{|\mathcal{Z}^{K_r}|}{|\mathcal{Z}^K|} L(\mathcal{Z}^{K_r}), \quad (1)$$

where $\mathcal{Z}^K$ is the structure points and $\mathcal{Z}^{K_l}, \mathcal{Z}^{K_r}$ are two children sets created by $\mathcal{Z}^K$ splitting at $w$, $L(\cdot)$ is the criterion of impurity (*e.g.* Shannon entropy or Gini index).

The other change is that in-spite of the deterministic rule, there is random selection of the splitting point on the basis of a multinomial distribution. In original random forest maximization of $V(w)$ used to be done where splitting point $w$ is intended. But here, Splitting points are selected randomly as per multinomial distribution $Q(\phi)$ made on the basis of the impurity decrease of all possible points. The certain setting up of $Q(\phi)$ is given below:

For all possible splitting points, suppose $V = (V_1, \cdots, V_m)$ be the vector of impurity decline and all candidate features. Firstly, normalize it as $\hat{V} = (\frac{V_1 - \min V}{\max V - \min V}, \cdots, \frac{V_m - \min V}{\max V - \min V})$, and then compute the probabilities $\alpha = (\alpha_1, \cdots, \alpha_m) = \text{softmax}(C\hat{V})$, where $C > 0$ is a hyper parameter.

For selecting the splitting point, to regulate the probabilities the hyper parameter $C$ plays an important role. If $C$ is larger, then there will be more probability to select largest impurity decrease point. The MRF totally becomes random forest when $t \to 0$ and $C = 0$, for the splitting point selection process. The MRF turns Breiman's random forests when $t \to 0$ and $C \to +\infty$.

From the above two processes, To grow a tree, Selection of one feature and its correlating splitting value is done. Structure points influence only the building of the tree while for prediction estimation points are being involved. The process of splitting is to be continued until the given halting criteria are fulfilled.

Just like random forests, MRF's halting condition is also associates to the minimum leaf size. More specifically, in every leaf, the number of estimation points is required to be larger than $g_e$ where $g_e \to \infty$ and $g_e/e \to 0$ as $e \to \infty$.

*3) Prediction:* When a tree has been developed utilizing the structure points as portrayed above, we re-decide the predicted values for the leaves utilizing the estimation points.

Suppose the unlabeled sample is $j$ and $f$ denotes a decision tree in MRF. The probability that the sample $j$ with label $d(d \in \{1, \cdots, N\})$ evaluated by this tree is

$$\beta^{(d)}(j) = \frac{1}{|\mathcal{O}_f^H(j)|} \sum_{(P,Q) \in \mathcal{O}_f^H(j)} Q = d, \quad (2)$$

and the prediction is given by maximizing $\beta^{(d)}(j)$:

$$\hat{b} = f(j) = \arg \max_d \{\beta^{(d)}(j)\}, \quad (3)$$

where $\mathcal{N}_f^H(j)$ represents the number of estimation points in the node containing $j$, $(\cdot)$ is the indicator function.

The last prediction from the MRF depends on the greater part vote:

$$\overline{\hat{b}} = \overline{f^{(Q)}(j)} = \arg \max_d \sum_{i=1}^{Q} f^{(i)}(j) = d, \quad (4)$$

where $Q$ is the number of separate trees in Multinomial random forest(MRF).

EDE uses adaptive $U \in [0.1, 1.0]$ and $CR \in [0, 1]$ On the basis of their last execution, the mutation method and the features of $CR$ are self-organized. Mutation is decays in two types depending upon the proposed technique. Best population is given by first one but the other does not have. In random fashion from these types, selection of two methods by proposed technique is done during evolution. Hence, the ensemble based DE performs better from variants of DE which are already existing.

*B. Standard Differential Evolution*

One of the meta-heuristic technique is Differential evolution (DE), which is a easy, efficient and powerful global optimization method. When compared with the competitive optimization methods in case of convergence speed along with robustness, DE performed better in many real world applications.

To optimize a given problem, DE uses recombination, mutation and selection. In the beginning, population is created randomly. To create new solutions, recombination, mutation and selection operators are used. For the evaluation of an optimistic trial vector, the selection operator is used for next iteration.

DE begins with a population of $e_t$ D-dimensional candidate solutions, that are presented as:

$$\Pi_{k,\delta}(k = 1, 2, \ldots, e_t) = n_{k,\delta}^1, n^2 k, \delta, \ldots, n_{k,\delta}^Z \quad (5)$$

where $k$ denotes the population with the $k^{th}$ solution. $\delta$ presents current generation. $\omega$ denotes the population dimension.

Initially, population focuses to consider every search domain and bounded as:

$$\Pi_{lb} = n_{lb}^1, n_{lb}^2, \ldots, n_{lb}^\omega \quad (6)$$

Also, population in the beginning is constrained to use maximum bound. Maximun bound is presented as:

$$\Pi_{ub} = n_{ub}^1, n_{ub}^2, \ldots, n_{ub}^\omega \quad (7)$$

Thus, initially, population can be rewritten as:

$$n_{k,0} = m_{lb} + r1 * (n_{ub} - n_{lb}) \quad (8)$$

Here, uniformly distributed random variable is represented as r1. where r1=$rand(0,1) \in [0,1]$.

*1) Mutation:* In the second step, Mutation is utilized to evaluate a trail vector as:

$$W_{k,\delta} = w_{k,\delta}^1, w_{k,\delta}^2, \ldots, w_{k,\delta}^\omega \quad (9)$$

Using different mutation methods, the $W_{k,\delta}$ is evaluated:

$Q_c^1$:
$$W_{k,\delta} = \Pi_{\text{best},\delta} + G.(\Pi_{h_1^k,\delta} - \Pi_{h_2^k,\delta}) \quad (10)$$

$Q_c^2$:
$$W_{k,\delta} = \Pi_{\text{best},\delta} + G.(\Pi_{h_1^k,\delta} - \Pi_{h_2^k,\delta}) + G.(\Pi_{h_3^k,\delta} - \Pi_{h_4^k,\delta}) \quad (11)$$

$Q_c^3$:
$$W_{k,\delta} = \Pi_{k,\delta} + G.(\Pi_{\text{best},\delta} - \Pi_{k,\delta}) + G.(\Pi_{h_1^k,\delta} - \Pi_{h_2^k,\delta}) \quad (12)$$

$Q_c^1$:
$$W_{k,\delta} = \Pi_{h_1^k,\delta} + G.(\Pi_{h_2^k,\delta} - \Pi_{h_3^k,\delta}) \quad (13)$$

$Q_c^2$:
$$W_{k,\delta} = \Pi_{h_1^k,\delta} + G.(\Pi_{h_2^k,\delta} - \Pi_{h_3^k,\delta}) + G.(\Pi_{h_4^k,\delta} - \Pi_{h_5^k,\delta}) \quad (14)$$

$Q_c^3$ :
$$\upsilon_{k,\delta} = \Pi_{k,\delta} + N.(\Pi_{h_1^k,\delta} - \Pi_{h_k,\delta}) + G.(\Pi_{h_2^k,\delta} - \Pi_{h_3^k,\delta}) \quad (15)$$

Here, $h_1^k, h_2^k, h_3^k, h_4^k, h_5^k \in [0, 1]$ represent mutually exclusive indexes. To constraint the amplification of DE, $G$ presents mutation scale factor.

*2) Recombination:* To maximize the target vectors diversity, the operator for recombination is used. Implementation of recombination process is given below:

$$W_{k,\delta} = w_{k,\delta}^1, w_{k,\delta}^2, \ldots, w_{k,\delta}^\omega$$
$$\Pi_{k,\delta} = \Pi_{k,\delta}^1, \Pi_{k,\delta}^2, \ldots, \Pi_{k,\delta}^\omega \quad (16)$$

To make a trial vector, Eq. 16 is utilized as:

$$\upsilon_{k,\delta} = \upsilon_{k,\delta}^1, \upsilon_{k,\delta}^2, \ldots, \upsilon_{k,\delta}^\omega \quad (17)$$

Thus, $(\upsilon_{k,\delta}^k)$ is computed as a new trial vector:

$$\upsilon_{k,\delta}^k = \begin{cases} w_{k,\delta}^k & \text{if } \text{rand}_k[0,1) \leq \text{CR or } (k = k_{\text{rand}}) \\ n_{k,\delta}^k & \text{others} \end{cases} \quad (18)$$

Here, the recombination constant is presented as $CR \in [0, 1]$. A randomly selected index is $jrand \in [1, \omega]$ which makes sure that $\upsilon_{k,\delta})$ will be different from $\Pi_{k,\delta}$ by not less than one parameter.

*3) Selection:* If $\upsilon_{k,\delta}^k$ is more than the upper or lower limits, then within the search range, random re-initialization will be done. Then, the fitness values of all trail vectors $(\upsilon_{k,\delta}^k)$ are evaluated as:

$$n_{k,\delta+1} = \begin{cases} \upsilon_{k,\delta} & \text{if } f(\upsilon_{k,\delta}) \leq f(n_{k,\delta}) \\ n_{k,\delta} & \text{otherwise} \end{cases} \quad (19)$$

In case, if better fitness is given by $\upsilon_{k,\delta}$ when compared to $n_{k,\delta}$, then replacement of $\upsilon_{k,\delta}$ will be done with $n_{k,\delta}$ and further proceed for succeeding generation; else continue .

## IV. PROPOSED TECHNIQUE

In this section, Proposed technique has been explained in detail. Algorithm of Proposed technique also been demonstrated .

### A. Ensemble Mutation Operator

Performance of DE based upon the above operators, which are recombination and mutation. Further, mutation operators are separated into two types (i.e, $Q_h^g$ and $Q_c^g$). The $Q_c^g$ is with best solution in in Eqs. (5-7) , and the $Q_h^g$ is without best solution including Eqs. (8-10). For the diversity of the population to be balanced and to balance the convergence speed, two mutation operators are considered by the proposed technique. The first one is obtained from $Q_h^g$, and second one is obtained from $Q_c^g$. Two mutation operators are randomly selected from $Q_h^g$ and $Q_c^g$, respectively.

**Algorithm 1: EDE based synergy prediction** Step 1:

1) First of all, parameters initialization of proposed ensemble based DE is done (i.e, population size $(e_t)$, function evaluations $(\rho)$, two mutation operators obtained from Eqs. (10-12) and Eqs. (13-15), respectively. Also, maximum number of $\rho$ are represented by $Q_\rho$

2) Initialize a counter $\delta = 0$. Randomly initialize a population with size $e_t$ (i.e, $\nu = \theta_{1,\delta}, \ldots, \theta_{e_t,\delta}$) with $\theta_{k,\delta} = n_{k,\delta}^1, \ldots, n_{k,\delta}^\lambda, k = 1, \ldots, e_t$ uniformly distributed within $[\theta_{lb}, \theta_{ub}]$, where $\theta_{lb} = n_{lb}^1, n_{lb}^2, \ldots, n_{lb}^\lambda$ and $\theta_{ub} = n_{ub}^1, n_{ub}^2, \ldots, n_{ub}^\lambda$.

3) Evaluate the fitness of each population and determine the solution with best fitness $(b_s)$.

4) $\rho = \rho + e_t$;

5) While $\rho \leq Q_\rho$
   for $k = 1 : e_t$
   ```
   Compute W_{k_1,δ} by the first
   technique
   Compute W_{k_2,δ} by the second
   technique
   Compute trial vector v_{k_1,δ} by
   Eq.(19)
   Compute trial vector v_{k_2,δ} by
   Eq.(19)
       ρ = ρ+2
   End for
   ```

6) If any variable is outside its limits, then re-initialization of trial vector $v'_{k,\delta}$ (including $v_{k\_1,\delta}$ and $v_{k\_2,\delta}$) within the search space is done randomly.

7) Selection Procedure
   for j=1 to $e_t$
   ```
       Evaluate the trial vector v_{k,δ}
         If  w(v'_{k,δ}) ≤ w(θ_{k,δ})
         θ_{k,δ+1} = v'_{k,δ}, w(θ_{k,δ+1}) = w(v'_{k,δ})
         P_{k,tbest} = v'_{k,δ}, w(P_{k,tbest}) = w(v'_{k,δ})
           If  f(v'_{k,δ}) < w(θ_{best,δ})
           θ_{best,δ} = v'_{k,δ}, w(θ_{best,δ}) = w(v'_{k,δ})
           t_g = v'_{k,δ}, f(t_g) = w(v'_{k,δ})
          End if
         End if
   End for
   ```

8) The generation count $\delta = \delta + 1$ is Incremented;

9) End while

In this paper, Accuracy(ACC) and correlation coefficient(CCO) are utilized to evaluate the fitness of every solution as given below:

$$Maximize(ACC, CCO) \qquad (20)$$

### B. Synergy Prediction using EDE

Various steps are portrayed that are needed to be implement to assess the synergy prediction in a productive way. Different steps are given below:

1) **Selection of framework :** For the completion of the presented procedure, multinomial random forest framework is used. Following parameters will be optimized with the use of proposed technique.

2) **Scaling:** It is used to ignore characteristics which are in highest numeric limit from the minimum numeric limit. It likewise decreases the computational complexity for machine learning methods. scaling of numeric attributes between either $[-1, 1]$ or $[0, 1]$ is done as follows in this paper:

$$w' = \frac{w - l_c^b}{u_c^b - l_c^b} \qquad (21)$$

Here, $w$ represents native value. $w'$ is a scaled value. $u_b$ and $l_b$ represent upper and lower limit of feature values, respectively.

3) **Training and testing data:** The synergistic information is deteriorated into training and testing information. To prepare the proposed random forest based machine learning technique, the training information is used. After that, For the proposed method, testing information is utilized to screen the viability.

4) **Differential evolution based RF:** The differential evolution based random forest is utilized for the tuning of the required attributes of random forests. As a fitness function, the Root mean squared blunder (RMSE) is used. The general target is to discover best parameters for random forest based synergy prediction strategy.

5) **Execution criteria:** At the point when the Execution criteria accomplished, the differential evolution based RF ends itself and then return tuned characteristics; else, for other iterations of DE continues .

6) **Method building:** The end results acquired from DE based random forest are utilized as Random forest technique's initial attributes. Then, it is utilized to make the trained Random forest based synergy estimation method.

7) **Performance analysis:** Here, trained method acquired utilizing DE based random forest. Random forest is put in on the training information to foresee the drug synergy esteems. Subsequently, Use these values for assessing the execution of presented machine learning method.

## V. PERFORMANCE ANALYSIS

This section gives the comparative analysis of proposed and existing machine learning strategies.The data of drug synergy score, which has been used for validation purpose in this research work, comprise of two different terms as Highest concentration of drug P and Highest concentration of drug Q. On distinct drug's distinct concentrations, the benefit of

drug synergy data set relies. Thus, Implemented the proposed model on the train dataset which comprises of different drugs in different concentration for drug synergy prediction.The different drugs in different concentrations has been assessed using Drug interaction coefficient (DIC), which is assessed as follows:

$$DIC = \frac{PQ}{P \times Q} \tag{22}$$

Here, PQ is the extent to the control combination from the two-drug mix combination and P or Q is the extent to the control combination from the single drug combination . The $DIC < 1$ shows synergism, particularly $DIC < 0.7$ demonstrates a fundamentally synergistic impact, $DIC = 1$ demonstrates additivity and DIC > 1 shows opposition.

In this paper, 15-overlap cross-approval is utilized for the testing of ensemble based machine learning exhibit at the time of preparation stage to assess over fitting issue. To accomplish, 15-overlap cross-approval, at first training information has been partitioned into 15 equal subsets (overlap). Approval set is used to keep the 1-overlap and keep other 14-overlaps in the cross approval preparing set. Demonstrated by utilizing the cross-approval set of training and assess the precision of the presented model by favoring the estimated values which opposes the endorsement set. Correspondingly, precision of each of the 15-overlaps have been assessed. To conquer the issue of over-fitting, mean of assessed ac-curacies have been assessed. For approval, each overlap has been utilized just once. In this way, 15-overlay cross-approval ensures, the presented trained model does not experiences the over-fitting problem. The proposed DERF and the existing methods are implemented in Python 3.6. Intel core i5 8th generation processor is used along with 8GB RAM and 2GB graphics card. Here, 20% to 90% ratio of air pollution dataset is taken for training purpose. Also, rest of dataset is used for testing purpose. The acceptance error is allowed only between , for accuracy evaluation.

### A. Experimental Setup

DERF with another methods have been executed with 16 GB RAM on Intel core $i7$ processor. For designing simulation environment, software, MATLAB $2017a$ is used with matplotlib library. In next section, datasets description(i.e., AstraZeneca-Sanger Drug Combination) [36] with its characteristics is given.

### B. Dataset

Various drug combinations are performed to evaluate the drug impacts on cells at various concentrations. The concentration space has been increased in the presence of two drugs: by dimension and the induced effect which has been described by a dose reaction surface that opposes to a curve. A dose reaction surface will generally resemble this: The synergy score information, which has been utilized for approval reason contain 2 terms for example (1) Highest concentration of drug P and (2) Highest concentration of drug Q. Thus, the benefit of synergy data is based on the different drugs combinations. DERF has been actualized on the train data which contains

the distinctive grouping of various medications to analyze the score of synergy . Various drugs with different concentrations has been assessed by utilizing drug interaction coefficient (DIC). Various features of drug synergy dataset are described in Table I.

TABLE I. CHARACTERISTICS OF DRUG SYNERGY DATASET

| Column name | Explanation |
|---|---|
| COMP_P | Drug P's name. |
| COMP_Q | Drug Q's name. |
| COMB_ID | Name of the combination of drug P and drug Q |
| C_L | Normalised cell line name. |
| IC50_P | Concentration in which half of the highest number of elimination is acquired with drug P. |
| H_P | Dose-reaction curve's slope for drug P. |
| Einf_P(Potency) | Highest number of cells killed (percentage) with drug P. |
| IC50_Q | Concentration in which half of the highest number of elimination is acquired with drug Q. |
| H_Q | Dose-reaction curve's slope for drug Q. |
| HIGH_CONC_P | Highest concentration of drug P. |
| HIGH_CONC_Q | Highest concentration of drug Q. |
| QA | Assurance of quality flag of combination assays |
| SYN_SCORE | Evaluated overall synergy of drug P and drug Q in mix. |
| Einf_Q | Highest number of cells killed (percentage) with drug Q. |

### C. Over and Under-Fitting Evaluation

In this section, Depending upon the same fraction of data, various machine learning methods are trained and tested. It is mainly used for the evaluate the over-fitting and under-fitting issue.

Tables II, III, and IV depict the performance analysis of the already existing and the proposed machine learning methods. It is found that DERF outperforms existing methods in terms of accuracy, coefficient of determination, and Root mean squared error, respectively.

TABLE II. TRAINING ACCURACY ANALYSIS

| Dataset | 20 % | 40 % | 60 % | 80 % |
|---|---|---|---|---|
| LR | 88.1 ±1.3 | 91.8 ±1.1 | 87.7 ±0.9 | 88.5 ±1.4 |
| DT | 89.3 ±1.9 | 92.4 ±1.1 | 88.4 ±1.8 | 89.8 ±2.2 |
| RF | 90.2 ±1.3 | 92.7 ±1.1 | 89.0 ±1.4 | 89.8 ±1.6 |
| SVM | 90.2 ±1.7 | 92.7 ±0.9 | 89.0 ±2.4 | 89.8 ±2.1 |
| ANN | 93.7 ±0.8 | 94.6 ±1.1 | 91.4 ±1.2 | 93.6 ±0.9 |
| ANFIS | 94.3 ±0.8 | 95.4 ±0.8 | 91.9 ±1.0 | 94.6 ±1.2 |
| DERF | 98.4 ±0.7 | 99.4 ±0.4 | 97.3 ±0.7 | 99.4 ±0.5 |

TABLE III. TRAINING COEFFICIENT OF DETERMINATION ANALYSIS

| Dataset | 20 % | 40 % | 60 % | 80 % |
|---|---|---|---|---|
| LR | 0.87 ±0.08 | 0.87 ±0.04 | 0.82 ±0.02 | 0.82 ±0.06 |
| DT | 0.88 ±0.01 | 0.87 ±0.02 | 0.83 ±0.06 | 0.83 ±0.04 |
| RF | 0.88 ±0.06 | 0.88 ±0.02 | 0.83 ±0.06 | 0.83 ±0.02 |
| SVM | 0.90 ±0.03 | 0.89 ±0.07 | 0.85 ±0.02 | 0.85 ±0.09 |
| ANN | 0.91 ±0.04 | 0.90 ±0.02 | 0.86 ±0.09 | 0.86 ±0.05 |
| ANFIS | 0.92 ±0.01 | 0.91 ±0.03 | 0.87 ±0.09 | 0.87 ±0.06 |
| DERF | 0.97 ±0.02 | 0.96 ±0.03 | 0.92 ±0.04 | 0.94 ±0.05 |

TABLE IV. ROOT MEAN SQUARED ERROR TRAINING ANALYSIS

| Technique | 20 % | 40 % | 60 % | 80 % |
|---|---|---|---|---|
| LR | 4.6 ±0.54 | 4.6 ±0.76 | 5.1 ±0.44 | 4.5 ±0.89 |
| DT | 5.5 ±0.82 | 4.3 ±0.41 | 4.5 ±0.69 | 4.0 ±0.79 |
| RF | 5.0 ±0.70 | 6.1 ±0.55 | 4.4 ±0.83 | 3.3 ±0.46 |
| SVM | 4.7 ±0.54 | 4.8 ±0.51 | 4.7 ±0.78 | 4.3 ±0.71 |
| ANN | 3.2 ±0.63 | 3.9 ±0.88 | 3.7 ±0.44 | 2.9 ±0.73 |
| ANFIS | 4.4 ±0.77 | 4.1 ±0.81 | 4.1 ±0.62 | 6.8 ±0.58 |
| DERF | 1.7 ±0.31 | 2.0 ±0.39 | 1.1 ±0.27 | 1.2 ±0.29 |

*D. Testing Analysis*

In this section, performance of the existing and the proposed machine learning methods are tested by considering the testing data. It is mainly used for evaluating the effectiveness of DERF over existing methods.

Tables V, VI, and VII depict the performance analysis of the existing and the proposed machine learning methods. It is found that DERF outperforms existing methods in terms of accuracy, coefficient of determination, and Root mean squared error, respectively.

TABLE V. TESTING ACCURACY ANALYSIS

| Dataset | 20 % | 40 % | 60 % | 80 % |
|---|---|---|---|---|
| LR | 91.3 ±1.7 | 89.4 ±2.1 | 90.5 ±1.9 | 89.8 ±1.3 |
| DT | 92.6 ±1.6 | 90.9 ±0.8 | 91.9 ±1.2 | 90.1 ±1.5 |
| RF | 93.1 ±1.8 | 90.6 ±2.0 | 91.3 ±1.3 | 91.7 ±2.2 |
| SVM | 93.1 ±1.0 | 90.6 ±1.7 | 91.3 ±2.1 | 91.7 ±1.9 |
| ANN | 93.5 ±1.1 | 93.9 ±1.3 | 92.6 ±1.6 | 95.9 ±1.0 |
| ANFIS | 94.2 ±1.0 | 94.7 ±1.3 | 93.4 ±2.1 | 96.3 ±1.1 |
| DERF | 98.1 ±0.8 | 98.4 ±0.8 | 98.2 ±0.9 | 99.1 ±0.7 |

TABLE VI. TESTING COEFFICIENT OF DETERMINATION ANALYSIS

| Dataset | 20 % | 40 % | 60 % | 80 % |
|---|---|---|---|---|
| LR | 0.81 ±0.08 | 0.84 ±0.07 | 0.82 ±0.07 | 0.88 ±0.07 |
| DT | 0.81 ±0.07 | 0.84 ±0.06 | 0.82 ±0.07 | 0.89 ±0.06 |
| RF | 0.82 ±0.08 | 0.85 ±0.07 | 0.83 ±0.06 | 0.89 ±0.06 |
| SVM | 0.83 ±0.07 | 0.86 ±0.08 | 0.84 ±0.06 | 0.91 ±0.05 |
| ANN | 0.84 ±0.09 | 0.87 ±0.04 | 0.85 ±0.11 | 0.92 ±0.06 |
| ANFIS | 0.85 ±0.10 | 0.88 ±0.09 | 0.86 ±0.11 | 0.93 ±0.05 |
| DERF | 0.95 ±0.04 | 0.96 ±0.03 | 0.96 ±0.02 | 0.98 ±0.01 |

TABLE VII. ROOT MEAN SQUARED ERROR TESTING ANALYSIS

| Dataset | 20 % | 40 % | 60 % | 80 % |
|---|---|---|---|---|
| LR | 5.1 ±0.63 | 5.4 ±0.78 | 3.7 ±0.47 | 3.0 ±0.68 |
| DT | 3.1 ±0.49 | 3.7 ±0.91 | 4.4 ±0.96 | 3.8 ±1.21 |
| RF | 5.6 ±0.41 | 4.8 ±0.37 | 3.3 ±0.81 | 5.5 ±0.76 |
| SVM | 4.6 ±0.47 | 3.3 ±0.91 | 4.7 ±0.74 | 6.3 ±0.35 |
| ANN | 4.4 ±0.65 | 3.9 ±0.61 | 3.8 ±0.45 | 4.9 ±0.47 |
| ANFIS | 6.3 ±0.39 | 6.5 ±0.48 | 5.4 ±0.51 | 3.5 ±0.55 |
| DERF | 2.5 ±0.24 | 1.9 ±0.27 | 1.6 ±0.19 | 1.9 ±0.28 |

## VI. CONCLUSION

The examination of drug synergy score needs well-organized regression models to decrease the prediction errors. The main objective of this paper is to design a novel differential evolution based multinomial random forest (DERF) approach. The proposed strategy has been assessed on the data set of drug synergy and furthermore contrasted with aggressive machine learning techniques. In experimental consequences, it has been seen that mean enhancement of proposed method over competitive methods in terms of accuracy and coefficient of determination are 2.3598 % and 1.8469 %, respectively. Hence, DERF is effective for planning a estimator of a real-time drug synergy. In this work, we have not considered the use of feature selection techniques. Therefore, to improve speed and accuracy rate, we may utilize some competitive FS techniques.

## REFERENCES

[1] H. Chen and J. Li, "Drugcom: Synergistic discovery of drug combinations using tensor decomposition," in *2018 IEEE International Conference on Data Mining (ICDM)*. IEEE, 2018, pp. 899–904.

[2] H. Singh, P. S. Rana, and U. Singh, "Prediction of drug synergy in cancer using ensemble-based machine learning techniques," *Modern Physics Letters B*, vol. 32, no. 11, p. 1850132, 2018.

[3] C. Zhang and G. Yan, "Synergistic drug combinations prediction by integrating pharmacological data," *Synthetic and systems biotechnology*, vol. 4, no. 1, pp. 67–72, 2019.

[4] P. Ding, R. Yin, J. Luo, and C. K. Kwoh, "Ensemble prediction of synergistic drug combinations incorporating biological, chemical, pharmacological and network knowledge," *IEEE journal of biomedical and health informatics*, 2018.

[5] R. Rahman and R. Pal, "A mathematical framework for analyzing drug combination toxicity for personalized medicine applications," in *2016 IEEE Healthcare Innovation Point-Of-Care Technologies Conference (HI-POCT)*. IEEE, 2016, pp. 13–16.

[6] P. K. Mukherjee, S. Banerjee, and A. Kar, "Exploring synergy in ayurveda and traditional indian systems of medicine," *Synergy*, 2018.

[7] M. Shin, D. Jang, H. Nam, K. H. Lee, and D. Lee, "Predicting the absorption potential of chemical compounds through a deep learning approach," *IEEE/ACM transactions on computational biology and bioinformatics*, vol. 15, no. 2, pp. 432–440, 2018.

[8] J. K. Tom, T. J. Albin, S. Manna, B. A. Moser, R. C. Steinhardt, and A. P. Esser-Kahn, "Applications of immunomodulatory immune synergies to adjuvant discovery and vaccine development," *Trends in biotechnology*, 2018.

[9] H. Li, S. Hu, N. Neamati, and Y. Guan, "Taiji: approaching experimental replicates-level accuracy for drug synergy prediction," *Bioinformatics*, vol. 10, 2018.

[10] M. Yang, M. P. Menden, P. Jaaks, J. Dry, M. Garnett, and J. Saez-Rodriguez, "Stratification and prediction of drug synergy based on target functional similarity," *bioRxiv*, p. 586123, 2019.

[11] L. Huang, D. Brunell, C. Stephan, J. Mancuso, X. Yu, B. He, T. C. Thompson, R. Zinner, J. Kim, P. Davies *et al.*, "Driver network as a biomarker: systematic integration and network modeling of multi-omics data to derive driver signaling pathways for drug combination prediction," *Bioinformatics*, 2019.

[12] C. Gilvary, J. R. Dry, and O. Elemento, "Multi-task learning predicts drug combination synergy in cells and in the clinic," *bioRxiv*, p. 576017, 2019.

[13] A. Jaiswal, B. Yadav, K. Wennerberg, and T. Aittokallio, "Integrated analysis of drug sensitivity and selectivity to predict synergistic drug combinations and target coaddictions in cancer," in *Systems Chemical Biology*. Springer, 2019, pp. 205–217.

[14] K. Savva, M. Zachariou, A. Oulas, G. Minadakis, K. Sokratous, N. Dietis, and G. M. Spyrou, "Computational drug repurposing for neurodegenerative diseases," in *In Silico Drug Design*. Elsevier, 2019, pp. 85–118.

[15] Z. B. Weinstein, A. Bender, and M. Cokol, "Prediction of synergistic drug combinations," *Current Opinion in Systems Biology*, vol. 4, pp. 24–28, 2017.

[16] R. Sałat and K. Sałat, "Modeling analgesic drug interactions using support vector regression: a new approach to isobolographic analysis," *Journal of pharmacological and toxicological methods*, vol. 71, pp. 95–102, 2015.

[17] W. Zhang, H. Zou, L. Luo, Q. Liu, W. Wu, and W. Xiao, "Predicting potential side effects of drugs by recommender methods and ensemble learning," *Neurocomputing*, vol. 173, pp. 979–987, 2016.

[18] L. Wang, Z.-H. You, X. Chen, S.-X. Xia, F. Liu, X. Yan, and Y. Zhou, "Computational methods for the prediction of drug-target interactions from drug fingerprints and protein sequences by stacked auto-encoder deep neural network," in *International Symposium on Bioinformatics Research and Applications*. Springer, 2017, pp. 46–58.

[19] T. Turki, Z. Wei, and J. T. Wang, "Transfer learning approaches to improve drug sensitivity prediction in multiple myeloma patients," *IEEE Access*, vol. 5, pp. 7381–7393, 2017.

[20] I. F. Tsigelny, "Artificial intelligence in drug combination therapy," *Briefings in bioinformatics*, 2018.

[21] R. Hemalatha and T. Devi, "Prognosticate the drugs for multiple myeloma patients by using gene expression technique with polyclonal and monoclonal samples." *ICTACT Journal on Image & Video Processing*, vol. 8, no. 3, 2018.

[22] Y. Yuan, F. Zheng, and C.-G. Zhan, "Improved prediction of blood–brain barrier permeability through machine learning with combined use of molecular property-based descriptors and fingerprints," *The AAPS journal*, vol. 20, no. 3, p. 54, 2018.

[23] I. Olier, N. Sadawi, G. R. Bickerton, J. Vanschoren, C. Grosan, L. Soldatova, and R. D. King, "Meta-qsar: a large-scale application of meta-learning to drug design and discovery," *Machine Learning*, vol. 107, no. 1, pp. 285–311, 2018.

[24] M. E. Levine, D. J. Albers, and G. Hripcsak, "Methodological variations in lagged regression for detecting physiologic drug effects in ehr data," *Journal of biomedical informatics*, vol. 86, pp. 149–159, 2018.

[25] R. Su, X. Liu, L. Wei, and Q. Zou, "Deep-resp-forest: A deep forest model to predict anti-cancer drug response," *Methods*, 2019.

[26] W. Li, X. Yan, J. Pan, S. Liu, D. Xue, and H. Qu, "Rapid analysis of the tanreqing injection by near-infrared spectroscopy combined with least squares support vector machine and gaussian process modeling techniques," *Spectrochimica Acta Part A: Molecular and Biomolecular Spectroscopy*, 2019.

[27] R. B. Mofrad, N. S. Schoonenboom, B. M. Tijms, P. Scheltens, P. J. Visser, W. M. van der Flier, and C. E. Teunissen, "Decision tree supports the interpretation of csf biomarkers in alzheimer's disease," *Alzheimer's & Dementia: Diagnosis, Assessment & Disease Monitoring*, vol. 11, pp. 1–9, 2019.

[28] S. S. Bashar, M. S. Miah, A. Z. Karim, M. A. Al Mahmud, and Z. Hasan, "A machine learning approach for heart rate estimation from ppg signal using random forest regression algorithm," in *2019 International Conference on Electrical, Computer and Communication Engineering (ECCE)*. IEEE, 2019, pp. 1–5.

[29] A. Ogunleye, Q.-G. Wang, and T. Marwala, "Integrated learning via randomized forests and localized regression with application to medical diagnosis," *IEEE Access*, vol. 7, pp. 18 727–18 733, 2019.

[30] J. Zhao, S. Gu, and A. McDermaid, "Predicting outcomes of chronic kidney disease from emr data based on random forest regression," *Mathematical biosciences*, vol. 310, pp. 24–30, 2019.

[31] Y. Feng, P. Zhang, M. Yang, Q. Li, and A. Zhang, "Short term load forecasting of offshore oil field microgrids based on da-svm," *Energy Procedia*, vol. 158, pp. 2448–2455, 2019.

[32] B. Kestenbaum, "Linear regression," in *Epidemiology and Biostatistics*. Springer, 2019, pp. 217–237.

[33] Y. Jin, H. Yu, Y. Zhang, N. Pan, and M. Guizani, "Predictive analysis in outpatients assisted by the internet of medical things," *Future Generation Computer Systems*, 2019.

[34] P. Rajalakshmy, E. Jacob, and T. J. Sharon, "Estimation of elbow joint angle from surface electromyogram signals using anfis," in *Computer Aided Intervention and Diagnostics in Clinical and Medical Images*. Springer, 2019, pp. 247–253.

[35] L. A. Hazelden, M. J. Newman, S. Shuey, J. M. Waldfogel, and V. T. Brown, "Evaluation of the head and neck cancer patient population and the incidence of hospitalization at an academic medical center," *Journal of Oncology Pharmacy Practice*, vol. 25, no. 2, pp. 333–338, 2019.

[36] Y. Guan, "Guanlab's solution to the 2016 astrazeneca-sanger drug combination prediction dream challenge," *Synapse Repository, Synapse ID: syn5614689*, 2016.

# Intruder Attacks on Wireless Sensor Networks: A Soft Decision and Prevention Mechanism

Iftikhar Hussain[1]
School of Computer Science
and Technology
University of Science and
Technology of China
Hefei 230000 China

Abrar Hussain[3]
Hefei National Laboratory
for Physical Science at the Microscale
University of Science and
Technology of China
Hefei 230026 China

Shahzad Haider[5]
School of Electronics Science and Technology
University of Science and
Technology of China
Hefei 230026 China

Samman Zahra[2]
Computer Science Department
COMSATS Islamabad
Islamabad 45550 Pakistan

Hayat Dino Bedru[4]
School of Software
Dalian University of Technology
Dalian 116620 China

Diana Gumzhacheva[6]
School of Management Science
Anhui University
Anhui 230031 China

*Abstract*—Because of the wide-ranging of applications in a variety of fields, such as medicine, environmental studies, robotics, warfare and security, and so forth, the research on wireless sensor networks (WSNs) has attracted much attention recently. WSNs offer economical, flexible, scalable and pragmatic solutions in many situations. Sensor nodes are tiny and have a limited, non-rechargeable battery source, small memory/computational abilities and low transmitter power. Energy resources are vital as once the battery is depleted, the node is no longer usable. Multiple medium access control (MAC) protocols are designed to increase the life cycle of a node by minimizing its unnecessary energy consumption. In some critical applications like the surveillance of enemy movements on a battlefield, opponents deploy adversary nodes to disturb the performance of WSNs by mainly depleting the battery sources of legitimate nodes. In this work, an intrusion detection mechanism has been adapted to detect different kinds of intruders' attacks in MAC protocols of WSN's. A soft decision mechanism has been implemented to detect collision and exhaustion attacks. A preventative mechanism has also been introduced, which helps a node to avoid these intrusive attacks. Results show how the lifetime of a node increases and network performance also increases with better throughput and reduced delay.

*Keywords*—*MAC protocols; S-MAC; wireless sensor networks; intrusion detection*

## I. Introduction

Medium access controls (MAC) perform the key action of developing coordination among nodes and managing the means for their effective communications with the help of an allocated medium. There are various procedures formulations in different situations. Such designed procedures also differ for varying applications and constraints [1]. MAC is one of the major concerns when designing wireless sensor networks. Realizing the concern over limited energy resources encountered in the domain of wireless sensor networks, MAC protocols are exclusively designed in a way to render them energy efficient. One primary cause of excessive energy consumption or wastage is the transfer of data by the two nodes sharing the same medium at the same time [2]. This data transfer leads to the collision of the data packets. In order to accommodate this problem in sensor networks, the MAC procedures are expected to aid the nodes in accessing the medium to avoid the packets' collisions [3]. MAC procedures or protocols are significant features for running any network embodying a shared medium and for attaining efficient performance of the network. In the context of wireless networks, MACs are studied extensively in [4].

Wireless protocols, such as time division multiple access (TDMA), code division multiple access (CDMA), and frequency division multiple access (FDMA), are conventional and are usually employed in conventional wireless networks. Since sensor nodes are powered by battery, these protocols cannot be directly applied to wireless sensor networks [5]. Due to this limitation, such protocols have minimal memory as well as computation power. It is also due to this limitation of wireless protocols that MAC protocols cannot be directly applied to wireless sensor nodes, and this is why they need obligatory alteration [6]. In order to design efficient MAC protocols for sensor networks, network scalability and the computation of energy efficiency are the primary and essential considerations. Features such as latency, bandwidth utility, and throughput are of secondary importance. In a nut-shell [7], priorities and considerations differ with distinct applications of sensor networks.

The remainder of the paper is designed as follows: Section II will give a brief overview of all the terminologies that are used in this research article; Section III is dedicated to the description of some common network attacks on S-MAC, while Sections IV and V will explain the function and simulation results of our proposed detection mechanism.

## II. Problem Description

### A. Sensor Medium Access Control (S-MAC)

Sensor Medium Access Control is a MAC protocol for sensor networks that is energy efficient. The major applications of S-MAC include tolerance latency and long idle listening. S-MAC's communication occurs among nodes, which also act as each other peers, rather than as solo base stations [8]. In

addition to the maintenance of collision avoidance, it is equally important to adjust scalability to maintain or improve energy efficiency. S-MAC [9] achieves energy efficiency by reducing energy usage from every main source that is responsible for the excessive use of energy. In return, it permits partial performance degradation in both latency and per-hop fairness.

### B. Synchronization in S-MAC

Synchronization mainly depends on all the three phases which includes:

- Listening period
- Wake-up period
- Sleep period

Following the standards of 802.11 IEEE, numerous efficient control methods such as sleeping, collision avoidance, synchronization, and listening are usually coupled with a contention-based MAC while implementing S-MAC. In principle, S-MAC employs a cyclic wake-up method wherein each node has its listening and sleeping period of the definite length according to its schedule [10]. In this method, every node switches to sleep mode for a defined period and then wakes up for subsequently fixed listening as shown in Fig. 1.

Listening Node - Listening mode is further broken down into three phases:

- SYNC phase
- RTS phase
- CTS phase



Fig. 1. Synchronization in S-MAC

Below is the description of each synchronization phase of S-MAC:

**SYNC phase**: In this type of phase, node A collects the transferred and corresponding packets from its adjacent B node and C node. It then generates a table of the listening periods of adjacent (A, B) nodes in a designed table. The A's SYNC stage is additionally separated into intervals of time [11]. When the adjacent node B is struggling to send a SYNC packet, it is picked at any given time, arbitrarily placed, and then begins the sending process in case no other signal is found in previous vacancies. If any signal is found, node B returns to sleep mode, and it waits in this mode until node A wakes up. In other words, A keeps count of B's timetable. It is possible for A to wake up for reasonable periods to send its SYNC data packet to B in the mode of broadcasting. This part is the synchronization.

**Request to send (RTS) phase**: Following synchronization, nodes start receiving transferred signals or packets from the adjacent nodes. RTS/CTS handshake is employed to avoid a signal collision.

**Clear to send (CTS) phase**: In this phase, CTS is transmitted upon receiving the RTS data packet from the adjacent nodes. Following this transmission is the beginning of the exchange of the packet, and this exchange continues for the optimum sleep time of A.

On synchronizing the schedules of A and its adjacent nodes, all nodes wake up at the same moment, and one and only one packet of SYNC is to be transmitted by A in order to reach all of its adjacent nodes. With the immense help of the S-MAC procedure, all the adjacent nodes reach an agreement on the same limited time table which then becomes the basis of shaping virtual clusters. The transmission of data packets [12] is not at all hindered by clustering as it only performs the exchange of schedules.

S-MAC protocols continue creating virtual clusters., Upon installing and switching, node A listens for a pre-defined and known synchronized time. If node A gets any SYNC packet from its adjacent nodes, say node B, it starts following its timetable and starts switching and transmitting packets the moment that node B enters into its listening period. Another scenario is that node A chooses a timetable for itself and starts transmission accordingly. In between the contention period of the transmission of the packet, if node A gets the timetable of any of its adjacent nodes, it strictly starts following it after dropping its schedule. If node A gets a signal that its adjacent nodes are following its schedule, then it prefers to stick to its schedule and, most likely, A will begin transmitting the data packets following both schedules. In this case, node A comes to know that neither of its adjacent nodes overlaps with its timetable, then it again starts following the timetable of its adjacent node and drops its own. However, node A always can get an invalid SYNC data packet and, in this case, it starts listening to its adjacent nodes to attain an entire synchronization period [13]. Due to this arrangement of virtual clusters, this instance is reasonably difficult.

The big multi-hop network is categorized into 'islands of timetable harmonization'. The nodes which are at the verge of a virtual cluster follow more than one timetable to forward their SYNC data packets. Hence, such nodes use more energy than the nodes with adjacent nodes of the identical timetable. Because of the cyclic wake-up method in S- MAC, the nodes spend most of their time in sleep mode which has a good reputation in terms of battery usage despite their latency. One of the major drawbacks of S- MAC is that it becomes quite challenging to catch up with the time duration while switching from the wake-up phase to the altering load states [14]. This is due to the limited time duration of the listening mode.

### C. Timing Relationships in S-MAC

The scheduling is done by sending the SYNC packet, which is a small packet that contains sender information such as its address and its next sleep time. When the sender starts transmitting the SYNC packet, its next sleep time is related to that moment. After the reception of the SYNC packet, the receiver gets the time from it and subtracts this time from the

transmission time of the packet [13]. For the nodes which are to receive both the data and SYNC packet, the listening interval time is divided into two parts where the SYNC packet gets the first part, and the data packet gets the second part of the time interval, as Fig. 2 shows. There is a contention window for each part which has time slots for the sender to do virtual carrier sensing before transmitting. If the SYNC packet [14] is going to be transmitted, the sender has to do carrier sensing when the receiver begins listening.



Fig. 2. Timing Relationship in S-MAC

1) The first part of the figure shows that sender 1 has done the carrier sensing and after getting access to the To end, the carrier sensing it selects any random time slot. If there is no transmission detected by the node at the end of the time slot that it got randomly, then it gets the contention window for the transmission of the SYNC packet and transmits the SYNC packet. Similarly, when the data packet is sending, the same procedure is adopted by the nodes. Fig. 2 shows three different scenarios in which there is communication between the sender and the receiver.
the channel, only the SYNC packet is sent for synchronization with the node.

2) In the second part of the figure, sender 2 has done the carrier sensing, sending the RTS packet and then the data packet if it received the CTS packet. It shows the transmission of a unicast packet.

3) The third part of the figure shows that sender 3 has done the carrier sensing and then sent the SYNC packet. After that, it again did the carrier sensing and sent the RTS packet. If it got the CTS packet, it would transmit the data packet. It shows the sending of both the SYNC and data packet.

### III. RELATED WORK

While talking about networks, with the recent advances in technology world Wireless Sensor Network (WSN) became one of the most capable network solution to most of the communication problems [29] that might be encountered in the field of Health, Military and Agriculture etc. It's obvious that there are two sides of a picture, along with this positive and promising side of WSN there is a weak and vulnerable side too. Because in WSN the involved sensor nodes are also vulnerable to some serious security attacks and the reasons vary from the gaps left in their deployments to some sensitive nodes that might be left unattended [30] because of being present in the areas where there is no one to properly check their security.

Our article has also focused on some serious security attacks on WSN. Authors in [20] have also focused on one of the most occurring attacks on WSN i.e. DoS Jamming attack. This attack works by sending a huge volume of illegitimate traffic to the node in order to jam the legitimate traffic and thus the network. The technique proposed in this article exponentially weighted moving average (EWMA) is used to detect abnormal changes in the intensity of jamming attack. While the authors in [21] have studied the protocols related to access control in WSN. Their work has focused on the authentication problem during accessing the networks. They have analyzed different access control protocols and also include discussion on replacement of expensive protocols with some affordable ones. The article [22] discusses about the Intrusion Detection System (IDS) proposed for WSN. They have done a survey on different IDS for WSN and also for Mobile Ad-Hoc Networks (MANET). Lastly, they have proposed an IDS scheme for WSN after comparing the existing schemes along with their weaknesses. An Intrusion Detection System has been proposed in [23] for the detection of sinkhole attack in WSN. The article included the attack implementation to check their proposed IDS followed by the in depth study of Sinkhole attack. Another work done on a commonly occurring Worm attacks in WSN in [24] focusing specifically on prevention of the sensor worms from propagating in the entire network. They have proposed an algorithm for assigning the relevant version of software to each sensor node in the sensor network to restrain the worm propagation. Researchers in [25] have done an analysis of a number of different security issues related to data integrity, data availability and data confidentiality. They have analyzed different security attacks on WSN i.e. a number of different Passive attacks, Denial of Service attacks, physical attacks, false node attack, etc.

By going through the related work it became pretty obvious that our work is different from other works. The work discussed above have either proposed an IDS scheme for a single security attack on WSN or they have just done a survey. in a way that we have not only focused on multiple attacks. But we have focused not only on multiple attacks with their in-depth study but also presented and IDS for detection of multiple WSN security attacks.

### IV. ATTACKS ON S-MAC

There are several kinds of common network attacks on S-MAC such as a Collision Attack, Unfairness Attack, Exhaustion Attack, Sinkhole Attack and Wormhole Attack [15]. We have briefly discussed all of the attacks with their detection methods and detection mechanisms. Below are the brief descriptions of network attacks on S-MAC:

- Collision Attack

- Unfairness Attack

- Exhaustion Attack

- Sinkhole Attack

- Wormhole Attack

#### A. Collision Attack

This attack occurs when legitimate nodes tend to communicate with each other, and the rival nodes start sending

data packets in these overlapping periods in order to hinder legitimate communication. In this way, the packet sent by the legitimate node gets lost, and the node has no other choice but to wait to find or acquire another medium for transmitting the RTS/ CTS packet. In order for a node to retransmit, it spends its energy over and over again for the very same packet, and this consumption eventually results in the reduction of energy. Normally, one byte is enough for making a CRC error in addition to disabling the data packet [16].

Advantages

- Power is consumed periodically in each data packet and is difficult to detect.

- This attack also culminates the ACK packet which results in exp. backoff message and wastes the battery.

- It can be launched anywhere in the entire network, and the attacker does not have extra capabilities.

- It weakens data integration in the MAC layer.

Detection Method

- Misbehaviour detection techniques

Defensive Mechanism

- All countermeasures of congestion attacks

### B. Unfairness Attack

In S-MAC procedures, control is entirely employed by each node. All nodes transmit RTS packets so that they can request to attain the medium in which a CTS packet is sent back to the demanded, desired node. The medium is competed for in a particular node in each time vacancy. The first node tends to attain the medium, but the illicit nodes also get the benefit of this method and transmit the data packet with a low waiting time duration. It repeatedly transmits such packets to attain the medium and, hence, causes a hindrance for the legitimate nodes to have maximum access to the medium [16].

Attack effects

- Decrease in the services of effective networks.

- Nodes are desperate to have access to the medium.

- Limited access of nodes to the medium and can paralyze the usual communication within the medium.

Detection Methods

- Misbehaviour detection techniques

Defensive Mechanism

- Employment of smaller frames

### C. Exhaustion Attack

The legitimate and rival nodes are installed in WSN in an open milieu. S-MAC procedures deal with the CTS/RTS method and are known for their capacity of transmission. Therefore, while attaining the medium, a node should transmit an RTS data packet and, in return, the receiver should send a CTS data packet [17]. The rival nodes [18] usually take

advantage of this method, and the demanded node constantly transmits CTS data packets, as a result of which the network gets weakened or gets overloaded with the amalgamation of both legitimate and illegitimate nodes, further resulting in an exhaustion attack.

### D. Sinkhole Attack

In such attacks, the rival nodes (or the compromised nodes) exhibit their attractiveness to the nodes to illustrate all of their traffic data from their constituency. Therefore, all the data packets' transfer is intended for the base station which in turn is drawn by the rival nodes. In order to have full autonomy over the data transfer [28], the compromised node aims for its adjacent nodes. The sinkhole attack [10] is started by the rival nodes from the adjacent nodes that are very close to the base station.

### E. Wormhole Attack

In a wormhole attack, a link is developed by the illegitimate/ compromised node between two specified points in the network and this link is known as the wormhole link. This kind of directness of the wormhole link can only be made with the help of wireless transmission, optical fibre or wireline. The moment this type of direct connection (wormhole link) [5] is developed, the communication is captured by the rival nodes, and they channel the nodes from the origin to the other endpoint, known as the destination point.

Attack Effects

- False routing information.

- Alteration in the network topology.

- Packets alteration by wormhole nodes.

- Alteration in the normal flow of messages [28].

## V. METHODOLOGY

In this section, we are going to discuss the analysis of Collision and Exhaustion attacks on the Wireless Sensor Nodes. For that, we have used a simulation model to implement a secure MAC in MATLAB. For the simulation model, we have used the following vital parameters.

### A. Simulation Model

Table I shows the parameters used for the simulation model along with their values:

TABLE I.     YEAR WISE TREATMENT FREQUENCIES

| Parameters | Values |
|---|---|
| Sensing Area Dimensions (X * Y sq. m) | 50 m x 50 m |
| Legal nodes | 14 |
| Intruder nodes | 5 |
| Sink node | 1 |
| Dimension of sink node (X * Y sq. m) | 25 x 25m. |
| Transmission Energy | $50\mu J$ |
| Received Energy | $30\mu J$ |
| Idle Energy Consumed | $5\mu J$ |
| Data Rates | 250kbits/s |

## B. WSN Deployment

Fig. 3 below shows the wireless sensor nodes' network deployment. There are 20 nodes in which 14 nodes are the transmitting nodes, and one is the sink node. The five (5) nodes are acting as intruders. We have deployed these nodes in the x and y plane. All the nodes are going to transmit the data to the sink node which is located in the x and y plane at the location of (25, 25) in meters. Different locations are given to the nodes. The nodes are deployed randomly.



Fig. 3. Nodes deployed in our wireless network

## C. Key Features of our Defensive Mechanism

**Efficient Network Energy Utilization** - Our exhaustion attack on the MAC layer node repeatedly sends RTS and CTS packets, due to which energy is consumed each time by sending these control packets again and again. Our defensive mechanism detects collision attacks which detect intruder nodes, and thus we can save network energy in this way.

**Network Quality** - S-MAC collision attacks are continuously being introduced in a WSN environment which loses the data packets and hence decreases the successful data packet transmission rate, which in turn affects network quality. Our defensive mechanism detects collision attacks efficiently. Decreasing the collision ratio in each node can improve network quality.

## D. Intrusion Detection Mechanism

We chose the following statistics as intrusion indicators in the intrusion detection part:

**Collision Ratio (Rc)** - It is defined as the detection of the collision time for a node per second.

**Probability of data packets' successful Transmission (PST)** - A successful transmission can be defined as the sending and receiving of a packet by a node correctly. The probability of data packets' successful transmission is the ratio between the successful transmissions to the total number of data packets transmitted.

**RTS packet arrival ratio (RRTS)** – Defined as the number of RTS packets successfully received by a node per second. We collected the values of all indicators and estimated the intrusion probability. According to these values, we can conclude whether there is an intruder or not. For this we use the soft function which is given below:

$$y(x) = \frac{1}{1 + \exp\left[-A \times (x - C)\right]} \tag{1}$$

In Equation 1, $A$ is the slope parameter. If the value of $A$ is bigger, the slope is steeper, and $C$ is the centre of the curve. From this equation, we can calculate the probabilities of the above-discussed intruder.

We generated the random values for collision, data packets' successful transmission and the RTS data packet arrival ratio (RRTS). We put the values of all these intruders in the above equation and found out the probability for each of them. The following are the results we achieved after putting random values in the above formula:

1) By inputting the values of the collision ratio in the soft function, we get the probability of collision, which is called PC.
2) By inputting the values of data packets' successful transmission ratio, we get the probability of the total, which is called PT.
3) By inputting the values of RTS data packets' arrival ratio, we get the probability of exhaustion, which is called Pe.

The shape of the curve can be adjusted by changing the $A$ and $C$ parameters. We can find the next values of $A(k)$ and $C(k)$ from the equations below:

$$A(k+1) = A(k) + \alpha \times \frac{\partial J}{\partial A} \tag{2}$$

Where $A(k)$ is the initial condition value; $\alpha$ is a value between 0 and 1.

$\frac{\partial J}{\partial A}$ is given us:

$$\frac{\partial J}{\partial A} = 2(yd - y)A(k) / \left(1 + \exp^{A(k)*(x - C(k))}\right)^2 \tag{3}$$

Where the actual value is $y$ and the desired value is $yd$.

$$C(k+1) = C(k) + \alpha \times \frac{dJ}{dC} \tag{4}$$

Where $C(k)$ is the initial condition value; $\alpha$ is a constant value between 0 and 1.

$$\frac{dJ}{dC} = 2(yd - y)\frac{-A(k)\exp\left(-A(k) \times [x - C(k)]\right)}{\left(1 + \exp(-A(k) \times [x - C(k)])\right)^2} \tag{5}$$

In Equation 5, the actual value is $y$, and the desired value is $yd$.

## E. Criteria for Attack Detections

*Criteria for detection of Exhaustion Attack:* Fig. 4 shows the criteria for the exhaustion attack. We got the value of exhaustion from the soft function by inputting the values of the RTS data packets' arrival ratio. The probability of success is added to the probability of exhaustion and compared with the threshold. The threshold is set, and then this summation

result is compared with the threshold. If the sum is greater than the threshold, then the attack is found. Otherwise, there was no attack. We set the threshold higher than the probability of success for our results.



Fig. 4. Criteria for detection of Exhaustion Attack

*Criteria for detection of Collision Attack::* Fig. 5 shows the criteria for the detection of a collision. In this figure, it is clear that when we got the probability of success from the soft function, then it is scaled by phi ($\phi$). Similarly, we got the probability of collision from the soft function by inputting the values of the collision ratio, then scaled by theta ($\theta$). The probability of success is summed individually with the probability of collision. As the above process shows, the probability of success is added to the probability of collision and then compared with the threshold, thus attaining a result. The threshold is set, and then this summation result is compared with the threshold. If the sum is greater than the threshold, then there was an attack. Otherwise, no attack occurred. We set the threshold higher than the probability of success for our results.



Fig. 5. Criteria for detection of Exhaustion Attack

## VI. EXPERIMENT RESULTS AND ANALYSIS

### A. When No Collision Attacks were Found

Fig. 6 shows the graph as it appears when the values of E and F are changing. The next values of E(N) and F(N) are defined above. The graphs are:

- The square graph (green) is the graph of the probability of collision.

- The circle (blue) graph shows the probability of success.

- The dashed graph shows the threshold.

- The (red) solid line graph shows the sum of probability of success after multiplying with phi ($\phi$) and the probability of collision after multiplying with theta ($\theta$).



Fig. 6. No Collisions Found

### B. When Collisions are Found

In Fig. 7, the (sum) solid line graph is greater than that of the threshold dashed graph, which means there is an intruder in the network.

### C. Comparison of Delay to show Collision Attack

Fig. 8 shows the behaviour of nodal PS+PE transmission when there is no intruder and when there is an intruder. The delay of nodes is not high when there is no intruder. As the solid graph (blue) shows, it is clear that the delay of nodes transmitting without an intruder in the network is less than that of the delay when there is an intruder in the network. The solid (blue) graph describes the delay of nodes when there is no intruder in the network. When there is an intruder in the network, then the delay rises abruptly and increasing



Fig. 7. Collisions Found

as compared to that of delays of nodes when no intruder is present. This unexpected increase in the delay shows that a collision occurred while transmitting the data. Therefore, the delay increased because of the retransmission of the packet. The (green) dashed graph shows the delay when there is an intruder in the network. This abrupt change in the graph is indicative of an intruder. When collisions occur, the node will retransmit the packet, and in this way, the delay is increased.



Fig. 8. Behaviour of nodes

### D. When no Exhaustion Attacks are found

Fig. 9 shows a comparison of the probability of success with the probability of exhaustion. This graph is also a multigraph in which:

- The square graph (green) shows the probability of exhaustion.

- The circle graph (blue) shows the probability of success.

- The dashed graph shows the threshold.

- The solid (red) one is the sum graph of the probability of success and probability of exhaustion.

We completed a comparison of the sum graph with the threshold graph. If the sum graph is greater than that of the dashed graph (threshold), then it means there is an intruder in the network. So, here there is no intruder.



Fig. 9. No Exhaustion Attacks Found

### E. When Exhaustion Attacks are found

Fig. 10 clearly shows that there is an intruder since the sum solid line graph (red) is greater than that of the dashed graph which represents the threshold.



Fig. 10. Exhaustion Attack Found

### F. When Exhaustion Attack is Detected

Fig. 11 clearly shows that the exhaustion attack is detected.



Fig. 11. Exhaustion Attack Detected

### G. Exhaustion Attack

The graph in Fig. 12 shows the behavior of sensory nodes before and after detection of an exhaustion attack.

- The circle graph (blue) shows the exhaustion attack (repeated RTS packet arrival) and the waste of energy in (response of CTS) the presence of the intruder.

- The square graph shows energy dissipation after the detection of the intruder.

## VII. DISCUSSION

Our work has targeted the WSN which on one side is a promising communication network and provides the user a number of benefits like lower cost, lower power consumption, and easy deployment [26] and also supports a number of important real-life applications. While on the other side, Security

Fig. 12. Behaviour of Nodes

is also becoming a concerning issue for the WSNs because mostly the networks' nodes are deployed in some antagonistic area where the nodes have small memory, limited amount of energy [27] and this results in occurrence of a number of security attacks which sometimes jam the legitimate network traffic and can also leak sensitive information.

Security being a sensitive issue for WSN and the security attack targeting the network layer resulting in disturbed communication in important real life environment like military, we chose to analyze the Intruder attacks like Unfairness attack, Exhaustion attack, Sinkhole attack and Wormhole attack on S-MAC in a sensor network. We have presented an Intrusion detection system (IDS) to detect the attacks and an Intrusion Prevention System (IPS) to prevent those attacks in future. But this is done at smaller level to first check whether it will give us the expected results or not. Our simulation results showed the results as expected. This proved the effectiveness of our soft IDS and IPS mechanisms. We might continue this work to make it able to be implemented in some real- time environment for the detection of Intruder attacks to make it more fruitful.

## VIII. CONCLUSION

In this research, we have focused on securing a Wireless Sensor Network against collision and exhaustion attacks. As in MAC protocol of WSN, the intruder can disturb the performance of the whole WSN by just depleting the battery source of legitimate nodes. We have adopted an Intrusion detection mechanism for the detection of different kinds of intruders' attacks. Then we developed a soft decision mechanism to detect intrusions and exhaustion attacks in WSN. A preventative mechanism has also been developed which helps the node to avoid such intruder attacks. The simulation results have proved the effectiveness of our mechanism by showing how a lifetime of a node has increased along with the network performance with reduced energy consumption, reduced delay on the node and successful transmission. In the future, we will be focusing on other emerging attacks on the MAC layer as well as S-MAC. In the future, we plan to extend our implemented detection mechanism to defend against other emerging attacks like DoS attack on WSN. This extension can lead to its implementation in different real-life scenarios like battlefields, health departments and also for other critical missions that need a cutting edge method for secure data transmission.

## REFERENCES

[1] F. Dong., J. Yang., C. Xiong., H. Ding, and Y. Zhang, "Research and implementation of a hybrid mac protocol for wireless sensor networks based on clustering structure," *Int J Recent Sci Res.*, vol. 9, no. 10, pp. 29 131–29 134, 2018.

[2] F. Z. Djiroun and D. Djenouri, "Mac protocols with wake-up radio for wireless sensor networks: A review," *IEEE Communications Surveys & Tutorials*, vol. 19, no. 1, pp. 587–618, 2016.

[3] K. A. Memon, M. A. Memon, M. M. Shaikh, B. Das, K. M. Zuhaib, I. A. Koondhar, and N. U. A. Memon, "Optimal transmit power for channel access based wsn mac protocols," *INTERNATIONAL JOURNAL OF COMPUTER SCIENCE AND NETWORK SECURITY*, vol. 18, no. 7, pp. 51–60, 2018.

[4] R. Ramya, G. Saravanakumar, and S. Ravi, "Mac protocols for wireless sensor networks," *Indian Journal of Science and Technology*, vol. 8, no. 34, p. 1, 2015.

[5] D. Mohammed, M. Omar, and V. Nguyen, "Wireless sensor network security: Approaches to detecting and avoiding wormhole attacks," *Journal of Research in Business, Economics and Management*, vol. 10, no. 2, pp. 1860–1864, 2018. [Online]. Available: http://scitecresearch.com/journals/index.php/jrbem/article/view/1413

[6] V. Casola, A. De Benedictis, A. Drago, and N. Mazzocca, "Analysis and comparison of security protocols in wireless sensor networks," in *2011 IEEE 30th Symposium on Reliable Distributed Systems Workshops*. IEEE, 2011, pp. 52–56.

[7] K. Chelli, "Security issues in wireless sensor networks: Attacks and countermeasures," in *Proceedings of the World Congress on Engineering*, vol. 1, no. 20, 2015.

[8] A. Rai, S. Deswal, and P. Singh, "Mac protocols in wireless sensor network: a survey," *International Journal of New Innovations in Engineering and Technology*, vol. 5, no. 1, pp. 95–101, 2016.

[9] Y. Rao, Y.-m. Cao, C. Deng, Z.-h. Jiang, J. Zhu, L.-y. Fu, and R.-c. Wang, "Performance analysis and simulation verification of s-mac for wireless sensor networks," *Computers & Electrical Engineering*, vol. 56, pp. 468–484, 2016.

[10] M. Sharma, A. Tandon, S. Narayan, and B. Bhushan, "Classification and analysis of security attacks in wsns and ieee 802.15. 4 standards: A survey," in *2017 3rd International Conference on Advances in Computing, Communication & Automation (ICACCA)(Fall)*. IEEE, 2017, pp. 1–5.

[11] S. Kaur and S. Sharma, "A review on various routing protocols based on clustering in wsn." *International Journal of Advanced Research in Computer Science*, vol. 8, no. 7, 2017.

[12] W.-M. Song, Y.-M. Liu, and S.-E. Zhang, "Research on smac protocol for wsn," in *2008 4th International Conference on Wireless Communications, Networking and Mobile Computing*. IEEE, 2008, pp. 1–4.

[13] G. Gautam and B. Sen, "Performance analysis of 802.11 and smac protocol under sleep deprivation torture attack in wireless sensor networks," *International Journal of Computer Sciences and Engineering*, vol. 3, no. 5, pp. 317–322, 2015.

[14] S. Otoum, M. Ahmed, and H. T. Mouftah, "Sensor medium access control (smac)-based epilepsy patients monitoring system," in *2015 IEEE 28th Canadian conference on electrical and computer engineering (CCECE)*. IEEE, 2015, pp. 1109–1114.

[15] K. Pelechrinis, M. Iliofotou, and S. V. Krishnamurthy, "Denial of service attacks in wireless networks: The case of jammers," *IEEE Communications surveys & tutorials*, vol. 13, no. 2, pp. 245–257, 2011.

[16] P. Sinha, V. Jha, A. K. Rai, and B. Bhushan, "Security vulnerabilities, attacks and countermeasures in wireless sensor networks at various layers of osi reference model: A survey," in *2017 International Conference on Signal Processing and Communication (ICSPC)*. IEEE, 2017, pp. 288–293.

[17] T. Borgohain, U. Kumar, and S. Sanyal, "Survey of security and privacy issues of internet of things," *arXiv preprint arXiv:1501.02211*, 2015.

[18] M. Tiloca, D. De Guglielmo, G. Dini, G. Anastasi, and S. K. Das, "Jammy: A distributed and dynamic solution to selective jamming attack in tdma wsns," *IEEE Transactions on Dependable and Secure Computing*, vol. 14, no. 4, pp. 392–405, 2017.

[19] M. Pawar and J. Agarwal, "A literature survey on security issues of wsn and different types of attacks in network," *Indian J. Comput. Sci. Eng*, vol. 8, pp. 80–83, 2017.

[20] Osanaiye, Opeyemi and Alfa, Attahiru and Hancke, Gerhard, "A statistical approach to detect jamming attacks in wireless sensor networks," *Sensors, Multidisciplinary Digital Publishing Institute*, vol. 18, pp. 1691, 2018.

[21] V. Mittal, S. Gupta, and T. Choudhury, "Comparative Analysis of Authentication and Access Control Protocols Against Malicious Attacks in Wireless Sensor Networks," pp. 255–262.

[22] I. Butun, S. D. Morgera, and R. Sankar, "Wireless Sensor Networks," vol. 16, no. 1, pp. 266–282, 2014.

[23] I. Krontiris, T. Dimitriou, and T. Giannetsos, "Intrusion Detection of Sinkhole Attacks," pp. 150–161, 2008.

[24] Chen, Honglong and Lou, Wei and Wang, Zhi and Wu, Junfeng and Wang, Zhibo and Xia, Aihua, "Securing DV-Hop localization against wormhole attacks in wireless sensor networks",*Elsevier, Pervasive and Mobile Computing*, pp. 22–35, vol. 16, 2015.

[25] K. S. Selvam and S. P. Rajagopalan, "Security Analysis with respect to Wireless Sensor Network – Review," vol. 6, no. 4, 2017.

[26] Sharma, Kalpana and Ghose, MK and others, "Wireless sensor networks: An overview on its security threats", *IJCA, Special Issue on "Mobile Ad-hoc Networks" MANETs* , pp. 42–45, 2010.

[27] Dinker, Aarti Gautam and Sharma, Vidushi, "Attacks and challenges in wireless sensor networks", *2016 3rd International Conference on Computing for Sustainable Global Development (INDIACom)*, pp. 3069–3074, 2016.

[28] M. Pawar and J. Agarwal, "A literature survey on security issues of wsn and different types of attacks in network," *Indian J. Comput. Sci. Eng*, vol. 8, pp. 80–83, 2017.

[29] Can, Okan and Sahingoz, Ozgur Koraytitle, "A survey of intrusion detection systems in wireless sensor networks", *IIEEE, 2015 6th International Conference on Modeling, Simulation, and Applied Optimization (ICMSAO)*, pp.1–6, 2015.

[30] Shahzad, Furrakh and Pasha, Maruf and Ahmad, Arslan, "A survey of active attacks on wireless sensor networks and their countermeasures",*arXiv preprint arXiv:1702.07136*,2017.

# Deep Learning Approaches for Data Augmentation and Classification of Breast Masses using Ultrasound Images

Walid Al-Dhabyani[1], Aly Fahmy[4]
Faculty of Computer and Information,
Cairo University, Cairo, Egypt

Mohammed Gomaa[2], Hussien Khaled[3]
National Cancer Institute,
Cairo University, Cairo, Egypt

*Abstract*—**Breast classification and detection using ultrasound imaging is considered a significant step in computer-aided diagnosis systems. Over the previous decades, researchers have proved the opportunities to automate the initial tumor classification and detection. The shortage of popular datasets of ultrasound images of breast cancer prevents researchers from obtaining a good performance of the classification algorithms. Traditional augmentation approaches are firmly limited, especially in tasks where the images follow strict standards, as in the case of medical datasets. Therefore besides the traditional augmentation, we use a new methodology for data augmentation using Generative Adversarial Network (GAN). We achieved higher accuracies by integrating traditional with GAN-based augmentation. This paper uses two breast ultrasound image datasets obtained from two various ultrasound systems. The first dataset is our dataset which was collected from Baheya Hospital for Early Detection and Treatment of Women's Cancer, Cairo (Egypt), we name it (BUSI) referring to Breast Ultrasound Images (BUSI) dataset. It contains 780 images (133 normal, 437 benign and 210 malignant). While the Dataset (B) is obtained from related work and it has 163 images (110 benign and 53 malignant). To overcome the shortage of public datasets in this field, BUSI dataset will be publicly available for researchers. Moreover, in this paper, deep learning approaches are proposed to be used for breast ultrasound classification. We examine two different methods: a Convolutional Neural Network (CNN) approach and a Transfer Learning (TL) approach and we compare their performance with and without augmentation. The results confirm an overall enhancement using augmentation methods with deep learning classification methods (especially transfer learning) when evaluated on the two datasets.**

*Keywords*—*Generative Adversarial Networks (GAN); Convolutional Neural Network (CNN); deep learning; breast cancer; Transfer Learning (TL); data augmentation; ultrasound (US) imaging; cancer diagnosis*

## I. INTRODUCTION

Medical imaging is a worthy tool to diagnose the presence of several diseases and the analyize of the experimental results [1]. Biomedical imaging is part of the foundations of overall cancer care. Breast cancer is well-known and widespread through women world-wide and it causes mortality rates. It is anticipated that more than eight percent of women will acquire breast cancer during their lifetime [2]. Digital Mammography (DM) is the most generally used and practical technique for breast cancer diagnosis [3]. Early detection is the most important factor in decreasing the costs of cancer management and mortality. DM imaging has some weaknesses in dense breasts where tumors can be hidden by surrounding tissue (where the dense tissue has a similar attenuation contrasted to the tumor). In practice, ultrasound (US) imaging is the best alternative to DM, which is applied as a complementary approach for breast cancer classification and detection due to its sensitivity, safety and versatility [4]. However, the weakness of US imaging is that it is hand-dependent which relies more on radiologists. Explaining US images needs specialist radiologists due to its difficulty and appearance of speckle noise. Therefore, Computer Aided Diagnosis (CAD) can help radiologists in the US-based classification and detection of breast cancer, reducing the influence of the hand-dependence of US imaging.

Some researches have studied the effect of CAD diagnostics [5], [6] and noted that CAD is a robust tool to enhance diagnostic specificity and sensitivity. Breast US research has a shortage of benchmark dataset which results in limiting the advancement of recent algorithms. Therefore, breast US images quality is extremely dependent on the acquisition process and there is a large variability between various US systems that affects the outputs achieved by algorithms. The output is also influenced by the size, location, and appearance of the tumor or micro-calcification.

Training a deep model on insufficient data regularly results in over-fitting because a model of high capacity is capable of "memorizing" the training set. Multiple methods have been presented to mitigate this problem, but none performed effectively as to be used exclusively. These techniques can be split into two large categories: (1) regularization techniques, pointing to limit the model's capacity (e.g., dropout and parameter norm penalty) and (2) data augmentation techniques, aiming to increase the size of the dataset [7]. In practice, most models improve from these two techniques. We concentrate on these two categories. GANs [8] are a family of unsupervised neural networks most generally utilized for image generation. Data augmentation has confirmed to be very efficient and is adopted universally in the field of deep learning [9], [10]. It is in fact so effective that it is being used even in tasks that include massive data [11]. The most common forms of augmentation include flipping, scaling, translating, rotating, blurring and sharpening. The goal of such transformations is to obtain a new image that contains the same semantic information as the original.

While augmentation most certainly helps neural networks

learn and generalize more effectively, it also has its drawbacks. In most cases, augmentation techniques are limited to minor changes on an image, as more "heavy" augmentations might damage the image's semantic content. Furthermore, the forms of augmentation one can use differ from problem to problem, making their application ad-hoc and empirical. For instance, medical images have to be mildly augmented as they follow strict standards (i.e., they are centered, their orientation and intensity vary little from image to image and many times they are laterally/horizontally asymmetric) [12]. Finally, augmentation techniques are applied to one image at a time and thus are unable to gather any information from the rest of the dataset.

### A. Problem Statement and Motivation

Although there are a lot of scientific researches in the process of classification and detection of cancer tumors using different types of modalities, breast US imaging has rare researches due to the shortage of public benchmark datasets. We utilize Data Augmentation Generative Adversarial Networks (DAGANs) [13] to make our dataset (BUSI) and dataset (B) larger. We particularly chose to use breast US imaging because US scan is safe for human body while DM and other screening technology may not achieve the same standard of safety as US imaging. Furthermore, We are proposing deep learning approaches for breast US imaging classification using state-of-the-art algorithms to improve the accuracy results using deep learning approaches proved to achieve promising results.

### B. Paper Contribution

- Due to the scarce number of datasets of ultrasound images for breast cancer, we believe that our dataset collection and data augmentation are an important contribution that can be a great seed for related studies. We plan to make our dataset publicly available for other researchers.

- We propose a novel augmentation technique that overcomes the above-mentioned limitations and is capable of augmenting any given dataset with realistic, high-quality images generated from scratch using DAGAN.

- We used two datasets which are our BUSI dataset and dataset B [14]. And we ran state-of-the-art deep learning models; CNN and TL as classification algorithms and they produced promising results.

- Finally, We merge the two datasets to overcome the limitation of the size of the dataset and compare the new results (the merged dataset) with the previous results (the two separate datasets). In addition, we will enlarge our datasets by combining them with traditional augmentation and DAGAN data to enhance the final results.

The remainder of this paper is divided as follows: Section II clarifies some related work in these fields. Subsequently, Section III illustrates the two Breast US datasets. Section IV discusses our methodology. Section V contains the results and discussion. And finally, Section VI has a conclusion and future work.

## II. RELATED WORK

In this section, related work for breast US image classification and data augmentation in medical images are reviewed. Furthermore, a brief introduction about deep learning for breast imaging is discussed.

### A. Breast US Image Classification

This section explains in brief three state-of-the-art approaches for tumor classification in breast US imaging.

*1) Convolutional Neural Networks (CNN):* Huynh et al. [15] assessed the performance of utilizing transferred features from pretrained CNNs [16] in classifying cancer in breast US images, and to examine this method of transfer learning with preceding methods including human-designed features. A breast US dataset composed of 1125 samples and 2392 Regions of Interest (ROIs) was utilized. Every ROI was annotated as malignant or benign. Features were extracted from each ROI using pre-trained CNNs and used to train Support Vector Machine (SVM) classifiers in the tasks of distinguishing benign vs malignant tumors. For a baseline comparison, classifiers were also trained on prior analytically-extracted tumor features. They conducted five-fold cross-validation with the area under the receiver operating characteristic curve (AUROC) as the performance metric. Classifiers trained on CNN-extracted features were comparable to classifiers trained on human-designed features. In the task of malignant versus benign, the SVM trained on both CNN-extracted features and human-designed features achieved an AUC of 90%. In the task of determining benign vs malignant, the SVM trained on human-designed features achieved an AUC of 85%, compared to the AUC of 85% achieved by the SVM trained on CNN-extracted features. The authors obtained great results using transfer learning to characterize ultrasound breast cancer images. This method allows them to instantly classify a little dataset of lesions in a computationally reasonable fashion without any hand-operated input. Current deep learning approaches are dependent on huge datasets and large computational resources, which are frequently difficult to access for clinical applications. It is important to highlight that, the dataset of this study [15] is not publicly available neither by request.

*2) Stacked Deep Polynomial Network (S-DPN):* Jun Shi et al. [17] proposed Deep Polynomial Network (DPN) [18] algorithm not just presents better performance on a massive dataset, but also has the possibility to learn strong characteristic representations from a comparatively little dataset. In their study, a S-DPN algorithm is suggested to further enhance the representation performance of the primary DPN, and S-DPN is then used to the task of texture feature learning for US classification of tumor with a little dataset. The task of tumor classification is achieved on two datasets, namely the prostate US elastography dataset and breast B-mode US dataset. On these two cases, results of the experiment confirm that S-DPN achieves the best classification performance with accuracies of 92.40% on breast US dataset and 90.28% on prostate US datasets. It is important to highlight that, the dataset of this study [17] is not publicly available.

*3) Shearlet-based Texture Feature Extraction:* Zhou et al. [19] augmented the classification accuracy of the US computer-aided diagnosis (CAD) for the detection of breast tumor

based on texture feature, they also offered to use Shearlet transform to achieve texture feature descriptors. Shearlet transform produces a scattered representation of high-dimensional data with especially higher directional sensitivity at different scales. Hence, texture feature descriptors of shearlet-based can strongly explain breast tumors. In order to accurately evaluate the achievement of Shearlet-based features, curvelet, contourlet, and wavelet-based texture feature descriptors are also obtained for comparison. All these features were then fed to two different classifiers, AdaBoost and support vector machine (SVM), to estimate the consistency. The results of the experiment of breast tumor classification presented that the classification accuracy, specificity, sensitivity, negative predictive value, positive predictive value and Matthew's correlation coefficient of shearlet-based method were 91.0%, 92.5%, 90.0%, 90.3%, 92.6%, 0.822% by SVM, and 90.0%, 90.0%, 90.0%, 89.9%, 90.1%, 0.803% by AdaBoost, respectively. Most of the results of the Shearlet-based significantly exceeded those of other approach based results under both classifiers. They suggested a new texture feature extraction approach based on Shearlet transform for describing breast tumor in US image. The comparative experiment results showed that the Shearlet-based texture feature can more efficiently identify breast tumors in US image than other features extracted from curvelet, contourlet, wavelet and Gray-Level Co-Occurrence Matrix (GLCM) approaches. It is important to highlight that, the dataset of this study [19] is not publicly available.

### B. Data Augmentation

GANs have been successfully used for data augmentation. Wang et al. [20] and Antoniou et al. [13], for example, use custom GAN architectures in a low-data setting to achieve consistently better results than traditionally augmented classifiers, while Perez et al. [21] devise a novel pipeline called Neural Augmentation which, through style transfer techniques, aims at generating images of different styles, performing equally as good as traditional augmentation schemes in a subsequent classification task. Additionally, Neff [22] proposes a generative model which learns to produce pairs of images and their respective segmentation masks in order to assist a UNet segmentation model, proving that in simpler datasets networks trained with a mix of synthetic and real images stay competitive with networks trained on strictly real data using usual data augmentation.

One field in which data augmentation is especially important is that of medical imaging, where the lack of available public data is a ubiquitous problem since access to individual medical records is heavily protected by legislation and appropriate consent must be given. In most cases, this process is hindered by bureaucracy and/or high costs, while the resulting collection is greatly imbalanced towards normal subjects. Several authors employ Machine Learning techniques to learn directly from the available data and surpass the state-of-the-art in problems as diverse as generating benchmark data, cross-modality synthesis, super-resolution or image normalization [23].

The medical field has only recently started adopting GAN-based methodologies for synthesizing images [24]. In particular, Bentaieb et al. [25] and Shaban et al. [26] proposed GAN-based style transfer approaches to stain normalization in histopathology images, with quite interesting results in various datasets. For tackling segmentation tasks, various authors have proposed custom GAN architectures and pipelines which are adversarially trained to produce proper segmentation masks from a given medical image dataset [27]–[29]. Regarding image translation between modes, the authors of [30] synthesize T2-weighted brain MRI images from T1-weighted ones, and vice versa, using a Conditional GAN model. Finally, many authors, such as [31] and [32], have attempted to generate counterfeit medical images in order to increase the size of the training set of different deep learning models, a task more closely related to the one examined in this study.

Supplementary to all of the above efforts, our approach aims to exploit the superior performance of GANs for the benefit of medical image classification. We explore the impact of GAN-assisted data augmentation on the diagnosis of breast cancer through US scans.

### C. Deep Learning for Breast Imaging

In general, the state-of-the-art classification methods are not robust, specifically the image processing based methods, relying on special assumptions and rule-based methods. Without necessitating such a powerful hypothesis, deep learning approaches have shown an improved accuracy in object classification and detection, which proposed that could also improve the state-of-the-art of tumor classification in breast ultrasound. Deep learning in medical imaging is usually represented by convolutional networks. GANs [8] are a family of unsupervised neural networks most usually used for image production. Each GAN is formed of two networks: a generator and a discriminator, playing against each other in a two-player game. These models have proven to be capable of creating realistic images and will serve as an assisting basis for this study. DAGAN is also used to make the dataset larger. Based on how we can train them, they can be frequently categorized into the following categories:

1) **CNNs approach.** This method trains the CNNs with images for training and testing [33], [34]. However, feeding every image to the network is time-consuming [35].
2) **Transfer learning approach.** Another approach that has been extensively used recently in biomedical research is the transfer learning technique [15], [36]. This method uses a pretrained model from natural images to overcome the lack of data in medical imaging study.
3) **Generative Adversarial Networks.** This method allows us to generate new images from our dataset. GAN [8] is a strong and new approach in image synthesizing.

In breast imaging, the majority of the current publications are focusing on using CNNs for MG. Dhungel et al. [37] have performed masses segmentation using deep learning; Mordang et al. [38] introduced the use of CNNs in microcalcification detection; and lately, Ahn et al. [39] suggested the use of CNNs in breast density evaluation. In breast US imaging, Huynh et al. [15] suggested the use of a transfer learning approach for breast US images classification. Yap et al. [14] proposed to use deep learning approaches for classification

of breast US tumor. As of the date of this publication, this is the only work the authors have found that handles breast ultrasound but it does not enhance the accuracy in tumor classification. Most of the aforementioned work focused on lesion detection. Furthermore, publications utilizing data augmentation with GAN are rare. In medical images, Frid-Adar et al. [31] proposed the use of DAGAN to enhance CNN performance in liver lesion classification. We are, in our consideration, the first to use DAGAN with breast US images. In this paper, we propose to use deep learning approaches for breast US tumors classification. To show the benefits of deep learning approaches, we compare the performances among all the deep learning approaches which are used in this paper for tumor classification. Furthermore, DAGAN and traditional augmentation are used to make the dataset larger and enhance the performance of our classification approaches.

## III. DATASETS

In general, to develop a healthcare system using deep learning, a dataset should be available. This study uses two different datasets of breast US images. Our dataset BUSI was collected and obtained from US systems with different specifications and at different times. The Dataset B [14] was requested from its owners. Examples of both datasets are shown in Fig. 1.

Dataset BUSI collected at baseline includes ultrasound breast images among women in ages between 25 to 75 years old. The number of patients is 600 female patients. It was collected in 2018 from Baheya Hospital for Early Detection and Treatment of Women's Cancer, Cairo (Egypt) with LOGIQ E9 ultrasound system and LOGIQ E9 Agile ultrasound. The data is categorized into three classes, which are normal, benign, and malignant. The dataset consists of 780 images from different women with an average image size of 500 x 500 pixels. Within the 780 tumor images, 133 were normal images without cancerous masses, 437 were images with cancerous masses and 210 were images with benign masses. Our dataset BUSI is available online[1] for studies.



Fig. 1. Samples of breast US images from both datasets where the first row contains images from dataset BUSI and the second row contains images from dataset B.

The other dataset is referred to as Dataset B [14]. It was collected in 2012 from the UDIAT Diagnostic Center of the

Parc Tauli Corporation, Sabadell (Spain). It has 163 images from different females. The average image size of the dataset is 760 x 570 pixels. The number of images in the dataset is 163 images where 53 images were with malignant masses and 110 images were with benign tumors. It was created for lesion detection not for classification while our study uses it for lesion classification.

## IV. METHODOLOGY

Our methodology is divided into two parts. In the first part, we discuss data augmentation using GAN and traditional augmentation. While the second part discusses classification techniques which are performed by using deep learning approaches Convolutional Neural Network (CNN) and a Transfer Learning (TL) on BUSI dataset, dataset B and merged datasets (BUSI+B). The whole model architecture is shown in Fig. 2. It is important to highlight that the classification algorithms were performed on four forms of data samples as follows: (1) without augmentation which means the real images(the blue dash line and arrows). (2) with traditional augmentation. (3) using DAGAN. (4) using traditional augmentation and DAGAN ( the orange box and arrows).

### A. Data Augmentation Generative Adversarial Networks (DA-GAN)

The second goal of this study was to produce realistic images for each of the classes on-demand while the first goal is to enhance the classification accuracy using deep learning approaches. Each GAN is composed of two networks: a generator and a discriminator, playing against each other in a two-player game. These models have proven to be capable of creating realistic images. To achieve this, a framework was performed where a single GAN was trained on each of the classes. A GAN architecture of sufficient capacity to understand and model the underlying distributions of each of the classes had to be selected. A GAN that satisfies the above goal should, after training, be able to produce realistic images of the class it was trained upon.

Furthermore, GAN [8] is formed of two networks which are the generator and the discriminator. The generator accepts a noise vector as input and produces fake data, which are then fed, along with real ones, to the discriminator, whose goal is to distinguish which distribution the samples were produced from. Conversely, the generator's goal is to learn the real distribution without witnessing it, in order to make its output indistinguishable from real samples. Both networks are trained simultaneously and adversarially until an equilibrium is reached. In order to combat instability issues during training, the Earth Mover's or Wasserstein distance was used, partially because it leads to convergence for a much broader set of distributions, but mostly because its value is directly correlated to the quality of the generated data [40]. The discriminator was initially achieved by clipping its weights by an arbitrary value Wasserstein GAN (WGAN) [40]. It was later shown that this technique led to sub-optimal behavior, which could be ameliorated with the inclusion of a gradient penalty term to the discriminator's loss function calculated on a random interpolation point between the real and the fake samples [41]. The resulting architecture WGAN gradient penalty (WGAN-GP) [41] is the one utilized in our study.

Fig. 2. The proposed methods for breast US image classification and data augmentation techniques.

*1) Generator:* An architecture with 11 layers was selected as the generator of the network. The architecture is depicted in Fig. 3A. The generator input is a vector of 128 random values in the range of (0,1). It is sampled from a uniformed distribution. A Fully Connected (FC) layer followed the input layer. The subsequent layers are regular 2D convolutions (Conv) and 2D transposed convolutions (Conv trans up), sometimes referred to as "deconvolution" layers. A 5x5 sized kernel and "same" padding were selected for both types of layers, while a stride of 2 was selected for the transposed convolutions. This performs in the doubling of the spatial dimensions of its input. A "Leaky ReLU" function activated all layers apart from the last layer. The final layer has a tangent hyperbolic (tanh) activation function because its output needs to be bound in order to be able to output an image. A tanh function was preferred over a sigmoid function because it is centered around 0, which helps during training [42]. Finally, after five alternations of convolution and transposed convolution layers (each of which doubles the size of its input), an image with a resolution of (192x160) and 1 channel is produced.



A) Generator          B) Discriminator

Fig. 3. DAGAN architecture: A) Generator structure and B) Discriminator structure

*2) Discriminator:* The discriminator is a usual CNN architecture intended towards binary classification. The one used in the present study consists of 11 layers can be seen in Fig. 3B. The input to the discriminator is a single-channel 192x160 image. This image is then passed five times through alternating layers of convolutions with a stride of 1 and 2 respectively; the latter is used for sub-sampling as there are no pooling layers present in the architecture. The last two layers are FC ones. All layers in the network are activated by a "Leaky ReLU" , besides the last one which has no activation function.

### B. Traditional Augmentation Techniques

Due to the nature of our datasets, we could only implement a limited range of visual transformations. In particular, we applied a horizontal flip, brightness, scaling and zooming. The number of augmented images that were obtained from traditional augmentation would increase by a factor of 2 for each augmented method.

### C. Convolutional Neural Network

Based on the Deep learning definition, it is a representation learning approach [43] that will automatically detect features satisfying a special task from the data. The feature extractors are task-specific, in that they are not fixed to a set of specific rules every time [44]. Each network contains multiple layers that lead to hierarchical features used in the learning process [43], [45].

CNNs [46] are a valuable technique in image analysis, particularly in recognition, detection or classification of faces [47], text [45], biological images [48] and human bodies [49]. For these reasons, we study the performance of deep learning in breast US tumor classification.

CNNs consist of convolutional layers and pooling layers [46], where the role of the former is to extract local features from a set of learnable filters and the role of the latter is to merge neighboring patterns, reducing the spatial size of the previous representation and adding spatial invariance to translation [43]. CNNs are hierarchical neural networks and their accuracy is based on the design of the layers and training models [50].

Some common CNNs are available which are AlexNet [16], LeNet [45] and GoogleNet [51]. We studied the use of two types of deep learning models for breast classification: AlexNet [16] and a transfer learning approach using Convolutional Networks [52].

*1) CNN-AlexNet:* As the ultrasound breast images in the datasets are gray-scale and the size of the breast tumor or micro-calcification is relatively small, AlexNet [16] was chosen as a suitable architecture to solve the classification problem of multi-classes. The training and validation images are input of the model containing all classes in the datasets. We split all datasets to 70%,15%, and 15% for training, validation, and testing, respectively. The AlexNet architecture is simple and was primarily built for digit classification [45]. Breast tumors include related gradients that can be presented through CNNs. The overall architecture is shown in Fig. 4, with the inputs consisting of images of breast tumors and normal tissue. The inputs are fed into the first convolution layer and max-pooling layer, which is repeated once and finalized with two fully connected layers. The final number of outputs are 2 neurons or 3 neurons, which are the activations generated for the two or three classes: (benign and malignant) or (normal, benign and malignant), respectively. The final part of the CNN is the output of class probabilities to measure how close the final fully connected parameters are with respect to the labels of the training and validation data. The loss was calculated using multinomial logistic loss with a softmax classifier. The output of our network is a prediction of whether the image is a tumor or healthy breast tissue. It is formed by two fully connected layers with the softmax function defined as

$$f_j(z) = \frac{e^{z_j}}{\sum_k e^{z_k}} \quad (1)$$

where $f_j$ is the *j*-th element of the vector of class scores $f$ and $z$ is a vector of random real-valued scores that are flattened to a vector of values between zero and one that sum to one. The loss function is defined so that having good predictions during training is equivalent to having a small loss. A Rectified Linear Unit (ReLU) layer is included at the first fully-connected layer. This element-wise operation is calculated and defined as

$$f(x) = max(0, x) \quad (2)$$

where the function $f$ thresholds the activations at zero.

*2) Transfer Learning:* Transfer Learning (TL) [53], [54] is a method where a CNN is trained to learn features for a broad domain after which the classification function is changed to optimize the network to learn features of a more specific domain. Under this setting, the features and the network parameters are transferred from the broad domain to the specific one. Furthermore, Transfer Learning (TL) is a method that provides a system to apply the knowledge learned of prior tasks to a new task domain that is somehow related to the prior domain. Our proposed transfer learning approach is based on VGG16 [55], ResNet [56], Inception [57], and NASNet [58]. These networks were primarily utilized for the classification of more than one thousand various objects of classes on the ImageNet dataset [16]. The default image sizes for TL models are shown in Table I.

### D. Implementation

*1) Preprocessing:* In this subsection, we focus on the preparation of the datasets and image augmentation. Additional

TABLE I. THE DEFAULT INPUT SIZE FOR TRANSFER LEARNING MODELS

| Models | Input Shape |
|---|---|
| VGG16 | 224x224 |
| ResNet | 224x224 |
| Inception | 299x299 |
| NASNetLarge | 331x331 |

preprocessing steps were taken to facilitate model training, such as resizing them to 192x160 with Lanczos interpolation. In addition, we randomly divided the dataset into training, validation and test sets, keeping intact the sequence of each image so that every image appears in only one of the aforementioned sets. We should note here that in our initial experiments we randomly shuffled and split all images without preserving each image sequences; this allowed the models to identify key features in each subject's morphology and achieve a perfect score on the test set (i.e., for each test set image, the model had been trained on another from the same image). Because of this, the study of the models' generalization on new, unseen patients, which is a necessary requirement for all medical applications, became infeasible. We train both datasets (BUSI and B) using DAGAN model. Trained models are saved and reused in generating new images. DAGAN runs in 700 epochs for each class (normal, benign and malignant). Samples of real image for datasets and augmented images are shown in Fig. 5. We generate 5000 images for each class using DAGAN model. All the images are added to our datasets.

*2) Classification Methods:* The proposed CNN approach in this paper is AlexNet model [16]. The breast US images are in grayscale. The datasets were split into 70%, 15%, and 15% for training, validation, and testing, respectively. The validation set (15%) is used for hyper-parameter tuning and early stopping. The network is trained by using Adam optimizer - with a learning rate of 0.0001. It uses 60 epochs (early stopping) with 0.30 of dropout rate. We used a stride of one and two pixels in max-pooling. To obtain the best performance for the state-of-the-art classification methods on the datasets, we use regularization techniques such as normalization and dropout.

For transfer learning, we used four pretrained models which are VGG16 [55], ResNet [56], Inception [57], and NASNet [58]. An Adam optimizer is used with learning rate 0.001. The number of epochs was 10 epochs. The output layer of TL models is altered and we train our data in it. The softmax activation function is utilized in TL experiments.

In order to measure the effectiveness of the proposed methodology, the following experiment was devised: Firstly, a Deep neural network architecture was selected, which is capable of achieving satisfactory performance on classifying the three classes (i.e. normal, benign and malignant). Secondly, DAGAN and traditional augmentation are used to enhance the performance of classification algorithms by generating more data-samples.

Our methodology is summarized in the following points:

1) First, we have two datasets BUSI and dataset B and a third one that was created by merging the two datasets (BUSI+B).
2) We perform two types of data augmentation to generate more data samples, the first type is traditional augmentation and the second type is DAGAN.

Fig. 4.   CNN architecture of AlexNet Model.



Fig. 5.   Samples of real images from dataset BUSI are in the first row and augmented images using DAGAN are in the second row.

3) Two deep learning classification approaches were used, CNN (AlexNet) and TL (VGG16, ResNet, Inception, and NASNet).

4) It is important to highlight that in our experiments we perform the classification algorithms on four forms of data samples as follows: (1) without augmentation which means the real images. (2) with traditional augmentation. (3) using DAGAN. (4) using traditional augmentation and DAGAN. As a result of this, the total number of 60 classification codes have been implemented (see Table II and Fig. 6)

### E. Implementation Environment

Our classification experiments are performed on Windows 10 operating system using Keras API library[2] version 2.0.1 (on top of TensorFlow[3]) using Python (version 3.5). In this study, training and classification are performed on Intel (R) Core (TM) i73630QM CPU @2.40MGz and GPU NVIDIA Quadro K2000M With 8GB of shared GPU memory, and 16 GB RAM.

Furthermore, DAGANs are performed in a powerful server which uses Ubuntu 18.04 operating system using as mentioned above, Tensorflow, Keras, and Python. The server specification

is Intel Xeon(R) CPU E5-26200 @2.00GHz×12, llvmpipe (LLVM 7.0, 256bits), and 50 GB RAM.

## V. RESULTS AND DISCUSSION

There are many parameters that affect the results in deep learning when used in medical images such as the type of algorithms, hyperparameters, and size of the dataset. We considered all of these parameters in our experiments.

Dataset BUSI was obtained from a modern US system, which offers new challenges for the current techniques in tumor classification. These US systems obtain high-resolution images that may cover other structures such as air in the lungs, ribs or pectoral muscle, making the tumor classification more difficult. Dataset B was collected from an older US system. Images are usually of a lower resolution. However, these differences did not affect our experiments.

The performance accuracies of our experiments are shown in Table II. Regarding our experiments, we found the following:

1) In all of our experiments, we found that increasing the number of data samples using data augmentation and datasets merging, significantly improve the classification accuracies. This is obvious in the results Table II and Fig. 6. Note that the classification results obtained from DAGAN outperform the traditional argumentation. While the results were the best when we combine DAGAN and traditional argumentation.

2) When we performed the experiments on datasets without data augmentation, they produce low accuracies (even if we combined the two datasets (BUSI+B)). This is due to the shortage of data.

3) We figured out that traditional augmentation is not very effective in our work due to the nature of medical images. In addition, medical images are not like natural images that are used in object classification. There are limited traditional augmentation techniques that can be used in medical images.

These results showed that the supervised deep learning methods were data-driven and the performance increased with more training dataset. We can confirm that the transfer learning approach achieved the best accuracy when trained with data augmentation through the use of DAGANs and traditional

---

[2]https://keras.io/
[3]https://www.tensorflow.org

TABLE II. COMPARISON OF THE ACCURACY OF DIFFERENT METHODS WHEN TESTING ON SINGLE AND COMBINED DATASET. THE BEST RESULTS IS INDICATED IN BOLD.

| Dataset | Method | Sub-Method | Without Augmentation | Traditional Augmentation | DAGAN Augmentation | Both Traditional and DAGAN augmentation |
|---|---|---|---|---|---|---|
| Dataset BUSI | | CNN-AlexNet | 58% | 62% | 73% | 78% |
| | TL | VGG16 | 70% | 74% | 84% | 88% |
| | | Inception | 68% | 73% | 82% | 85% |
| | | ResNet | 79% | 82% | 89% | 93% |
| | | NASNet | 83% | 85% | 91% | **94%** |
| Dataset B | | CNN-AlexNet | over-fitting | 56% | 75% | 80% |
| | TL | VGG16 | 68% | 72% | 80% | 82% |
| | | Inception | 65% | 70% | 77% | 80% |
| | | ResNet | 75% | 79% | 86% | 90% |
| | | NASNet | 79% | 82% | 90% | **92%** |
| Datasets (BUSI+B) | | CNN-AlexNet | 60% | 65% | 82% | 84% |
| | TL | VGG16 | 72% | 75% | 86% | 88% |
| | | Inception | 70% | 73% | 84% | 87% |
| | | ResNet | 76% | 79% | 88% | 92% |
| | | NASNet | 84% | 88% | 96% | **99%** |



Fig. 6. The chart illustrates the performance accuracies in all the performed methods with three datasets in (without: no data augmentation, with TA: with traditional augmentation, with DAGAN: using generated images, with DAGAN and TA: used traditional augmentation and generated images).

augmentation in combined dataset (BUSI+B). The final result is 99% (when training on TL NASNet pretrained model).

The use of a powerful generative model for producing images (e.g., DAGAN) has many advantages over traditional augmentation schemes. The most important advantage is the quality of the produced images and the capability of generalizing beyond the limits of the original dataset to produce new patterns. The proposed technique is especially useful in low-variance datasets where the images follow a very strict format. We would like to point out that other studies reached 92% accuracy in classification methods using their own datasets while we reached 99% using our datasets.

On the other hand, there are some limitations in our work which are:

- The training process is time consuming and requires high computer resources.

- There is not a sufficient number of real images that have been collected to avoid the classification errors

in the augmented images.

- We can not synthesize high-resolution images using DAGAN.

## VI. CONCLUSIONS AND FUTURE WORKS

This paper investigated the use of two deep learning classification approaches (particularly CNN (AlexNet) and Transfer Learning approaches). Two datasets were used which are our Dataset BUSI and Dataset B. Furthermore, we combined them obtaining a third one which is dataset (BUSI+B). We used a novel methodology for data augmentation with the use of GAN. It involves training a GAN for each of the classes of the original datasets and then using it to produce a number of synthetic images. All models were trained on breast US images datasets to classify cancerous and non-cancerous images.

To study the impact of this augmentation strategy for classification methods, four experiments were conducted. Firstly, CNN and transfer learning models were trained on all datasets

on a form of baseline. Secondly, the same models were trained with traditional and thirdly by the proposed GAN augmentation techniques. Fourthly, by both forms of augmentation(traditional and DAGAN).

The performances were evaluated on the three datasets (BUSI, B, and BUSI+B). Amongst the various methodologies presented in this paper, the transfer learning NASNet achieved the best results (99%) in Dataset (BUSI+B) when it is used with DAGAN and traditional augmentation. Deep learning methods are adaptable to the specific characteristics of any dataset since these are machine-learning based and particular models are constructed for each dataset. Experiments confirm that augmentation through GANs outperforms traditional augmentation methods when used with CNN and transfer learning.

Finally, the models trained with the proposed methods using GAN augmentation methodology outperform the ones with a traditional one by a large margin. In fact, because of the nature of the images, the traditional techniques gave no enhancement over the baseline experiments. The final experiments, which combined both forms of augmentation exceeded the rest, pointing that while traditional augmentation could not function on its own, it performs well when it is combined with GAN augmentation.

In the future, we believe that deep learning approaches could be adjusted to other medical imaging techniques such as 3D ultrasound or other modalities. Mass classification is the initial step of a CAD system. Hence, in our future work we plan to do breast ultrasound lesion detection and segmentation, and evaluate the performance of the complete CAD system. Because of the improved results of our experiments using DAGAN, multiple future research areas could be spawned. We are planning to experiment with different structures for further developments on the data quality, either within the WGAN-GP by utilizing a more robust discriminator, or by using a newer, more modern framework that leads to enhanced experimental performance, such as the Progressive Growing GANs [59] or the Auxiliary Classifier GANs [60].

## REFERENCES

[1] E. Hall, "Radiobiology for the radiologist, radiation research, vol. 116, no. 1, 1988."

[2] H.-D. Cheng, J. Shan, W. Ju, Y. Guo, and L. Zhang, "Automated breast cancer detection and classification using ultrasound images: A survey," *Pattern recognition*, vol. 43, no. 1, pp. 299–317, 2010.

[3] O. Akin, S. B. Brennan, D. D. Dershaw, M. S. Ginsberg, M. J. Gollub, H. Schöder, D. M. Panicek, and H. Hricak, "Advances in oncologic imaging: update on 5 common cancers," *CA: a cancer journal for clinicians*, vol. 62, no. 6, pp. 364–393, 2012.

[4] A. T. Stavros, D. Thickman, C. L. Rapp, M. A. Dennis, S. H. Parker, and G. A. Sisney, "Solid breast nodules: use of sonography to distinguish between benign and malignant lesions." *Radiology*, vol. 196, no. 1, pp. 123–134, 1995.

[5] M. H. Yap, E. Edirisinghe, and H. Bez, "Processed images in human perception: A case study in ultrasound breast imaging," *European journal of radiology*, vol. 73, no. 3, pp. 682–687, 2010.

[6] K. Drukker, N. P. Gruszauskas, C. A. Sennett, and M. L. Giger, "Breast us computer-aided diagnosis workstation: performance with a large clinical diagnostic population," *Radiology*, vol. 248, no. 2, pp. 392–397, 2008.

[7] J. Kukačka, V. Golkov, and D. Cremers, "Regularization for deep learning: A taxonomy," *arXiv preprint arXiv:1710.10686*, 2017.

[8] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in neural information processing systems*, 2014, pp. 2672–2680.

[9] D. C. Cireşan, U. Meier, L. M. Gambardella, and J. Schmidhuber, "Deep, big, simple neural nets for handwritten digit recognition," *Neural computation*, vol. 22, no. 12, pp. 3207–3220, 2010.

[10] C. N. Vasconcelos and B. N. Vasconcelos, "Increasing deep learning melanoma classification by classical and expert knowledge based image transforms," *CoRR, abs/1702.07025*, vol. 1, 2017.

[11] R. Wu, S. Yan, Y. Shan, Q. Dang, and G. Sun, "Deep image: Scaling up image recognition," *arXiv preprint arXiv:1501.02876*, 2015.

[12] Z. Hussain, F. Gimenez, D. Yi, and D. Rubin, "Differential data augmentation techniques for medical imaging classification tasks," in *AMIA Annual Symposium Proceedings*, vol. 2017. American Medical Informatics Association, 2017, p. 979.

[13] A. Antoniou, A. Storkey, and H. Edwards, "Data augmentation generative adversarial networks," *arXiv preprint arXiv:1711.04340*, 2017.

[14] M. H. Yap, G. Pons, J. Martí, S. Ganau, M. Sentís, R. Zwiggelaar, A. K. Davison, and R. Martí, "Automated breast ultrasound lesions detection using convolutional neural networks," *IEEE journal of biomedical and health informatics*, vol. 22, no. 4, pp. 1218–1226, 2018.

[15] B. Huynh, K. Drukker, and M. Giger, "Mo-de-207b-06: Computer-aided diagnosis of breast ultrasound images using transfer learning from deep convolutional neural networks," *Medical physics*, vol. 43, no. 6Part30, pp. 3705–3705, 2016.

[16] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.

[17] J. Shi, S. Zhou, X. Liu, Q. Zhang, M. Lu, and T. Wang, "Stacked deep polynomial network based representation learning for tumor classification with small ultrasound image dataset," *Neurocomputing*, vol. 194, pp. 87–94, 2016.

[18] R. Livni, S. Shalev-Shwartz, and O. Shamir, "An algorithm for training polynomial networks," *arXiv preprint arXiv:1304.7045*, 2013.

[19] S. Zhou, J. Shi, J. Zhu, Y. Cai, and R. Wang, "Shearlet-based texture feature extraction for classification of breast tumor in ultrasound image," *Biomedical Signal Processing and Control*, vol. 8, no. 6, pp. 688–696, 2013.

[20] Y.-X. Wang, R. Girshick, M. Hebert, and B. Hariharan, "Low-shot learning from imaginary data," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7278–7286.

[21] L. Perez and J. Wang, "The effectiveness of data augmentation in image classification using deep learning," *arXiv preprint arXiv:1712.04621*, 2017.

[22] T. Neff, C. Payer, D. Stern, and M. Urschler, "Generative adversarial network based synthesis for supervised medical image segmentation," in *Proc. OAGM and ARW Joint Workshop*, 2017.

[23] A. F. Frangi, S. A. Tsaftaris, and J. L. Prince, "Simulation and synthesis in medical imaging," *IEEE transactions on medical imaging*, vol. 37, no. 3, p. 673, 2018.

[24] X. Yi, E. Walia, and P. Babyn, "Generative adversarial network in medical imaging: A review," *arXiv preprint arXiv:1809.07294*, 2018.

[25] A. BenTaieb and G. Hamarneh, "Adversarial stain transfer for histopathology image analysis," *IEEE transactions on medical imaging*, vol. 37, no. 3, pp. 792–802, 2018.

[26] S. Kazeminia, C. Baur, A. Kuijper, B. van Ginneken, N. Navab, S. Albarqouni, and A. Mukhopadhyay, "Gans for medical image analysis," *arXiv preprint arXiv:1809.06222*, 2018.

[27] H.-C. Shin, N. A. Tenenholtz, J. K. Rogers, C. G. Schwarz, M. L. Senjem, J. L. Gunter, K. P. Andriole, and M. Michalski, "Medical image synthesis for data augmentation and anonymization using generative adversarial networks," in *International Workshop on Simulation and Synthesis in Medical Imaging*. Springer, 2018, pp. 1–11.

[28] W. Dai, X. Liang, H. Zhang, E. Xing, and J. Doyle, "Structure correcting adversarial network for chest x-rays organ segmentation," Sep. 27 2018, uS Patent App. 15/925,998.

[29] Y. Xue, T. Xu, H. Zhang, L. R. Long, and X. Huang, "Segan: Adversar-

ial network with multi-scale l 1 loss for medical image segmentation," *Neuroinformatics*, vol. 16, no. 3-4, pp. 383–392, 2018.

[30] S. U. Dar, M. Yurt, L. Karacan, A. Erdem, E. Erdem, and T. Çukur, "Image synthesis in multi-contrast mri with conditional generative adversarial networks," *IEEE transactions on medical imaging*, 2019.

[31] M. Frid-Adar, I. Diamant, E. Klang, M. Amitai, J. Goldberger, and H. Greenspan, "Gan-based synthetic medical image augmentation for increased cnn performance in liver lesion classification," *Neurocomputing*, vol. 321, pp. 321–331, 2018.

[32] P. Costa, A. Galdran, M. I. Meyer, M. Niemeijer, M. Abràmoff, A. M. Mendonça, and A. Campilho, "End-to-end adversarial retinal image synthesis," *IEEE transactions on medical imaging*, vol. 37, no. 3, pp. 781–791, 2018.

[33] D. Ciresan, A. Giusti, L. M. Gambardella, and J. Schmidhuber, "Deep neural networks segment neuronal membranes in electron microscopy images," in *Advances in neural information processing systems*, 2012, pp. 2843–2851.

[34] T. Kooi, G. Litjens, B. Van Ginneken, A. Gubern-Mérida, C. I. Sánchez, R. Mann, A. den Heeten, and N. Karssemeijer, "Large scale deep learning for computer aided detection of mammographic lesions," *Medical image analysis*, vol. 35, pp. 303–312, 2017.

[35] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.

[36] H. Ravishankar, P. Sudhakar, R. Venkataramani, S. Thiruvenkadam, P. Annangi, N. Babu, and V. Vaidya, "Understanding the mechanisms of deep transfer learning for medical images," *arXiv preprint arXiv:1704.06040*, 2017.

[37] N. Dhungel, G. Carneiro, and A. P. Bradley, "Deep learning and structured prediction for the segmentation of mass in mammograms," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2015, pp. 605–612.

[38] J.-J. Mordang, T. Janssen, A. Bria, T. Kooi, A. Gubern-Mérida, and N. Karssemeijer, "Automatic microcalcification detection in multi-vendor mammography using convolutional neural networks," in *International Workshop on Breast Imaging*. Springer, 2016, pp. 35–42.

[39] C. K. Ahn, C. Heo, H. Jin, and J. H. Kim, "A novel deep learning-based approach to high accuracy breast density estimation in digital mammography," in *Medical Imaging 2017: Computer-Aided Diagnosis*, vol. 10134. International Society for Optics and Photonics, 2017, p. 101342O.

[40] M. Arjovsky, S. Chintala, and L. Bottou, "Wasserstein gan," *arXiv preprint arXiv:1701.07875*, 2017.

[41] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. C. Courville, "Improved training of wasserstein gans," in *Advances in Neural Information Processing Systems*, 2017, pp. 5767–5777.

[42] Y. LeCun, L. Bottou, G. B. Orr, K.-R. Müller *et al.*, "Neural networks: Tricks of the trade," *Springer Lecture Notes in Computer Sciences*, vol. 1524, no. 5-50, p. 6, 1998.

[43] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning. nature 521 (7553): 436," *Google Scholar*, 2015.

[44] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, "Caffe: Convolutional architecture for fast feature embedding," in *Proceedings of the 22nd ACM international conference on Multimedia*. ACM, 2014, pp. 675–678.

[45] Y. LeCun, L. Bottou, Y. Bengio, P. Haffner *et al.*, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.

[46] Y. LeCun, Y. Bengio *et al.*, "Convolutional networks for images, speech, and time series," *The handbook of brain theory and neural networks*, vol. 3361, no. 10, p. 1995, 1995.

[47] Y. Taigman and M. Yang, "Marc'aurelio ranzato, and lior wolf. deepface: Closing the gap to human-level performance in face verification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 1701–1708.

[48] F. Ning, D. Delhomme, Y. LeCun, F. Piano, L. Bottou, and P. E. Barbano, "Toward automatic phenotyping of developing embryos from videos," *IEEE Transactions on Image Processing*, vol. 14, pp. 1360–1371, 2005.

[49] P. Sermanet, K. Kavukcuoglu, S. Chintala, and Y. LeCun, "Pedestrian detection with unsupervised multi-stage feature learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 3626–3633.

[50] D. C. Ciresan, U. Meier, J. Masci, L. M. Gambardella, and J. Schmidhuber, "Flexible, high performance convolutional neural networks for image classification," in *Twenty-Second International Joint Conference on Artificial Intelligence*, 2011.

[51] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1–9.

[52] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3431–3440.

[53] R. Caruana, "Multitask learning," *Machine learning*, vol. 28, no. 1, pp. 41–75, 1997.

[54] S. Thrun, "Is learning the n-th thing any easier than learning the first?" in *NIPS*, 1995.

[55] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[56] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[57] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2818–2826.

[58] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le, "Learning transferable architectures for scalable image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 8697–8710.

[59] T. Karras, T. Aila, S. Laine, and J. Lehtinen, "Progressive growing of gans for improved quality, stability, and variation," *arXiv preprint arXiv:1710.10196*, 2017.

[60] A. Odena, C. Olah, and J. Shlens, "Conditional image synthesis with auxiliary classifier gans," in *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR. org, 2017, pp. 2642–2651.

# Feature Fusion for Negation Scope Detection in Sentiment Analysis: Comprehensive Analysis over Social Media

Nikhil Kumar Singh[1], Deepak Singh Tomar[2]
Department of Computer Science,
Maulana Azad National Institute of Technology,
Bhopal, India

*Abstract*—**Negation control for sentiment analysis is essential and effective decision support system. Negation control include identification of negation cues, scope of negation and their influence within it. Negation can either shift or change the polarity score of opinionated word. This paper present a framework for feature fusion of text feature extraction, negation cue and scope detection technique for enhancing the performance of recent sentiment classifier for negation control. Explore text feature POS, BOW and HT with negation cue and scope detection techniques for classification technique over social media data set. This paper has included the evaluation of sentiment classification (Support vector machine, Navies Bayes, Linear Regression and Random Forest) and Nine feature fusion over presented prepossessing framework. This paper yield interesting result about collective response of feature fusion for negation scope detection and classification technique. Feature Fusion vector significantly increase the polarity classification accuracy of sentiment classification technique. POS with Grammatical dependency tree can detect negation with better accuracy as compared to other feature fusion.**

*Keywords*—*Sentiment analysis; feature fusion; negation cues; scope detection; conjunction analysis; punctuation mark; grammatical dependency tree; Navies Bayes; linear regression; random forest; SVM*

## I. Introduction

Sentiment analysis (SA) is the computational analysis of the opinion, attitudes, emotions of speaker/writer towards some topic and identification of non- trivial, subjective information from text repository. Before the term sentiment analysis came into existence [1], this area was recognized as opinion mining, point of view and subjectivity. At the present time, SA is speedily growing field due to the rise of online message spreading platform such as blogs, social media and commercial website. On regular basis billions of people share their experiences, knowledge and views on latest trend of politics, economics and other global- critical issue. In current time Sentiment Analysis, subjectivity and Opinion mining enthralled significant interest from both the research community and Marketing Agency [2]. The main purpose of sentiment analysis is to rank the opinion according to its level of positive, negative or neutral polarity [3]. Sentiment analysis have many applications, ranging from product analysis [4]to improving sales and marketing strategies, predicting stock market fluctuations [5], identifying changes ideological in political issues [6] , in the prediction of film critics [7] and in Electronic Government. -regulation [8],

that is to say, the opinion of the citizens on a law before its approval. Although there has been a lot of work done in the area of sentiment analysis, there are still open challenges related to SA's multilingual strategy, classifying the sentence with slangs, symbols, misspelled words and expressions. idioms, SA sarcastic sentences and handle negation and identify polarity. mark in negative feelings [3]. Here in this paper we have summarized the effect of negation cues over sentiment analysis and introduced a comparative analysis of recent text feature extraction, negation cues and scope detection technique. This paper present a framework securitizing and preprocessed social media data set and formulate the supervised classification technique with feature fusion for negative sentiment analysis. The rest of the paper is organized as follows: Section 2 presents over view of Negative sentiment analysis; Section 3 covers related work on negation handle mechanism for sentiment analysis and polarity detection over social media data set. Section 4 present a framework for securitizing and preprocessed social media data set and subsection 4(A-C) explain how social media data are processed, step for prepossessing, negation cue and scope technique for efficient SA and experimental Contents for performance evaluation respectively. Section 5 describe the experimental setup for comparative evaluation of different scope detection technique with classification approach for sentiment analysis over social media and finally, Section 6 concludes the paper and outlines the founding and future work.

## II. Negation Sentiment Analysis

Negation can be defined as a linguistic event. It acts as polarity influence which can effect the meaning or the semantic of sentence for e.g the polarity of sentence changes from positive to negative eventually which can swing the polarity strength. To overcome this, necessary action for negation in SA are required. Author in [9] state that negation is a complex phenomenon that studied under different disciplines. In NLP, negation is considered as operator and scope is a principle feature of operators, i.e. negation influence the meaning of other phase of the sentence within their scope. Negations can not only change the meaning of single words or phase of words but also reduce the polarities of opinionated word. For example consider following sentences $S_1$, $S_2$ and $S_3$.

**Sentence** ($S_1$):*This Sunscreen Lotion is not costly but it*

*suits me.*

**Sentence ($S_2$):** *The product doesn't have nice packaging but really effective.*

**Sentence ($S_3$):** *This sunscreen is less relevant for fairer skin.*

Where in Sentence $S_1$, scope of negation 'Not' is only limited to the next word after negation i.e. 'costly'. Where negation only invert the meaning of word "Suits". Whereas in Sentence S2, Scope of negation "not" is till the end of sentence. On other hand in sentence S3 uses diminisher "Less" to reduce the polarities of opinionated words instead of completely reversing the polarities. Method to handle negation in sentiment analysis is depend upon type of negative linguistic patterns and class negative word used in respective negative sentence as shown in Fig. 1 and Table I. Table I contain list of negations which serve as an indicator of the presence of a negation with different linguistic patterns.

Depending upon assertion linguistic patterns, negation in negative sentences may be occur explicitly (with explicit clues such as not, no etc.) and implicitly (with implicit clues such as scarcely, hardly, few, seldom, little, only, etc.). For expressing the negative opinion, if negation encoded opinionated word has been used then its implicit negation whereas if standard negation cues are used with opinionated word then it is explicit negation. The list of explicit and implicit negation cue are listed in Table I. For example consider following sentences $S_4$, $S_5$ and $S_6$.

**Sentence ($S_4$):** *This music system is not good.*

**Sentence ($S_5$):** *My personal experience to use this music system is horrible.*

**Sentence ($S_6$):** *Sound system of this music system is superb, I'm suffering from headache after enjoying the song!!!*

**Sentence ($S_7$):** *This music system is irrelevant for oldies!!!*

For instance Sentence $S_4$ have explicit negative sentiment about the music system whereas sentence $S_5$ use "horrible" as opinionated that encode negative sentiment about music system. On other hand sentence $S_6$ use irony to reflect its negative sentiment about the respective product. Whereas at structural level negative sentence may be appear with morphological, syntactic, contrast, compound and non-negative negations. In Morphological negation, negative meaning is carried out by modifying opinionated word either by prefix (e.g. ir-, non-, un- etc.) or suffix (e.g. -less). Whereas in Syntactic negations, explicit negation cues are used to revise the polarity of a single opinionated word or a sequence of words. For instance sentence $S_7$ use morphological negation to show its negative concern about the cell phone. Whereas $S_1$, $S_2$, $S_4$ and $S_5$ syntactic negative sentence. In contrast negation, negative expression show contrast or manage opposition between opinionated terms. While compound negation express comparison or inequality between opinionated term. Whereas in non-negative negation that's used for interrogative and conditional sentences, negative cues and opinionated term may not contain any opinion or sentiment. For instance, sentence $S_8$, $S_9$ and $S_{10}$ shows contrast, compound and non-negative negation respectively.

**Contrast negation ($S_8$):-** *I brought this cell phone not for camera resolution but for its MP3.*

**Compound negation ($S_9$):-** *Touchscreen of cell phone is not better than other..*

**Non-negative Negation ($S_{10}$):-** *Is Sound quality of this cell phone is not good?*

Intensifier and diminisher phase of word use as valance shifter in negation. Valance shifter usually degrade or upgrade the polarity strength instead of inverting the polarity of opinionated word. For instance, sentence $S_{11}$ and $S_{12}$ shows intensifier and diminisher based valance shifter in negation. Where the term "Very much" in sentence $S_{11}$ degrade the negative polarity orientated by the term "not relevant" while the term "less" used in $S_{12}$ shift positive polarity of front camera towards little bit negative.

**Negative Intensifier ($S_{11}$):-** *This Sunscreen is not very much relevant for me.*

**Negative Diminishes ($S_{12}$):-** *Effect of this Sunscreen is less relevant for beach outing.*

## III. RELATED WORK

Handle Negation in sentiment analysis required to identification negation term as cue detection and recognize its linguistic influence as scope detection. Recently researcher focus to identify negative cues grammatical structure for framing supervised syntactic rule through for training purpose [10], [11], [12], [13]

Ghiassi et.al. [14] applied supervised rules for polarity score calculation and tagged opinionated term with six different polarity level i.e. "XP" (extremely positive), "VP" (very positive), "SP" (somewhat positive), "SN" (somewhat negative), "VN" (very negative) or "XN" (extremely negative) by using information gain feature extraction technique. Whereas Apple et al. [15] present fuzzy set theory based probabilistic classifier for categorizing polarity intensity up to five level from Mild to most intensive as Poorly slight, Moderate, very and Most intensive sentiment word.

Garcia et al. [16] present probabilistic classifier to highlight the negativity, Korkontzelos et al. [17] use part of speech (POS) to evaluate grammatical dependency among negation cue and opinionated word in medical area. Diamantini et al. [18] use depth-first search (DFS) strategy for building grammatical dependency tree to identify of negation cues. Tian Kang et al. [19] use Conditional Random Fields (CRF) for 'BIO' tagging to represent the boundaries of negation cues. Prollochs et al. [20] use manually labeled dataset for predicting negation cue and it scopes by reinforcement learning and machine learning technique.

Polarity shift via negative cues affect sentiment analysis performance. Recent research has been focus over arithmetical techniques to discriminate explicit and implicit polarity shifts valuation. Tellez et al. [21]use rule-based method to spot polarity shifts in explicit negations and contrasts. Ghiassi et al. [14] use BOW to handle valence shifter such as intensifiers, diminishers and sarcasm.

Jimenez-Zafra et al. [22] use SFU review -NEG corpus for the supervised polarity classification system. AL-Sharuee et al. [23] handling intensifiers and negation using SentiWordNet and use antonym dictionary to replace adjectives and adverbs

Fig. 1.   Classification of Negative Cues

TABLE I.   NEGATION CUE WITH DIFFERENT LINGUISTIC PATTERNS

| Negation Class | Negative Word |
| --- | --- |
| Explicit Negation | no, not, rather, never, none, nobody, nothing, neither, nor, nowhere (in all tense class). |
| Implicit Negation | scarcely, hardly, few, seldom, little, forget, fail, doubt, deny and etc. |
| Diminisher | hardly, few, , little, less . |
| Syntactic | no, not, rather, never, none, nobody, nothing, neither, nor, nowhere (in all tense class). |
| Morphological | Prefixes: de-, dis-, il-, im-, in-, ir-, mis-, non-, un- , Suffix: -less |
| Intensifiers | Absolute, badly, biggest, epic, specially, eternally, exceptionally, extremely, freak in, fuckin, hella, huge, incredibly, major, massive, mighty, most, deadly, ever, really, ridiculous, significant, So, such, super, truly, ultimate, undoubtedly, very. |

that follow negation terms with their opposite sentiment words.

## IV.   COMPARATIVE ANALYSIS

This paper present a four tier framework for feature fusion of text feature extraction (POS, BOW and HT) and negation scope detection technique. Comparative analysis are present interesting and useful facts regarding the state-of-the-art of four benchmark sentiment classifier with feature fusion (as mention in Table II). Proposed framework use to comparing the performance supervised sentiment classifier after preprocessing and feature fusion for negation sentiment classification as shown in Fig. 2.

### A. Social Media Massage Pre-Processing

Social media post and tweets contain high rich of domain specific slag language, emoticons, symbols, idioms and sarcastic sentences. For accurate sentiment analysis proposed framework explored the unique properties social media data and try to refine by sentence splitting, slag replacement, word normalization and negation control pre-processing step for better sentiment classification.

*1) Sentence Splitting:* In proposed framework delimiter ('.','?','!',',',';') are used to split social media post into different sentence level. For example consider the review of cell phone posted by reviewer $R_1$.

### Cell Phone- User Review ($R_1$)

*"I bought a dual camera cell phone last week. Camera resolution is awesome, having lower battery life, but its ok for me. I m loooooving it."*

Sentence splitting phase split the review $R_1$ into five different sentence as sentence $S_1$ , $S_2$,$S_3$, $S_4$ and $S_5$ .

$S_1$:- I bought a dual camera cell phone last week.

$S_2$:- Camera resolution is awesome.

$S_3$:- having lower battery life.

$S_4$:- but its ok for me.

$S_5$:- I m loooooving it.

Fig. 2.    Feature extraction and Pre-processing

*2) Slag Replacement:* Proposed Framework use domain specific Slag and emoticon corpus for slag replacement. For example consider the unprocessed comment $C_1$ and $C_2$ where tokens 'Ur' and 'lol' are compared to entries in slag corpus and return processed comment $C_3$ and $C_4$ with token 'Your' and 'laughing out loud'.

**Unprocessed comment** $C_1$: Ur sound is really pleasant.

**Unprocessed comment** $C_2$: It's Really Good . lol!.

**Processed comment** $C_3$: Your sound is really pleasant.

**Processed comment** $C_4$: It's Really Good. laughing out loud!.

*3) Word Normalization:* Proposed Framework use Rogets Thesaurus corpus for word normalization by keyword matching. For normalization phase of post are match with entries in Rogets Thesaurus. If missed, repeated letters are sub sequentially compact until it's not matched. For example consider the unprocessed comment $C_5$ where the token 'gooooood' are compared to entries in Rogets thesaurus and return refine one i.e. 'good' with processed comment $C_6$.
**Unprocessed comment** $S_5$: Its really Gooooood.
**Processed comment** $S_6$: Its really Good.

*B. Text Feature Extraction from Social Media Post*

Once the social media massages are preprocessed, processed massages are passed for sentiment classification. For relevant classification this paper deploys bag-of-words (BoW), feature hashing (FH), and POS feature extraction technique to extract and select text features.

*1) Parts of Speech (POS) tagger:* POS Taggers provide syntax analysis of social media posts or tweets, and annotated each word as noun, verb, adjective, adverb, coordinating conjunction etc with a grammatical tagger . In sentiment analysis POS tagger used for Phrase identification, entity extraction and word sense disambiguation. POS Tagger employed probabilistic approach to evaluate the grammatical tagger and annotated highest probable tagger as shown in equation 1.

$$P\left(tag^i\middle|phase^j\right) = \lambda_1(phase^j)\frac{n(tag^i,phase^j))}{n(phase^j)}$$
$$+\lambda_2(phase^j)\frac{n_m(tag^i)}{n_m()} \qquad (1)$$

Where

- $P(tag^i|phase^j)$ is the probability of tagger i annotated over phase j.

- $n(tag^i, phase^j)$ is number of times phase j appears with grammatical tagger i.

- $n(phase^j)$ is number of times phase j appears.

- $n_m(tag^i)$ is number of times a phase that had never been seen with grammatical tagger i gets grammatical tagger i.

- $n_m()$ is number of such occurrences in total.

$$\lambda_1(phase^j) = \begin{cases} 1 & if\ n(phase^j) \geq 1 \\ 0 & otherwise \end{cases} \qquad (2)$$

For Sentiment analysis, adjectives (grammatical tagger) are fine source of polarity for opinioned word in message. Consider the unprocessed comment $C_7$ and their resultant pos tagger provided by Stanford parser [http://nlp.stanford.edu:8080/ parser/index.jsp]. In processed comment $C_8$ word "*nice*" is adjective that shown polarity of comment $C_8$ about the entity "*Camera*".

**Unprocessed Comment** ($C_7$)**:**It is not a nice camera.
**Processed Comment** ($C_8$)**:** it/**PRP** is/**VBZ** not/**RB** a/**DT** nice/**JJ** camera/**NN**

*2) Bag-of-Words (BoW):* For sentiment analysis, bag-of-words is use to transforms social media post or tweets into weighted vectors that contain relative polarity score of each word in massage. BoW independently tackle each word (token) in a tweet as order-invariant collection of features as shown in equation 3.

$$M = t_1, t_2, t_3, ,,,,, t_n \qquad (3)$$

In sentiment analysis short phase of word should capture better sentiment then single word. Bag-of-word work over that principle and consider bigram, trigram or n-gram phase of words for polarity score with help of sentiment lexicon. Consider the review tweet $C_9$ about the quality of phone . Unigram work over single word token "*bad*" whereas bigram take two word phase token i.e. " *very bad*" for calculating the polarity of comment.

**Comment** ($C_9$)**:** *It is a very bad phone.*

Word phase " *very bad*" defiantly has higher negative polarity value than "*bad*"

*3) Feature Hashing (FH):* Hashtag is a opinioned term that labeled itself by social media user at end of their tweets to convey their sentiment and opinion. Generally social media user use hashtag to convey their sarcastic.
Consider the review comment $C_{10}$ about movie, which is not positive but reviewer labeled their positive sentiment at end of tweet to convey their actual feeling .

**Comment** ($C_{10}$)**:**Movie is unpredictable. *#awesome*

*C. Negation Feature Extraction*

Negation control in sentiment analysis are involve two sub task specifically negative cues and scope detection. Negative cues detection is responsible to recognize the negative influence phase or term in sentences.For negation control, proposed framework use rule based keyword matching technique for negative cue detection and conjunction analysis, punctuation mark identification and grammatical dependency tree for scope detection technique.

*1) Negation cues detection:* Negation cues are the term or the phase of word that reflect negativity in review post. Proposed framework identify the negation cues by keywords matching technique from negation words corpus and replaced by token "NEGATION" as shown in negation feature

extraction section in Fig. 2.

For example consider the comment $C_{11}$. Where negation word 'not' identify by keyword matching and replace by taken 'NEGATION' for further treatment as shown in comment $C_{12}$.

**Sentence** $C_{11}$**:**Battery Life of this cell phone is not long but I am happy with its camera resolution.

**Sentence** $C_{12}$**:**Battery Life of this cell phone is NEGATION long but I am happy with its camera resolution.

*2) Negation Scope detection:* - Scope detection technique figure out the linguistic impact of negation cues in opinion sentences. Proposed framework use conjunction word analysis (CWA), Punctuation mark identification (PMI) and Grammatical dependency tree (GDT) scope detection technique to figure out the linguistic coverage of "NEGATION" token labeled by negation cues detection phase.

(a) **Conjunction Analysis:** Conjunction words determine and fixed the influence of negative word that comes before and after the occurrence of "NEGATION" token. For example consider the Comment $C_{13}$ where one lady post different opinion about different aspect of beauty product. Lady have negative opinion about price but positive opinion about it quality. In comment $C_{13}$ conjunction word "but" help to figure out the influence of two opposite sentiment orientated opinioned word "good" and "expensive" before and after its appearance.

Comment ($C_{13}$):- "This Sunscreen lotion is really good but it's too expensive."

Some other Conjunction word such as "expect", "however", "whereas", "although", "and", "or", "unless", "nevertheless" help to figure out the influence of negative token in sentence. Conjunction word "AND" some time fail to figure out the scope of negation. For example consider the sentence $C_{14}$ where negation word "doesn't" invert the polarity of both "good" and "nice" sentiment word.
Comment ($C_{14}$):- "This cell phone doesn't have good battery backup and nice camera quality."

Whereas as proposed feature fusion, Text feature extraction technique POS [24], [25], BOW [21], [16], [19], [13], [26], [27] and Hashtag help the overcome the limitation of conjunction word "AND" through grammatical marking, sentiment word and sarcasm identification respectively and simultaneously lead to evaluate polarity score of different part of sentence.

(b) **Punctuation Mark Identification:** - Punctuation Mark (",", "!", ";") limit the influence of negation between "NEGATION" token and next punctuation mark. For example consider the production manager comment over company last year production in sentence $C_{15}$. Where manager is really upset about current year production but he hopeful for next year. In this comment comma "," is use to separate out these two sentiment of production manager.

TABLE II.    FEATURE FUSION CASE

| Feature Fusion | Text feature | Scope Detection |
|---|---|---|
| Case 1 | POS | CWA |
| Case 2 | POS | PMI |
| Case 3 | POS | GDT |
| Case 4 | BOW | CWA |
| Case 5 | BOW | PMI |
| Case 6 | BOW | GDT |
| Case 7 | HT | CWA |
| Case 8 | HT | PMI |
| Case 9 | HT | GDT |

Sentence $(C_{15})$:- "The production of this year is not up to mark, we are hopeful for next year."

Punctuation Mark "," some time fail to figure out the scope of negation. For example consider the comment $C_{16}$ where negation word "doesn,t" invert the polarity of both "good" and "nice" sentiment word.

Sentence $(C_{16})$:- "This cell phone doesn't have good battery backup, nice camera quality and touchscreen."

Whereas as proposed feature fusion, Text feature extraction technique POS [24], [25], BOW [21], [16], [19], [13], [26], [27] and Hashtag help the overcome the limitation of punctuation mark "," through grammatical marking, sentiment word and sarcasm identification respectively and simultaneously lead to evaluate polarity score of different part of sentence.

(c) **Grammatical Dependency Tree:** - Grammatical dependency between orders of occurrence of sentiments oriented word and negative cue help to figure out influence of NEGATION token [18]. Grammatical dependency parser build syntactic tree [28] and their lowest level are help to figure out scope of negation. Text feature extraction technique POS [24], [25], BOW [21], [16], [19], [13], [26], [27] and Hashtag help for grammatical marking lead to evaluate lowest level of grammatical syntactic relationship.

### D. Sentiment Classification

After examine the text feature extraction (POS, BOW, HT) and scope of negation (CWA, PMI, GDT) technique, proposed framework present nine one too many feature fusion case from Text feature to Scope of negation as shown in table. Feature fusion improve the performance of feature extraction by overcome the limitation of their subordinate. This paper evaluate the performance of Classifiers SVM, Naives Bayes, Linear regression and random Forest after incorporating the different feature fusion case as shown in table.



Fig. 3.    Support Vector Machine For SA

*1) Support Vector Machine::* In proposed framework, SVM determine the optimal hyper plane $(W_{ff}P_S + b)$ based on feature fusion to maximize feature margin $(f_m)$ between positive and negative polarity social media post and tweets as shown in Fig. 3.

Support vector machine for sentiment classification [3], classifier the preprocessed massage dataset $M_{ff}$ after feature fusion. Where the performance of polarity classification depend upon type of feature fusion applied. After incorporating feature fusion technique for negative sentence sentiment analysis, SVM treat all the token in scope of negation as feature fusion vector space as shown in equation 4.

$$W_{vs}^{ff} = \{(W_i^s, \{W_{sn}\}) \, n_t \in m_{ff}\} \qquad (4)$$

Where

- $m_{ff}$ is pre-processed text data set after incorporating Feature fusion.

- $W_{vs}^{ff}$ is Feature fusion vector space.

- $W_i^s$ is the sentiment word in negative scope.

- $W_{sn}$ is the set of word in scope of negation.

- $n_t$ is negative token.

Preprocessed massage data set $(M_{ff})$ is the set of n couple of element $(t_i, P_c)$, where $t_i$ is associated with token within the $M_{ff}$ and $P_c$ indicate their respective polarity class (+ve ,

-ve) as shown in equation. $t_i$ can be capture by using feature fusion technique as shown in equation 5.

$$M_{ff} = \{(t_i, p_c)\, t_i \in ff, p_c \in \{+ve, -ve\}_{i=1}^n\} \quad (5)$$

The Feature fusion vectors that define the hyper plane are the support sentiment feature fusion vectors (ffv) as shown in equation 6.

$$ffv = \{(Superb, +ve), (Best, +ve), (Horrible, -ve)\} \quad (6)$$

In proposed framework, SVM is needed to maximize the width of the feature margin ($f_m$). Where

$$(W_{ff}.P_c + b_1) \geq +ve\forall \quad Positive\ Sentiment\ over$$
$$the\ Positive\ hyperplane \quad (7)$$

$$(W_{ff}.P_c + b_2) \geq -ve\forall \quad Negative\ Sentiment\ over$$
$$the\ Negative\ hyperplane \quad (8)$$

Feature margin between positive and negative hyber plane is

$$f_m = \{(W_{ff}.P_c + b_1) - (+ve), (W_{ff}.P_c + b_2) - (-ve)\} \quad (9)$$

$$f_m = \frac{|b_1 - b_2|}{\|W_{ff}\|} \quad (10)$$

To maximized the feature margin ($f_m$) , it's needed to minimized weight of sentiment feature vector space ($W_{ff}$).

*2) Naïve Bayes:* In proposed framework Naïve Bayes determine the polarity class (+ve,-ve) of any preprocessed massage data set $M_{ff}$ after feature fusion on the basis of maximum posterior probability as shown in equation 11 and 12 [3].

$$M_{ff}^p = argmax_{p\in(Positive,Negative)}P(p|M_{ff}) \quad (11)$$

$$P(p|M_{ff}) = \frac{P(M_{ff}|p)P(p)}{P(M_{ff})} \quad (12)$$

Where P ($p|M_{ff}$) is final posterior probability and P($M_{ff}|p$) is the probability of sentence $M_{ff}$ belong to polarity class $P_c$. Whereas P(p) and P($M_{ff}$) is the independent probability polarity class $P_c$ and sentence $M_{ff}$.Whereas after incorporating feature fusion vector ($ff_v$) as a relevant feature for negative sentiment analysis, NB treat all the token in $ff_v$ as independent probability entity as shown in equation 13.

$$P(n|ff_v) = P(n|M_{ff}^1) * P(n|M_{ff}^2) * P(n|M_{ff}^3) *$$
$$* ..... * P(n|M_{ff}^n)) \quad (13)$$

Where P ($n|ff_v$) are independent given the polarity Class ($P_c$) and each word in scope of negation substitute their individual probability for exploring polarity classes.

*3) Random forest:* In proposed framework Random forest predict the polarity class (+ve, -ve) for preprocessed massage data set ($M_{ff}$) after incorporating feature fusion. Random forest predict the sentiment polarity class of sentence ($M_{ff}$) by building randomized regression trees $\{ff_n$ (c,$p_c$,$M_{ff}$)m≥1$\}$ based relationship between polarity class and sentences as shown in equation 14.

$$\overline{ff_n}(c, M_{ff}) = E_{pc}[ff_n(c, p_c, M_{ff})] \quad (14)$$

Where $E_{pc}$ is exception on polarity class ($P_c$) classification with random feature fusion parameter (ff) on condition c and data set ($M_{ff}$). Whereas incorporation of Feature fusion vector of negation as conditional parameter 'c' lead to minimized exception (Epc) on polarity class and increase classification rate.

*4) Linear regression::* In proposed framework linear regression find a feature fusion based decision boundary that linearly distinct positive and negative polarity classes as shown in equation 15.

$$P_c(M_{ff}) = \begin{cases} +ve & if\ C * M_{ff} \geq 0\ (Polarity\ score) \\ -ve & if\ C * M_{ff} < 0\ (Polarity\ score) \end{cases} \quad (15)$$

$$P_c(M_{ff}) = threshold\ C * M_{ff} \quad (16)$$

Where P passing the polarity function $C * m_{ff}$ through the threshold function as shown in equation 16.

## V. Environment Setup Result Analysis

For performance analysis of recent benchmark classification technique (NB, SVM, RF and LR) over five different social media data set from two different source total nine different experiment has been carried out.Nine different experiment belong to one to many nine different feature fusion case that build in proposed framework as shown in Table II. All the nine different experiment has been carried out over 5 different social media post and tweets data set. First two data set is scraped by twitter API i.e. Stanford data set (TSCDS) [29] and Sanders Twitter Sentiment Corpus data set (TSDS) [30]. Stanford data set contain 160000 training tweets accompanied by 80000 both positive and negative tweets. Whereas Sanders Twitter Sentiment data set contain 570 positive and 654 negative tweets. However last three data set has been carried out from amazon online product reviews data set of smartphone (ASPR), movies (AMR) and book (ABR) [31]. Detail description of data set composition is summarized in Table III.

Performance evaluation of benchmark sentiment classifier with and without feature fusion for negation control are described in Table IV. Performance of classifier has been increased after incorporating feature fusion over negative social media post or tweets.

The baseline classifier (SVM, Nave Bayes, Random Forest and Linear Regression) without feature fusion for negation control can achieve approximate 45%- 55% and 50% - 65%

TABLE III.     DATA SET DESCRIPTION

| Ref. | Platform | Data set Name | Total number of tweets/Review | Positive (tweets/Review) | Negative (tweets/Review) |
|---|---|---|---|---|---|
| [29] | Twitter | TSCDS | 160000 | 80000 | 80000 |
| [30] | Twitter | TSDS | 1224 | 570 | 654 |
| [31] | Amazon | ASPR | 17500 | 12500 | 5000 |
| [31] | Amazon | AMR | 35000 | 30000 | 5000 |
| [31] | Amazon | ABR | 90000 | 81000 | 9000 |

TABLE IV.     COMPARATIVE ANALYSIS OF SENTIMENT ANALYSIS TECHNIQUE

| CT | FET | NSDT | Twitter API | | Amazon | | |
|---|---|---|---|---|---|---|---|
| | | | TSCDTS | Stanford Dataset | SmartPhone Review | Movie Review | Book Review |
| NB | – | – | 53.77 | 48.45 | 55.67 | 61.46 | 65.68 |
| | POS | CWA | 89.45 | 84.23 | 87.46 | 90.26 | 92.78 |
| | | PMI | 90.12 | 86.68 | 89.24 | 91.56 | 93.48 |
| | | GDT | 91.74 | 88.65 | 90.42 | 92.87 | 95.68 |
| | BOW | CW | 82.22 | 78.64 | 80.2 | 84.66 | 88.4 |
| | | PM | 84.77 | 82.26 | 83.12 | 87.9 | 90.66 |
| | | GDT | 88.42 | 84.52 | 86.78 | 90.64 | 92.44 |
| | FHT | CW | 78.21 | 72.2 | 70.44 | 78.78 | 80.82 |
| | | PM | 80.87 | 76.36 | 76.89 | 82.94 | 85.69 |
| | | GDT | 84.28 | 80.68 | 82.25 | 86.86 | 88.27 |
| SVM | - | - | 49.71 | 45.23 | 53.78 | 60.85 | 64.66 |
| | POS | CW | 87.28 | 82.36 | 86.68 | 90.12 | 92.48 |
| | | PM | 88.45 | 86.24 | 88.99 | 91.2 | 93.38 |
| | | GDT | 90.85 | 87.56 | 89.86 | 92.78 | 95.42 |
| | BOG | CW | 82.88 | 79.28 | 80.58 | 85.42 | 88.88 |
| | | PM | 85.48 | 83.24 | 84.2 | 88.22 | 91.25 |
| | | GDT | 89.62 | 85.12 | 87.48 | 91.18 | 93.28 |
| | HT | CW | 76.44 | 71.88 | 69.86 | 78.28 | 80.24 |
| | | PM | 78.46 | 75.68 | 75.12 | 82.74 | 85.28 |
| | | GDT | 82.58 | 79.64 | 81.88 | 86.14 | 88.16 |
| RF | - | - | 54.51 | 51.37 | 52.67 | 58.68 | 63.22 |
| | POS | CW | 84.58 | 82.08 | 83.46 | 84.38 | 88.2 |
| | | PM | 85.26 | 86.12 | 86.04 | 86.42 | 90.12 |
| | | GDT | 88.78 | 84.89 | 87.67 | 88.28 | 92.27 |
| | BOG | CW | 80.12 | 77.48 | 78.29 | 82.26 | 86.64 |
| | | PM | 83.26 | 81.12 | 82.2 | 85.02 | 88.94 |
| | | GDT | 86.34 | 83.84 | 84.67 | 87.22 | 90.86 |
| | HT | CW | 75.28 | 71.88 | 69.28 | 77.2 | 79.68 |
| | | PM | 77.37 | 75.68 | 74.56 | 81.98 | 85.24 |
| | | GDT | 82.24 | 78.42 | 81.28 | 85.86 | 87.48 |
| LR | - | - | 55.12 | 52.23 | 50.98 | 58.42 | 64.42 |
| | POS | CW | 83.78 | 82.55 | 84.65 | 88.64 | 90.2 |
| | | PM | 84.26 | 86.04 | 87.28 | 90.68 | 91.58 |
| | | GDT | 86.24 | 85.28 | 88.98 | 91.78 | 93.12 |
| | BOG | CW | 79.25 | 77.89 | 76.2 | 81.98 | 87.24 |
| | | PM | 82.14 | 82.26 | 79.68 | 86.99 | 89.68 |
| | | GDT | 85.62 | 84.26 | 82.12 | 86.98 | 90.88 |
| | HT | CW | 78.98 | 74.64 | 72.88 | 79.14 | 82.22 |
| | | PM | 81.28 | 78.88 | 77.28 | 83.22 | 85.36 |
| | | GDT | 84.88 | 81.27 | 84.2 | 87.42 | 89.24 |

accuracy rate over twitter and Amazon social media data set respectively as shown in Fig. 4. Where linear regression achieve better performance over twitter data set and lead by approximate 1.4% over other classifier. Whereas over amazon review data set NB lead the performance by approximate 2% over rest.

The performance of baseline classifier is significantly boost up after incorporating feature fusion for negation control. In feature fusion case 1 i.e. incorporate POS with conjunction word analysis scope detection technique, NB (89.45%, 84.23%, 87.46%, 90.26% and 92.78%) , SNM (87.28%, 82.36%, 86.68%, 90.12%, 92.48%), RF (84.58%, 82.08%, 83.46%, 84.38%,88.2%) and LR (83.78%, 82.55%, 84.65%,

88.64%, 90.20%) significantly boost the performance by approximate 40.84%-73.85% , 43.03% -82.10%, 39.52%-59.79% and 40.02%-58.06% respectively over five different variant social media data set as shown in Fig. 5(a) & (b).

Classifier gain higher improvement over twitter data set is significantly due to presence of higher number of negative post i.e. approximate 50% and 53 .43 % in TSCDS and TSD data set respectively. Correspondingly lower improvement in Amazon data set due to presence of lower number of negative tweets i.e. approximate 28.57%, 14.28% and 10% in ASP, AMR and ABR data set respectively. Whereas with feature fusion case 1 (POS+CWA), NB lead the performance and SVM gain highest improvement as shown in Fig. 5(a) & (b).

Fig. 4. Sentiment Analysis with benchmark algorithm

In feature fusion case 2 (POS+PMI) , NB (90.12%, 86.68%, 89.24%, 91.56% and 93.48%) , SNM (88.45%, 86.24%, 88.99%, 91.20%, 93.38%) , RF (85.26%, 86.12%, 86.04%, 86.42%,90.12%) and LR (84.26%, 86.04%, 87.28%, 90.68%, 91.58%) significantly boost the performance by approximate 41.90%-78.91% , 44.42% -90.67%, 42.55%-67.65% and 42.17%-71.21% respectively over five different variant social media data set as shown in Fig. 6(a) &(b). In feature fusion case 2(POS+CWA), NB lead the performance. Whereas SVM and LR gain highest improvement twitter and Amazon data set respectively as shown in Fig. 6(a) &(b).

In feature fusion case 3 (POS+GDT) , NB (91.74%, 88.65%, 90.42%, 92.87% and 95.68%) , SVM (90.85%, 87.56%, 89.86%, 92.78%, 95.42%) , RF (88.78%, 84.89%, 87.67%, 88.28%,92.27%) and LR (86.24%, 85.28%, 88.98%, 91.78%, 93.12%) significantly boost the performance by approximate 45.24%-82.98% , 47.58% -93.59%, 45.96%-66.46% and 44.56%-74.54% respectively over five different variant social media data set as shown in Fig. 7(a) & (b). In feature fusion case 3(POS+GDT), NB lead the performance. Whereas SVM and LR gain highest improvement twitter and Amazon data set, respectively as shown in Fig. 7(a) & (b).

In feature fusion case 4 (BOW+CWA), NB (82.22%, 78.64%, 80.2%, 84.66%, 88.4%), SVM (82.88%, 79.28%, 80.58%, 85.42%, 88.88%), RF (80.12%, 77.48%, 78.29%, 82.26%, 86.64%) and LR (79.25%, 77.89%, 76.2%, 81.98%, 87.24%) significantly boost the performance by approximate 34.19%-62.32% , 37.46% -75.29%, 37.05%-50.83% and 35.43%-49.48% respectively over five different variant social media data set as shown in Fig. 8(a) & (b). In feature fusion case 4(BOW+CWA), NB lead the performance and SVM gain highest improvement as shown in Fig. 8(a) & (b).

In feature fusion case 5 (BOW+PMI), NB (84.77%, 82.26%, 83.12%, 87.90%, 90.66%), SVM (85.48%, 83.24%, 84.20%, 88.22%, 91.25%) , RF (83.26%, 81.12%, 82.20%, 85.02%, 88.94%) and LR (82.14%, 82.26%, 79.68%, 86.99%, 89.68%) significantly boost the performance by approximate 37.62%-69.79% , 41.13% -84.04%, 40.69%-57.92% and 39.22%-57.50%, respectively over five different variant social

media data set as shown in Fig. 9(a) & (b). In feature fusion case 5(BOW+PMI), NB lead the performance. Whereas SVM and LR gain highest improvement twitter and Amazon data set, respectively as shown in Fig. 9(a) & (b).

In feature fusion case 6 (BOW+GDT), NB (88.42%, 84.52%, 86.78%, 90.64%, 92.44%), SVM (89.62%, 85.12%, 87.48%, 91.18%, 93.28%), RF (86.34%, 83.84%, 84.67%, 87.22% 90.86%) and LR (85.62%, 84.26%, 82.12%, 86.98%, 90.88%) significantly boost the performance by approximate 40.32%-74.45% , 44.27% -88.20%, 43.73%-63.21% and 41.08%-61.33%, respectively over five different variant social media data set as shown in Fig. 10(a) & (b). In feature fusion case 6(BOW+GDT), NB lead the performance and SVM gain highest improvement as shown in Fig. 10(a) & (b).

In feature fusion case 7 (HT+CWA), NB (78.21%, 72.20%, 70.44%, 78.78%, 80.82%), SVM (76.44%, 71.88%, 69.86%, 78.28%, 80.24%), RF (75.28%, 71.88%, 69.28%, 77.20%, 79.68%) and LR (78.98%, 74.64%, 72.88%, 79.14%, 82.22%) significantly boost the performance by approximate 22.68%-49.02%, 24.1% -58.93%, 26.04%-39.93% and 27.64%-42.91%, respectively over five different variant social media data set as shown in Fig. 11(a) & (b). In feature fusion case 7(HT+CWA), LR lead the performance. Whereas SVM and LR gain highest improvement over twitter and Amazon data set respectively as shown in Fig. 11(a) &(b).

In feature fusion case 8 (HT+PMI), NB (80.87%, 76.36%, 76.89%, 82.94%, 85.69%), SVM (78.46%, 75.68%, 75.12%, 82.74%, 85.28%), RF (77.37%, 75.68%, 74.56%, 81.98%, 85.24%) and LR (81.28%, 78.88%, 77.28%, 83.22%, 85.36%) significantly boost the performance by approximate 30.07%-57.61%, 31.89% -67.33%, 34.84%-47.33% and 32.51%-51.03% respectively over five different variant social media data set as shown in Fig. 12(a) & (b). In feature fusion case 8(HT+PMI), LR lead the performance. Whereas SVM and LR gain highest improvement over Twitter and Amazon data set respectively as shown in Fig. 12(a) & (b).

In feature fusion case 9 (HT+GDT), NB (84.28%, 80.68%, 82.25%, 86.86%, 88.27%), SVM (82.58%, 79.64%, 81.88%, 86.14%, 88.16%), RF (82.24%, 78.42%, 81.28%, 85.86%, 87.48%) and LR (84.88%, 81.27%, 84.20%, 87.42%, 89.24%) significantly boost the performance by approximate 33.99%-66.53%, 36.35% -76.08%, 38.38%-52.66% and 38.53%-65.17% respectively over five different variant social media data set as shown in Fig. 13(a) & (b). In feature fusion case 9(HT+GDT), LR lead the performance. Whereas SVM and LR gain highest improvement over Twitter and Amazon data set respectively as shown in Fig. 13(a) & (b).

With different angle of evaluating the performance of classifier over all nine feature fusion technique. It is observed that classifiers gives better performance with feature fusion case 3 (POS+GDT) .

Naïve Bayes gain 96.68% accuracy with Feature Fusion Case 3 over different variant of data set as shown in Fig. 14(a). Naïve Bayes achieved highest improvement with Case3 i.e. approximate 82.98% over different variant of data set as shown in Fig. 14(b).

SVM gain 95.42% accuracy with Feature Fusion Case 3 over different variant of data set as shown in Fig. 15(a). SVM

Fig. 5. Feature Fusion Case 1:- POS and Conjunction Word Analysis for SA with Negation Scope detection Technique



Fig. 6. Feature Fusion Case 2:- POS and PMI based Feature Extraction for SA with Negation Scope detection Technique

achieved highest improvement with Case3 i.e. approximate 93.59% over different variant of data set as shown in Fig. 15(b).

RF gain 92.27% accuracy with Feature Fusion Case 3 over different variant of data set as shown in Fig. 16(a). RF achieved highest improvement with Case 2 i.e. approximate 67.65% over different variant of data set as shown in Fig. 16(b).

LR gain 93.12% accuracy with Feature Fusion Case 3 over different variant of data set as shown in Fig. 17(a). LR achieved highest improvement with Case 3 i.e. approximate 74.54% over different variant of data set as shown in Fig. 17(b).

After evaluating the performance baseline sentiment classi-

fier with feature fusion following outcome has been acquired. POS+GDT is best suited feature extraction and Scope detection technique to identify the range of influence marked by negation for negative sentiment Analysis. Whereas other gives biased result. NB is best suited sentiment classification approach under negation for case 1 to case 6 but for case 7 to case 9 LR achieve highest performance.Whereas SVM achieved highest improvement after encapsulating feature fusion with classification.

## VI. CONCLUSION

This paper present a framework for comparative analysis to analysis the performance of benchmark supervised sen-

Fig. 7. Feature Fusion Case 3:- POS and GDT based Feature Extraction for SA with Negation Scope detection Technique



Fig. 8. Feature Fusion Case 4:- BOW and CWA based Feature Extraction for SA with Negation Scope detection Technique

timent classifier (NB, SVM, RF, LR) for negation control. Proposed framework incorporate text feature extraction, Negation cue and scope detection technique as feature fusion for significantly improve the performance for negation control. This paper present a comparative analysis of incorporation of feature fusion with supervised sentiment classification technique for negative control over social media data set. Social media post or tweets may contain noise, misspelled word, emoticon and Slag language that required to be preprocess before feature extraction and sentiment analysis. Proposed framework initially preprocessed social media post or tweets to tackle noise, misspelled words and slag languages. And finally classify the tweets according to their polarity score after

incorporating feature fusion for negation control. For the negation control feature fusion case 3 (POS+GDT) is best suited feature extraction technique that improve the performance of NB by approximate 45.24% to 82.98%, NB by approximate 47.58% -93.59%, RF by approximate 45.96% -66.46% and LR by approximate 44.56% -63.28% over different variant of social media data set. It is observed that NB is best suited sentiment classification approach under feature fusion for negation whereas SVM achieved highest improvement over different variant of social media data set.

Fig. 9.  Feature Fusion Case 5:- BOW and PMI based Feature Extraction for SA with Negation Scope detection Technique



Fig. 10.  Feature Fusion Case 6:- BOW and GDT based Feature Extraction for SA with Negation Scope detection Technique

## REFERENCES

[1] G. Katz, N. Ofek, and B. Shapira, "Consent: Context-based sentiment analysis," *Knowledge-Based Systems*, vol. 84, pp. 162 – 178, 2015.

[2] F. H. Khan, U. Qamar, and S. Bashir, "Swims: Semi-supervised subjective feature weighting and intelligent model selection for sentiment analysis," *Knowledge-Based Systems*, vol. 100, pp. 97 – 111, 2016.

[3] N. K. Singh, D. S. Tomar, and A. K. Sangaiah, "Sentiment analysis: a review and comparative analysis over social media," *Journal of Ambient Intelligence and Humanized Computing*, May 2018.

[4] E. A. Stepanov and G. Riccardi, "Detecting general opinions from customer surveys," in *2011 IEEE 11th International Conference on Data Mining Workshops*, Dec 2011, pp. 115–122.

[5] L.-C. Yu, J.-L. Wu, P.-C. Chang, and H.-S. Chu, "Using a contextual entropy model to expand emotion words and their intensity for the sentiment classification of stock market news," *Knowledge-Based Systems*, vol. 41, no. Supplement C, pp. 89 – 97, 2013.

[6] R. Feldman, "Techniques and applications for sentiment analysis," *Commun. ACM*, vol. 56, no. 4, pp. 82–89, Apr. 2013.

[7] G. Mishne and N. Glance, *Predicting Movie Sales from Blogger Sentiment*, 01 2006.

[8] C. Cardie, C. Farina, and T. Bruce, "Using natural language processing to improve erulemaking: Project highlight," in *Proceedings of the 2006 International Conference on Digital Government Research*, ser. dg.o '06.  Digital Government Society of North America, 2006, pp. 177–178.

[9] E. Cambria and A. Hussain, *Sentic Computing: A Common-Sense-Based Framework for Concept-Level Sentiment Analysis*, 1st ed.  Springer Publishing Company, Incorporated, 2015.

[10] W. Li, K. Guo, Y. Shi, L. Zhu, and Y. Zheng, "Dwwp: Domain-specific new words detection and word propagation system for sentiment

Fig. 11.   Feature Fusion Case 7:- HT and CWA based Feature Extraction for SA with Negation Scope detection Technique



Fig. 12.   Feature Fusion Case 8:- HT and PMI based Feature Extraction for SA with Negation Scope detection Technique

analysis in the tourism domain," *Knowledge-Based Systems*, vol. 146, pp. 203 – 214, 2018.

[11]  M. A. Mirtalaie, O. K. Hussain, E. Chang, and F. K. Hussain, "Extracting sentiment knowledge from pros/cons product reviews: Discovering features along with the polarity strength of their associated opinions," *Expert Systems with Applications*, vol. 114, pp. 267 – 288, 2018.

[12]  F. Wu, Y. Song, and Y. Huang, "Microblog sentiment classification with heterogeneous sentiment knowledge," *Inf. Sci.*, vol. 373, no. C, pp. 149–164, Dec. 2016.

[13]  M. Abdul-Mageed, "Modeling arabic subjectivity and sentiment in lexical space," *Information Processing and Management*, 2017.

[14]  M. Ghiassi and S. Lee, "A domain transferable lexicon set for twitter sentiment analysis using a supervised machine learning approach," *Expert Systems with Applications*, vol. 106, pp. 197 – 216, 2018.

[15]  O. Appel, F. Chiclana, J. Carter, and H. Fujita, "A hybrid approach to the sentiment analysis problem at the sentence level," *Knowledge-Based*

*Systems*, vol. 108, pp. 110 – 124, 2016, new Avenues in Knowledge Bases for Natural Language Processing.

[16]  V. García-Díaz, J. P. Espada, R. G. Crespo, B. C. P. G-Bustelo, and J. M. C. Lovelle, "An approach to improve the accuracy of probabilistic classifiers for decision support systems in sentiment analysis," *Applied Soft Computing*, vol. 67, pp. 822 – 833, 2018.

[17]  I. Korkontzelos, A. Nikfarjam, M. Shardlow, A. Sarker, S. Ananiadou, and G. H. Gonzalez, "Analysis of the effect of sentiment analysis on extracting adverse drug reactions from tweets and forum posts," *J. of Biomedical Informatics*, vol. 62, no. C, pp. 148–158, Aug. 2016.

[18]  C. Diamantini, A. Mircoli, D. Potena, and E. Storti, "Social information discovery enhanced by sentiment analysis techniques," *Future Generation Computer Systems*, 2018.

[19]  T. Kang, S. Zhang, N. Xu, D. Wen, X. Zhang, and J. Lei, "Detecting negation and scope in chinese clinical notes using character and word

Fig. 13. Feature Fusion Case 9:- HT and GDT based Feature Extraction for SA with Negation Scope detection Technique

embedding," *Computer Methods and Programs in Biomedicine*, vol. 140, pp. 53 – 59, 2017.

[20] N. Pröllochs, S. Feuerriegel, and D. Neumann, "Negation scope detection in sentiment analysis: Decision support for news-driven trading," *Decision Support Systems*, vol. 88, pp. 67 – 75, 2016.

[21] E. S. Tellez, S. Miranda-Jiménez, M. Graff, D. Moctezuma, O. S. Siordia, and E. A. Villaseñor, "A case study of spanish text transformations for twitter sentiment analysis," *Expert Systems with Applications*, vol. 81, pp. 457 – 471, 2017.

[22] S. M. Jiménez-Zafra, M. T. Martín-Valdivia, M. D. Molina-González, and L. A. Ureña-López, "Relevance of the sfu reviewsp-neg corpus annotated with the scope of negation for supervised polarity classification in spanish," *Information Processing and Management*, vol. 54, no. 2, pp. 240 – 251, 2018.

[23] M. T. AL-Sharuee, F. Liu, and M. Pratama, "Sentiment analysis: An automatic contextual analysis and ensemble clustering approach and comparison," *Data and Knowledge Engineering*, vol. 115, pp. 194 – 213, 2018.

[24] I. Habernal, T. Ptáček, and J. Steinberger, "Supervised sentiment analysis in czech social media," *Information Processing and Management*, vol. 50, no. 5, pp. 693 – 707, 2014.

[25] A. Ortigosa, J. M. Martín, and R. M. Carro, "Sentiment analysis

in facebook and its application to e-learning," *Computers in Human Behavior*, vol. 31, pp. 527 – 541, 2014.

[26] J.-C. Na, C. Khoo, and P. H. J. Wu, "Use of negation phrases in automatic sentiment classification of product reviews," *Library Collections, Acquisitions, and Technical Services*, vol. 29, no. 2, pp. 180 – 191, 2005.

[27] S. M. Jimenez-Zafra, M. T. M. Valdivia, E. M. Camara, and L. A. Urena-Lopez, "Studying the scope of negation for spanish sentiment analysis on twitter," *IEEE Transactions on Affective Computing*, pp. 1–1, 2018.

[28] A. Kennedy and D. Inkpen, "Sentiment classification of movie reviews using contextual valence shifters," *Computational Intelligence*, vol. 22, pp. 110–125, 2006.

[29] A. Go, R. Bhayani, and L. Huang, "Twitter sentiment classification using distant supervision," vol. 150, 01 2009.

[30] D. Ziegelmayer and R. Schrader, "Sentiment polarity classification using statistical data compression models," in *2012 IEEE 12th International Conference on Data Mining Workshops*, Dec 2012, pp. 731–738.

[31] H. Cho, S. Kim, J. Lee, and J.-S. Lee, "Data-driven integration of multiple sentiment dictionaries for lexicon-based sentiment classification of product reviews," *Knowledge-Based Systems*, vol. 71, pp. 61 – 71, 2014.

Fig. 14.    Performance of SVM over different feature fusion for SA with Negation Scope detection Technique



Fig. 15.    Performance of RF over different feature fusion for SA with Negation Scope detection Technique

Fig. 16.   Performance of NB over different feature fusion for SA with Negation Scope detection Technique



Fig. 17.   Performance of LR over different feature fusion for SA with Negation Scope detection Technique

# Use of Blockchain in Healthcare: A Systematic Literature Review

Sobia Yaqoob[1], Muhammad Murad Khan[*2], Ramzan Talib[3], Arslan Dawood Butt[4], Sohaib Saleem[5],
Fatima Arif[6], Amna Nadeem[7]

Department of Computer Science, Government College University Faisalabad, Faisalabad, Pakistan[1,2,3,5,6]
Department of Computer Science and Information Technology, University of Okara, Okara, Pakistan[1,5,7]
Department of Electrical Engineering, Government College University Faisalabad, Faisalabad, Pakistan[4]

*Abstract*—**Blockchain is an emerging field which works on the concept of a digitally distributed ledger and consensus algorithm removing all the threats of intermediaries. Its early applications were related to the finance sector but now this concept has been extended to almost all the major areas of research including education, IoT, banking, supplychain, defense, governance, healthcare, etc. In the field of healthcare, stakeholders (provider, patient, payer, research organizations, and supply chain bearers) demand interoperability, security, authenticity, transparency, and streamlined transactions. Blockchain technology, built over the internet, has the potential to use the current healthcare data into peer to peer and interoperable manner by using a patient-centric approach eliminating the third party. Using this technology, applications can be built to manage and share secure, transparent and immutable audit trails with reduced systematic fraud. This study reviews existing literature in order to identify the major issues of various healthcare stakeholders and to explore the features of blockchain technology that could resolve identified issues. However, there are some challenges and limitations of this technology which are needed to be focused on future research.**

*Keywords*—*Issues; healthcare; blockchain; systematic review*

## I. Introduction

The blockchain technology was devised by an unidentified person "Satoshi Nakamoto" in October 2008. He proposed a peer-to-peer, non-intermediated, electronic cash system introducing the first digital currency named as Bitcoin. This Distributed Ledger Technology (DLT) is time-stamped chain of transactional blocks, sealed with a cryptographic hash function and digital signature implementing trustless protocol [1] [2] [3]. Bitcoin was the first application of blockchain technology implemented in 2009 [4], [5].

Blockchain (BC) is in its early stages but it captivated immense response and interest from the community, domain experts and researchers in various fields like banking, Internet of Things, manufacturing, governance, education, and healthcare (HC). Blockchain, distributed ledger technology, can contribute numerous benefits to the healthcare industry. Blockchain technology's diverse features like decentralization, immutability, robustness, security, privacy, currency and consensus via cryptographic algorithms [6], [4], [5]; has the capacity to address the current issues of medical and healthcare sector. It can alleviate reliance on a single centralized authority which is more vulnerable to inaccuracy and insecurity. The interoperable infrastructure of blockchain technology will enhance the data exchange among various healthcare peers to improve coordination, quality of care, innovations, and market

competitions positively [7], [8], [9], [10]. This disruptive technology could have the potential to resolve the issues of counterfeit drugs, claims adjudication and patient billing management as it provides historical information to track the origin of transactions making all the actions transparent [9]. This technology can convert the current costly systems to cost-saving or even money generating systems [11] as the users are rewarded with digital currency as an incentive for their contribution [10], [12], [13]. A blockchain based infrastructure can be envisaged for improved decentralized record keeping, health data exchange, reliable drug supply, claim processes, and cost-effective systems. Blockchain has a tremendous potential to transform the current healthcare infrastructure. However, there are several challenges that have been identified and more research must be carried out to address these challenges such as scalability, security (threat of 51% attack), anonymity and fraud hype, disclosure of confidentiality, and environmental unsustainability.

This review paper explores how blockchain technology can revolutionize our current healthcare infrastructure. The sector comprises of different healthcare players (providers, payers, patients, vendors, manufacturers, and research institutes, etc.) performing different roles with different needs. This paper attempts to provide a deeper insight by discovering the scope of blockchain for all healthcare players in a single study and to identify the challenges of blockchain technology in the said sector.

This study is divided into six sections. Section I introduces the technology and its significance in the healthcare sector. Section II explores literature reviewed for this research which presents this trans-disciplinary study into following parts: first, it describes the blockchain technology and its features; second, it strives to find out studies highlighting most common issues of the healthcare sector and associated applications of blockchain technology in the sector; third, this work is compared with review papers to show the significance and need of this work. Section III narrates the research methodology comprising sections labeled as; the need of conducting a systematic literature review, motivation and research questions, search strategy, inclusion and exclusion criteria, classification criteria, and data extraction in accordance with the selected approach for conducting this systematic literature review. Section IV compiles the results to answer research questions heading to Section V which describes the probability of threats to the validity of this research. To end with Section VI, conclusion and future directions are documented for use of blockchain

technologies in the healthcare sector.

The novelty of this work is that it will give insights to the readers to understand the features of this technology to develop and deploy the Blockchain based applications in healthcare sector. In this work, applications of blockchain technology are discussed in depth focusing all the healthcare players which will attract the readers to trace the maximum potentials of the technology in the healthcare sector. Furthermore, highlighted challenges of the technology described by this study will set directions for future research.

## II. Literature Review

This section provides an overview of fundamental concepts related to the blockchain technology and several applications of this technology in the area of healthcare. Comparison of existing review papers has also been tabulated and discussed in this section.

### A. Blockchain

Blockchain, a distributed ledger, is a chain of time-stamped blocks containing a specific number of validated transactions. Blocks are linked cryptographically using the hash value of the previous block. Each transaction generated by a user or node is digitally signed using a private key and broadcasted to the network. A validation/mining node takes up that transaction and encloses it into a block then block is broadcasted to the network [14]. Each node of the network checks the validation of the block by implementing the consensus protocol. The validated block is appended to the chain then updated ledger is replicated throughout the permissioned nodes of the network. Consensus protocol replaces the trusted third party or the central authority. Fig. 1 illustrates the difference between centralized and distributed ledger [1] [15] [16] [17]. The ledger provides security, auditability and anonymity-based transparency.



Fig. 1. Centralized vs. Distributed Ledger

*1) Blockchain Evolution:* Melanie Swan [15] categorizes the blockchain technology evolution into three phases: Blockchain 1.0, 2.0, and 3.0. Blockchain 1.0 is for the decentralization of money or known as "Internet of Money". This first application established a peer to peer digital payment systems without reliance on a third party. This tier of technology implements Proof of work, the consensus protocol for validating a block to embed in the ledger [4]. A digital

reward is given to the successful miner for mining the block, for his contribution to the ledger. The second generation is blockchain 2.0 which is the application of decentralization of smart property and smart contracts, came in 2014. It aims to transfer any unit of value using the concepts of smart contracts for automated administration and supervision. The smart contract is a script which triggers after meeting the conditions encoded within it [17] [13]. Ethereum and Eris blockchains come under this category. The third blockchain technology refers to blockchain 3.0 which targets the welfare of society and is particularly recommended to register and transfer public records in the areas of government, health, science, literacy, and art. Examples are healthcoin, learningcoin, gridcoin, etc. [15]. Blockchain 2.0 and 3.0 are also known as non-financial applications. Alternative consensus protocols and alternative crypto-currencies have also been introduced.

*2) Blockchain Ownership:* There are two basic types of blockchain i.e. Permissioned blockchains and Permissionless blockchains. Permissioned blockchain is a custom-built setup by a single authority or a consortium. The verification process can be done by a central authority or a set of trusted preselected parties (consortium). This private setup restricts data access to the group of users or a set of groups that controls the blockchain. A smaller number of participants provides efficiency and scalability [18] [12]. These blockchains ultimately have a central authority. This centralization of the setup can pave way for tampering as the 51% majority is required to get the consensus and it can be done easily in this controlled setup [5]. Eris, Ripple, and Hyperledger are examples [1]. Permissionless blockchains are fully decentralized to a large number of nodes and low in efficiency [1]. These blockchains require no prior authorization of participants for mining the transaction blocks. Anyone can contribute his/her computational power for network tasks and can get a monetary reward in return. This blockchain gives the public access to read and write transactions to the blockchain which is visible to everyone so also known as Public blockchains [18] [12]. Examples of permissionless blockchains include Bitcoin and Ethereum [4].

*3) Features of blockchain technology:*
**Decentralization:** The blockchain is a distributed digital ledger composed of a chain of blocks containing transactions. The decentralized database is shared and open to all parties throughout the nodes of the network [6], [19], [20], [12]. BLockchain-based networks provide fault-tolerant architecture as end-to-end replications remove the reliance on a single point of failure.

**Consensus mechanism:** Blockchain is a peer-to-peer distributed network without any intermediary. Each digitally signed block is sent to the mining pool where it is taken over by network nodes called miners and verified using the consensus algorithm [4]. The winner from the miners broadcast the block to all other nodes which confirm and validate the block with consensus and append the block in their ledger. The winner also gets a financial reward for its work [21]. Many alternative consensus protocols have also been proposed, such as proof-of-stake, proof-of-burn, proof-of elapsed-time [6], [12], [17]. Data integrity is maintained by computing these consensus algorithms as a substitute to third trusted party.

**Immutability:** Blockchain is immutable and tamper-proof thus ultimately provides security [15]. The hash function

makes the blockchain as a tamper-resistant ledger. A hash value is calculated by implementing some hashing algorithms (SHA-256, RSA, RIPEMD-160, etc.) on a block of transactions [20]. The hash value is further used to create a chain of blocks. In this way blockchain provides robustness. If someone intends to alter previous transactions, then it will require a change in the hash value which further needs the consensus of network and high computational power which is an unrealistic approach in this computational model. The hash value is also used to represent a user concealing real identity which is used for privacy purposes [17], [22].

**Traceability:** Blockchain is a digital ledger consisting of continuously growing sequence of blocks. A block is comprised of a complete list of transaction records. In this chain of blocks, every block has a parent block. The first block in the chain is known as the genesis block. Hash code of genesis block is added to the header of the second block then hash code for the second block is computed over the hash of genesis block and transactions of the block jointly. Hash of the second block becomes the block header of the third block and so on. In this way, the blocks are linked with each other having a time stamp as well. This link can be chased back to the origin or genesis block [1], [20]. This feature of blockchain provides data provenance to keep chronical track of activities and may also help to investigate backward throughout the chain.

**Smart Contract:** A smart contract is a computerized computational logic or terms of the contract. It automatically triggers transactions between parties after fulfilling encoded logic. This implementation makes the blockchain flexible and programmable [1], [19] [13]. Smart contracts are programmed for management and administration purposes [6]. The smart contract can be implemented in supply chains, claim insurance [6] and clinical trials. Clinical trials usually encompass a sequence of dependent phases to get specific outcomes. Each phase can be encoded in a smart contract which will be triggered after the consensus of network nodes [23], [22]. So, smart contracts can enforce traceability and transparency with full control over associated processes.

**Open Source:** Blockchain projects are mostly open source. Developers can make contributions to it. Blockchain technology has the potential to accommodate the evolution to be brought by the future [8]. Transformation of financial blockchain to non-financial block-chains has already been announced in the big interest of the community. Ethereum and Counter-party have profound interest to build up more value-added services for the future architecture of blockchain technology [15].

**Currency:** The Bitcoin blockchain was first implemented in peer to peer digital payment system which also provides rewards in digital currency to the users for their contribution i.e. Proof of Work, to the network [6], [4]. Bitcoin was the first digital currency. In this evolution, many alt-crypto-currencies have been springing up. The major alternative currencies (alt-currencies) are monitored at http://coinmarketcap.com/. Over 1597 crypto-currencies have been developed in the digital cash system (Bitcoin, ethereum, ripple, litecoin, etc.). In different fields of life an economic layer can be embedded to give rewards in response to digital asset contribution and use e.g. learningcoin in education systems, journalcoin for the

contribution of reviewers and editors, healthcoin to get national health services, etc. [15].

### B. Use of Blockchain in Healthcare

Healthcare sector always remains one of the most popular areas of research from the last few decades keep on finding innovative and more reliable ways to help the community and healthcare industry. Different stakeholders (practitioners, medical specialists, hospitals, therapists, patients, payers, etc.) need to organize, access and share health records without any modification in a secure and interoperable way. Data provenance is also essential to prove the authenticity of records. Blockchain technology is being implemented in different scenarios and has the potentials to address the key issues of the healthcare sector. However, it needs more research to be focused to deploy real-time applications of this technology. Following are some applications of this technology in the healthcare sector.

MedRec platform provides decentralized record management, authorization and data sharing among healthcare stakeholders. Patients can save their data and can also grant and withdraw permissions to their records. This framework provides full confidentiality as the records are not stored on blockchain instead pointers to the data storage locations, logs and permissions are only stored in this blockchain [24], [4]. Gem in collaboration with Philips Blockchain Lab has been introduced Gem Health Network using Ethereum blockchain. This framework is established to address operational costs. This shared infrastructure provides interoperability among various providers accessing the same information to boost better patient care [8]. Guardtime healthcare platform creates a non-intermediated relationship between patient and provider in Estonia. Guardtime blockchain enabled transparent information sharing among the patient, provider, and payer which promises secure, reliable and auditable records [8]. Patient's health data is being required by research organizations. In this context, Healthbank has been providing a platform for patients to save and share their health data with research organizations which can be used for academic research and pharmaceuticals. This platform is also incentivizing patients with financial rewards for their contributions [8]. [25] designed Blockchain based Data Sharing (BBDS) access control system using permission blockchain. Owners of data can access their EMRs from a shared data pool. This secure and scalable system identifies, authenticates and authorizes users using cryptographic keys and digital signatures acquiring an edge over HDG (Healthcare Data Gateways) which is a smartphone application built over blockchain cloud [26]. Fast Healthcare Interoperability Resources: FHIRchain [27] was developed by the Health Level Seven International (HL7) organization for exchanging clinical data. FHIR increases efficiency and interoperability.

### C. Compared Secondary Studies

Blockchain has been appeared a decade before in computing while it first came in the healthcare sector in 2014 with the advent of the non-financial version of the technology. Researchers are found enthusiastic to explore this unique technology to know the potentials and challenges. In this regard, six secondary studies were found that discuss the implementation of this technology is the healthcare sector.

TABLE I.     COMPARED SECONDARY STUDIES

| | Discussion Points | | | | | |
|---|---|---|---|---|---|---|
| **References** | **Features of BC** | **Benefits of BC in HC (All Stakeholders)** | **Challenges and Issues to BC Implementation** | **BC Applications in HC** | **Research Methodology** | **Cloud-based BC applications and platform in HC** |
| T.Kuo et.al. [6], K. Rabah [9], S. Angraal [28] | Partially Yes | No | Yes | Partially Yes | No | No |
| M. Mettler [8] | Partially Yes | No | No | Partially Yes | No | No |
| M. Benchoufi et. al. [23] | Partially Yes | No | No | Partially Yes | No | No |
| D. Randall et. al. [29] | Partially Yes | No | No | No | No | No |
| Y.Sobia et. al. | Yes | Yes | Yes | Yes | Yes | Yes |

Table I presents five different aspects which have been reviewed by existing researchers and also compares to this study.

The first aspect in Table I is "features of blockchain technology". All the secondary studies partially narrate the features of blockchain as compared to this research. The second aspect discussed by earlier studies is "benefits of the blockchain technology in healthcare". It is extracted that M. Mettler [8], K. Rabah [9], as well as this work, made an analysis of this aspect with respect to all the stakeholders of healthcare sector whereas other four review papers discussed few of the stakeholders. The third mentioned aspect, Challenges and Issues to BC Implementation, has been highlighted by T. Kuo et al. [6], K. Rabah [9], S. Angraal [28] and this work whereas M. Mettler [8], M. Benchoufi et al. [23] and D. Randall et al. [29] did not address this aspect of concern.

Rapid development of this technology is being observed in healthcare market however after comparing these reviews, it is found that few of primary studies have been included on the subject of the fourth aspect, BC Applications in HC, except D. Randall et al. [29]. Whereas this study classifies more primary studies to analyze Blockchain Applications in the healthcare sector. Furthermore, qualitative research methodology has been opted to extract results in existing research. Hence, we used a quantitative research approach making a distinction to non-structured review process i.e. systematic literature review, to classify the primary studies to make a deeper enclosure of all the aspects compared in Table I.

### III.   RESEARCH METHODOLOGY

A systematic literature review (SLR) is a type of review or secondary study which firstly defines specific research questions and then uses a well-defined methodology to collect, classify and extract all existing research to answer those questions [30] [31] [32]. Various guidelines are available for writing a Systematic Literature Review. However, the steps recommended by Barbara Kitchenham [30] [31] [32] are followed to conduct current research. This methodology is being followed by review papers published in high impact factor journals [33] [34]. This process has been formulated specifically for conducting systematic reviews for computing research.

### A.  Need of Conducting SLR

Table I presents the aspects discussed by existing research in accordance with the scope of this study. After analyzing earlier review papers, it is concluded that some of the aspects have been discussed briefly or ignored so far. Therefore, we focused to elaborate on the features of this technology to trace the maximum potentials in the healthcare sector. Furthermore, applications of blockchain technology are needed to be discussed in depth focusing all the healthcare players to fully transform the system. Thirdly, this systematic literature review executes a detailed and specialized sequence of activities to extract results which make the distinction to previous non-structured reviews.

### B.  Research Question (RQ) and Motivation

In response to the first phase of systematic literature review, the following research questions are formulated to cover the gaps found, as shown in Table II.

TABLE II.     RESEARCH QUESTION AND MOTIVATION

| Sr. # | **Research Question** | **Motivation** |
|---|---|---|
| 1 | What are the major issues pertaining to the Healthcare Stakeholders? | The objective is to highlight major issues obstructing the success of Healthcare sector. |
| 2 | What Blockchain features are used to resolve the identified issues? | The aim is to explore the emerging technology that resolves the pertinent issues and accelerate the said field. |
| 3 | What are the challenges and issues to Blockchain implementation? | The aim is to find out those Blockchain implementation issues that are still unaddressed. |

### C.  Search Strategy

Pursuing research questions, the following search queries were put to collect maximum literature for review:
"Issues in the Healthcare sector",
"Blockchain",
"Healthcare Blockchain"
"Systematic Review".
After careful exploration of different databases and journals, two hundred twenty-seven studies were collected.

*D. Inclusion and Exclusion Criteria*

During this phase, some studies were found to be exactly aligned with the research area i.e. Blockchain and Healthcare, while others were found to be partially or completely out of context. Inclusion criteria limited the studies to thirty-nine since their titles and mentioned keywords were found to be similar to searched keywords. Studies published by famous publishers and impact factor journals were included. First of all, we analyzed the abstract of each shortlisted study according to research questions, methodology and findings of these papers and categorized them accordingly. Remaining searched articles were excluded because they did not include searched keywords in their titles and abstracts. Duplicate, irrelevant studies and publications are written in other than the English language were also excluded.

*E. Classification Criteria*

The aim of this research is to analyze existing research work implementing applications of blockchain technologies in the healthcare sector. Shortlisted studies were classified according to research questions.

*F. Data Extraction*

After classifying all studies to be included in the systematic review, the next step was to extract and analyze the information contained in those studies. Data extraction tables were designed using the spreadsheet to collect all information needed to address research questions i.e.
- Issues in Healthcare
- Features of Blockchain addressing those issues
- Challenges of Blockchain implementation.

Data extraction tables were filled in for each included study; further, those tables were used to accumulate the information.

## IV. RESULTS

This section provides the details about research questions discussed under section III and further divided into three subsections. Section A identifies common issues of healthcare sector classified under different players of the said sector. Section B identifies blockchain features that can solve the current issues of the healthcare sector. Section C classifies the studies which highlight the issues and challenges of blockchain implementation which are to be fixed in the future.

*A. RQ1: What are the major issues pertaining to Healthcare Stakeholders?*

A system is a combination of some elements or objects which work together to make an effective output. Healthcare Sector is composed of five important players forming the ecosystem i.e. providers, patients, payers, supply chain bearer (manufacturers, vendors, pharmacy) and research organization. Fig. 2 provides an overview of the relationship between five major stakeholders in the healthcare sector. Some of the common relationships (triangles) are being identified in the following diagram where a common color of entities describes the relationship among them. Provider (hospital, doctor, specialist, etc.) plays a key role and usually acts as a third party. Each player is being encountered by some issues which are major

areas of concern. These issues are discussed below under each scenario:

**Scenario 1. Providers:** A provider is a primary player in the healthcare sector. Patient record management is essential for successful outcomes for providers and patients. But there are numerous challenges to manage and curate the patient records as mentioned in Table III. Implementation of Electronic Medical Records (EMRs) and Electronic Health Records (EHRs) is one of the major issues confronted by healthcare providers as it is analyzed to be costly in terms of time and money both. Administrative costs of the Medicaid incentive program are estimated at approximately $30 Billion i.e. 56% of total spend. It also includes deployment, maintenance and administrative costs [35], [36], [37]. The other major issue is fragmented patient records which increase poor communication among medical staff and a decrease in quality of care. Lack of interoperability standards to exchange health data among laboratories and hospitals is liable to higher overheads in time and resources [36] [38]. Furthermore, current document-centric and legacy IT systems in the healthcare sector are being run into a major issue relevant to the security and privacy of patients' data. Health data, maintained by a trusted third party in centralized databases, can be accessed by any malicious party, hacker, insider, or outsider [20], [16]. Providers act as a trusted third party for patients, payers, and pharmacies as shown in Fig. 2.

**Scenario 2. Patients:** A patient is the most important entity of the healthcare ecosystem. Patient health data is recorded by providers however patients do not have the rights to grant and revoke access to their medical records as per their desire [39]. The HIPAA (Health Insurance Portability and Accountability Act) Privacy Rules put restrictions on the usage of such information of patient's health. Patients are mostly concerned about the protection of their health data but their privacy is victimized in many ways. It is compromised when organization encrypts and decrypts the data. The other sources of collection of health-related personal data may include mobile apps, wearables, smart devices, etc. Social media networks are also found to be observed to collect users' personal data, actions, and habits without users' consent. Many private organizations collect, analyze and sell patients' personal health data to different companies for commercial benefits. Data may be used for drug marketing, research, public health care or some objectionable purposes [40], [41], [16]. Sharing of such personal data without the consent of a patient is the violation of HIPPA Privacy Rules [7], [3], [22]. Moreover, trusted third party is getting monetary benefits by selling patients' digital assets to outsiders as shown in Fig. 2 by using green color. Fragmented silos of patient health records or non-interoperability also leverage lack of communication between healthcare teams which results in poor quality of care, more time consumption and higher cost of reinvestigations [42].

**Scenario 3. Payers:** Insurance claim payments on behalf of patients by the payer (insurer or employer) are also needed to be verified from centric IT systems which are highly vulnerable to insecurity and intended fraudulent alterations [43]. Records (bills and prescriptions etc.) may be falsified by fake medical practitioner credentials, misbilling, and bogus testing, etc. Auditing and data provenance is crucial to track genuine insurance claim processes [44].

Fig. 2. Overview of the relationship among five major stakeholders in the healthcare sector (Suppliers, manufacturer, and pharmacies form a single unit of pharma supply chain)

**Scenario 4. Research Organizations:** Public health data is required by research organizations and pharmaceutical companies to track new diseases, invent their treatments and drug discovery. Health care data sharing is prerequisite for national data collection, clinical audit, and research [8]. Patient data recorded by a provider is shared with these organizations without the consent of patients, which is the violation of patient's privacy. Usually, patients are ignorant of the fact that their personal data is being shared with other parties without their consent which is the violation of patient's privacy and needed to be addressed [45]. Moreover, a chain of historical events describing exact switching points is crucial to track accelerated clinical trials and researches. Unavailability of data provenance contributes to undermining the reproducibility of outcomes. Maladministration lifts up the data snooping and misinterpreted reportage in disclosing the outcomes of trials in order to show the effects of new drugs discovered [23]. Patients' consent management, data management, data integrity, and transparent results are the major concerns of a trustful and productive system.

**Scenario 5. Pharmaceutical Supply Chain Management:** In the pharmaceutical supply chain, major players are suppliers, brokers, manufacturers, and pharmacies. It is a complex system including a diverse range of activities of acquisition of raw material, production, storage, distribution, etc. Proper management and monitoring are required to ensure reliability. Drug ingredients may travel through many sources where origin and quality of ingredients are not certified. Substandard constituents and illegitimate sources can play an important role in the production of counterfeit medicine [46]. So data provenance is obligatory to track genuine records, verified product information, and ownership in the drug supply chain. But current systems lack to maintain historical information to tackle the issues of counterfeit drugs [47].

**Scenario 6. Prescription Management at Pharmacies:**

Pharmacy is one of the supply chain stakeholders which forms a triangular relationship with provider and patient which is shown in Fig. 2 by using red color. Prescriptions forwarded by providers may be mishandled or un-deciphered at pharmacies and as a result, wrong medicine may be dispensed to a patient or it can also be given to an un-intended patient when prescriptions are jumbled up [48] [47]. The National Health Service (NHS)'s solution of electronic prescriptions rectifies many of the problems of the manual system. However, electronic prescriptions sent to community pharmacies from general practitioners expose patient confidentiality and need intense prescription management [49]. Moreover, main logistics inefficiencies executed by the hospital pharmacy may result in a big threat to the patients' safety. These inefficient activities may include incorrect inventory management, medicine shortage, long procurement cycles, time-consuming product recalls and improper use of technology [50]. Record management, data exchange among other supply chain stakeholders, security of records and privacy violation of patients' data are the major issues exhibited by this player of the healthcare sector.

In above-mentioned scenarios, some issues are observed to be common for all healthcare stakeholders. Major issues of each player are tabulated in Table III.

*B. RQ2: What Blockchain features are used to resolve the identified issues?*

Blockchain technologies have applicability in different scenarios which are the key benefit of this technology. The rapid advancement of healthcare sector can be envisioned by implementing the following features of blockchain technology, combating the major issues of the sector:

**Distributed Digital Ledger resolves the issue of record management:** Record management is inevitable to the success of every organization, which requires intense resource consumption in terms of human resources as well as software and hardware. Blockchain, the distributed ledger technology with no central authority promises to transform the current costly systems to inexpensive and easy to implement systems with higher efficiency and productivity [29] [17]. Decentralized communication among all medical stakeholders can be achieved as the same information is replicated throughout the permissioned nodes without any intermediary. Controlled access and real-time updates in records can increase fraud detection, combatting counterfeit drugs, claim adjudication and verified reproducibility in research [6]. The proposed work [45] can altogether lessen the turnaround time for EMR sharing intended for research and supervision of patient care, enhances basic management and reduces the general cost.

**Interoperability resolves the issue of data exchange:** Blockchain distributed ledger technology can resolve the issue of incompatible fragmented patient health records resulting in improved coordination and quality of care. This synchronized infrastructure will fasten real-time patients' notification, latest health and treatment information sharing and faster product innovation as the exchange of data will enhance market competitions positively [10] [51] [52] [38]. It helps in avoiding a single point of failure and it can facilitate in medical research by providing reliable data to different institutes so that better solution of patient care can be discovered [53]. Claim processing is usually affected by the complications of traditional

TABLE III.    COMMON MAJOR ISSUES OF ALL HEALTHCARE STAKEHOLDERS

| Issues | Healthcare Stakeholders | | | | |
|---|---|---|---|---|---|
| | **Provider** | **Patient** | **Payer** | **Research Organizations** | **Supply chain** |
| Record Management | Yes | Yes | Yes | Yes | Yes |
| Data Exchange | Yes | Yes | Yes | Yes | Yes |
| Security | Yes | Yes | Yes | No | Yes |
| Lack of Data Provenance | No | Yes | Yes | Yes | Yes |
| Privacy | No | Yes | No | No | No |
| Monetization | No | Yes | No | No | No |

distributed records, can be addressed by blockchain technology's interoperability feature [6] [46]. Shared immutable prescriptions can enhance medication reconciliation and quality of care.

**Consensus mechanism and cryptography resolves the issue of security and privacy:** Auditability and transparency is the magic demonstrated by immutable blockchain that has the ability to captivate the attention of users. Patient-controlled secure access is guaranteed by implementing consensus protocol and cryptographic keys in blockchain technology as the only private key can decrypt the data. Patient owns his/her data and also has the right to grant and revoke access to other persons (providers, payers or researchers) [6] [39] [17] [22] [5]. Security is one of the major issues encountered by almost every organization due to the reliance on a central authority. However, alteration in transactions in a blockchain based system will require the consensus of network and high computational power so data tampering is restricted.

**Traceability and time stamping resolves the issue of lack of data provenance:** Time stamped, verified records for claim qualification can save payers from hazardous financial loss. Findings of clinical trials require clarity, no data snooping; accurate endpoint switching, etc. These issues can be conquered by traceability feature of blockchain that provides the historical information by maintaining a chain of time-stamped blocks [23] [45] [22]. Drug manufacturing and distribution process can be tracked to detect counterfeit medicines using blockchain as it provides data provenance.

**Digital Currency resolves the issue of monetization:** In a blockchain based network, miners are rewarded with cryptocurrency as an incentive for consuming their computation powers and to serve and run the system. The cost has to be determined when a DApp provides services for patient and provider [28] [22]. Healthbank [8] has been providing a platform for patients and research organizations. They give financial rewards to the patient for their contributions and use health data for research and academic purposes.

*C. RQ3: What are the challenges of Blockchain technology after its implementation in the Healthcare sector?*

This new born technology has the potential to accelerate the healthcare sector however some implementation challenges are also discussed below:

**Scalability:** Blockchain powered healthcare system when uses sensor devices for patient care, faces storage limitation problem and requires heavy load computation that needs to be resolved [18]. Moreover, Blockchain network executes the transactions too slow and consequently system results in slow throughput as it may take days to complete a single transaction [54]. The transaction time is very long i.e. Bitcoin blockchain executes seven transactions per second (1 MB block size), depending upon the protocol (proof of work), in comparison with the Visa and Twitter networks, which perform 2000 and 5000 transactions per second respectively. Such a speed issue may limit the scalability of the blockchain network [6], [54]. All blocks are stored on every node which exists in the distributed system of the blockchain, also creates the issue of speed and scalability [07]. Real-time blockchain based health care applications may suffer from speed and scalability issue in a larger network.

**Security (threat of 51% attack):** Blockchain runs on consensus. Majority or more than half of total nodes i.e. 51%, malicious miners may occupy the network and reject the blocks from other honest miners, their greater computation power may seize the precious information or coins [6] [54] [5]. However, there are fewer chances of this threat in a larger network.

**Disclosure of confidentiality:** Open source nature of the blockchain database pulls out another limitation which is "transparency discloses confidentiality". It is more critical for healthcare records of patients and biomedical applications because patients related records are highly sensitive [6]. (due to transparent nature of blockchain, healthcare or other's sensitive data are shown to everyone on the blockchain)

**Anonymity and data privacy:** Another challenge to blockchain is that the uncertainty and fraud hype surrounds it [20]. Criminals may use cryptocurrencies taking advantage of anonymity under the blockchain network. People may buy illegal drugs on the Dark web by using cryptocurrencies. Fraud

hype is also associated with blockchain networks when hackers use "Ransomware" to seize computer networks and demand payment in cryptocurrencies [39]. Data privacy is compromised in a public blockchain network [55] [57] whereas patient is concerned about his/her privacy of health data. Moreover, blockchain based cloud environment exhibit limitations on access control methods for privacy and security of health data [56].

**Environmental unsustainability:** Environmentally unsustainability and inefficiency is another issue with current blockchain implementations. The "proof-of-work" requirements in current versions of the blockchain require massive amounts of electricity; the energy cost of a single Bitcoin transaction could power 1.5 American homes for a day [43]. As the ledgers get longer, the math gets harder, and the amount of power being used increases. Thankfully, there are less computing-intensive versions of the blockchain in development [39]. But the associated challenge of storage intensive medical records are needed to be resolved.

## V. Threats to Validity

This SLR tries to compare and classify the blockchain technologies for the healthcare sector. Systematic literature reviews are considered reliable in general but this review can have some potential limitations. Expected limitations are restricted to related studies, identification and selection, insufficient data extraction, and unconcluded results.

### A. Threats to Identification and Selection of Primary Studies

To provide a deeper insight to use blockchain in the healthcare sector, we try to gather as many primary studies as possible for extraction of knowledge to avoid biases. But as this new domain is in its early stages and researchers are exploring it very keenly and eagerly, so in the process of publication, many related research studies may become available in near future which is not presented here. A classification criterion is designed to shortlist and appropriately classifies as many primary studies as possible. We included all the related studies and did not evaluate these studies by giving quality scores.

### B. Threats to Data Extraction

We gathered as many articles as possible related to our domain including primary and secondary articles, no quality score was assigned to select studies which help in prioritizing the result outcomes and research trends, is a major threat to data extraction. Another threat is that the data extracted from these articles based on our perspective of research questions and motivations. There are clear chances that readers and researchers may find some points that need to be considered to make this study more influential.

### C. Threats to Synthesis and Results

As quality score is not assigned to collected studies it may lead to less quality results and synthesis.

## VI. Conclusion and Future Work

Blockchain Technology is relatively new in the field of computing and healthcare as well. This technology has great potential in the sub-sectors of the healthcare field solving the major issues with its features and properties. Technology has the potential to revolutionize the whole ecosystem. Providers, patients, and research organizations are more focused on its initial journey and needs more research however the intense research work must be carried out in health insurance and pharmaceutical supply chains. Blockchain technology is also observed with some challenges while an implementation that needs to be solved with further research. Threats to the validity of our study discussed above may lead to better future work of this research work.

## References

[1] Zheng, Zibin, Shaoan Xie, Hongning Dai, Xiangping Chen, and Huaimin Wang. "An overview of blockchain technology: Architecture, consensus, and future trends." In 2017 IEEE International Congress on Big Data (BigData Congress), pp. 557-564. IEEE, 2017.

[2] Nakamoto, Satoshi. "Bitcoin: A peer-to-peer electronic cash system," http://bitcoin. org/bitcoin. pdf." (2008).

[3] Kuo, Tsung-Ting, Hyeon-Eui Kim, and Lucila Ohno-Machado. "Blockchain distributed ledger technologies for biomedical and health care applications." Journal of the American Medical Informatics Association 24, no. 6 (2017): 1211-1220.

[4] Daniel, Jeff, Arman Sargolzaei, Mohammed Abdelghani, Saman Sargolzaei, and Ben Amaba. "Blockchain Technology, Cognitive Computing, and Healthcare Innovations." Journal of Advances in Information Technology Vol 8, no. 3 (2017).

[5] Mettler, Matthias. "Blockchain technology in healthcare: The revolution starts here." In 2016 IEEE 18th International Conference on e-Health Networking, Applications and Services (Healthcom), pp. 1-3. IEEE, 2016.

[6] Rabah, Kefa. "Challenges & opportunities for blockchain powered healthcare systems: A review." Mara Research Journal of Medicine & Health Sciences-ISSN 2523-5680 1, no. 1 (2017): 45-52.

[7] Zhang, Peng, Michael A. Walker, Jules White, Douglas C. Schmidt, and Gunther Lenz. "Metrics for assessing blockchain-based healthcare decentralized apps." In 2017 IEEE 19th International Conference on e-Health Networking, Applications and Services (Healthcom), pp. 1-4. IEEE, 2017.

[8] Park, Jin, and Jong Park. "Blockchain security in cloud computing: Use cases, challenges, and solutions." Symmetry9, no. 8 (2017): 164.

[9] Ahram, Tareq, Arman Sargolzaei, Saman Sargolzaei, Jeff Daniels, and Ben Amaba. "Blockchain technology innovations." In 2017 IEEE Technology & Engineering Management Conference (TEMSCON), pp. 137-141. IEEE, 2017.

[10] Kuo, Tsung-Ting, Hugo Zavaleta Rojas, and Lucila Ohno-Machado. "Comparison of blockchain platforms: a systematic review and healthcare examples." Journal of the American Medical Informatics Association 26, no. 5 (2019): 462-478.

[11] Risius, Marten, and Kai Spohrer. "A blockchain research framework." Business & Information Systems Engineering 59, no. 6 (2017): 385-409.

[12] 45. Hölbl, Marko, Marko Kompara, Aida Kamišalić, and Lili Nemec Zlatolas. "A systematic review of the use of blockchain in healthcare." Symmetry 10, no. 10 (2018): 470.

[13] Xia, Q. I., Emmanuel Boateng Sifah, Kwame Omono Asamoah, Jianbin Gao, Xiaojiang Du, and Mohsen Guizani. "MeDShare: Trust-less medical data sharing among cloud service providers via blockchain." IEEE Access 5 (2017): 14757-14767.

[14] Banerjee, Mandrita, Junghee Lee, and Kim-Kwang Raymond Choo. "A blockchain future for internet of things security: A position paper." Digital Communications and Networks 4, no. 3 (2018): 149-160.

[15] Swan, Melanie. Blockchain: Blueprint for a new economy. " O'Reilly Media, Inc.", 2015.

[16] Zhang, Peng, Douglas C. Schmidt, Jules White, and Gunther Lenz. "Blockchain technology use cases in healthcare." In Advances in Computers, vol. 111, pp. 1-41. Elsevier, 2018.

[17] Dwivedi, Ashutosh Dhar, Gautam Srivastava, Shalini Dhar, and Rajani Singh. "A decentralized privacy-preserving healthcare blockchain for iot." Sensors 19, no. 2 (2019): 326.

[18] Peters, Gareth W., and Efstathios Panayi. "Understanding modern banking ledgers through blockchain technologies: Future of transaction processing and smart contracts on the internet of money." In Banking beyond banks and money, pp. 239-278. Springer, Cham, 2016.

[19] Mamoshina, Polina, Lucy Ojomoko, Yury Yanovich, Alex Ostrovski, Alex Botezatu, Pavel Prikhodko, Eugene Izumchenko et al. "Converging blockchain and next-generation artificial intelligence technologies to decentralize and accelerate biomedical research and healthcare." Oncotarget9, no. 5 (2018): 5665.

[20] Tama, Bayu Adhi, Bruno Joachim Kweka, Youngho Park, and Kyung-Hyune Rhee. "A critical review of blockchain and its current applications." In 2017 International Conference on Electrical Engineering and Computer Science (ICECOS), pp. 109-113. IEEE, 2017.

[21] Nugent, Timothy, David Upton, and Mihai Cimpoesu. "Improving data transparency in clinical trials using blockchain smart contracts." F1000Research 5 (2016).

[22] Griggs, Kristen N., Olya Ossipova, Christopher P. Kohlios, Alessandro N. Baccarini, Emily A. Howson, and Thaier Hayajneh. "Healthcare blockchain system using smart contracts for secure automated remote patient monitoring." Journal of medical systems 42, no. 7 (2018): 130.

[23] Benchoufi, Mehdi, and Philippe Ravaud. "Blockchain technology for improving clinical research quality." Trials 18, no. 1 (2017): 335.

[24] Azaria, Asaph, Ariel Ekblaw, Thiago Vieira, and Andrew Lippman. "Medrec: Using blockchain for medical data access and permission management." In 2016 2nd International Conference on Open and Big Data (OBD), pp. 25-30. IEEE, 2016.

[25] Xia, Qi, Emmanuel Sifah, Abla Smahi, Sandro Amofa, and Xiaosong Zhang. "BBDS: Blockchain-based data sharing for electronic medical records in cloud environments." Information 8, no. 2 (2017): 44.

[26] Yue, Xiao, Huiju Wang, Dawei Jin, Mingqiang Li, and Wei Jiang. "Healthcare data gateways: found healthcare intelligence on blockchain with novel privacy risk control." Journal of medical systems40, no. 10 (2016): 218.

[27] Zhang, Peng, Jules White, Douglas C. Schmidt, Gunther Lenz, and S. Trent Rosenbloom. "Fhirchain: applying blockchain to securely and scalably share clinical data." Computational and structural biotechnology journal 16 (2018): 267-278.

[28] Angraal, Suveen, Harlan M. Krumholz, and Wade L. Schulz. "Blockchain technology: applications in health care." Circulation: Cardiovascular Quality and Outcomes 10, no. 9 (2017): e003800.

[29] Randall, David, Pradeep Goel, and Ramzi Abujamra. "Blockchain applications and use cases in health information technology." J Health Med Informat 8, no. 276 (2017): 2.

[30] Kitchenham, Barbara, O. Pearl Brereton, David Budgen, Mark Turner, John Bailey, and Stephen Linkman. "Systematic literature reviews in software engineering–a systematic literature review." Information and software technology 51, no. 1 (2009): 7-15.

[31] Khan, Muhammad Murad, Roliana Ibrahim, and Imran Ghani. "Cross domain recommender systems: a systematic literature review." ACM Computing Surveys (CSUR) 50, no. 3 (2017): 36.

[32] Yumna, Hafiza, Muhammad Murad Khan, Maria Ikram, and Sabahat Ilyas. "Use of Blockchain in Education: A Systematic Literature Review." In Asian Conference on Intelligent Information and Database Systems, pp. 191-202. Springer, Cham, 2019.

[33] Aly, Mohab, Foutse Khomh, Mohamed Haoues, Alejandro Quintero, and Soumaya Yacout. "Enforcing Security in Internet of Things Frameworks: A Systematic Literature Review." Internet of Things (2019): 100050.

[34] Alomar, Noura, Mansour Alsaleh, and Abdulrahman Alarifi. "Social authentication applications, attacks, defense strategies and future research directions: a systematic review." IEEE Communications Surveys & Tutorials 19, no. 2 (2017): 1080-1111.

[35] Till, Brian M., Alexander W. Peters, Salim Afshar, and John G. Meara. "From blockchain technology to global health equity: can cryptocurren-cies finance universal health coverage?." BMJ global health 2, no. 4 (2017): e000570.

[36] Anderson, James G. "Social, ethical and legal barriers to e-health." International journal of medical informatics 76, no. 5-6 (2007): 480-483.

[37] Lau, E. "Decoding the hype: Blockchain in Healthcare-A Software Architecture for the provision of a patient summary to overcome interoperability issues." Master's thesis, 2018.

[38] Liang, Xueping, Juan Zhao, Sachin Shetty, Jihong Liu, and Danyi Li. "Integrating blockchain for data sharing and collaboration in mobile healthcare applications." In 2017 IEEE 28th Annual International Symposium on Personal, Indoor, and Mobile Radio Communications (PIMRC), pp. 1-5. IEEE, 2017.

[39] Hoy, Matthew B. "An introduction to the blockchain and its implications for libraries and medicine." Medical reference services quarterly 36, no. 3 (2017): 273-279.

[40] Karafiloski, Elena, and Anastas Mishev. "Blockchain solutions for big data challenges: A literature review." In IEEE EUROCON 2017-17th International Conference on Smart Technologies, pp. 763-768. IEEE, 2017.

[41] Kaplan, Bonnie. "Selling health data: de-identification, privacy, and speech." Cambridge Quarterly of Healthcare Ethics 24, no. 3 (2015): 256-271.

[42] Skiba, Diane J. "The potential of Blockchain in education and health care." Nursing education perspectives 38, no. 4 (2017): 220-221.

[43] Brodersen, C., B. Kalis, C. Leong, E. Mitchell, E. Pupo, A. Truscott, and L. Accenture. "Blockchain: Securing a New Health Interoperability Experience." Accenture LLP (2016).

[44] Engelhardt, Mark A. "Hitching healthcare to the chain: An introduction to blockchain technology in the healthcare sector." Technology Innovation Management Review 7, no. 10 (2017).

[45] Dubovitskaya, Alevtina, Zhigang Xu, Samuel Ryu, Michael Schumacher, and Fusheng Wang. "Secure and trustable electronic medical records sharing using blockchain." In AMIA Annual Symposium Proceedings, vol. 2017, p. 650. American Medical Informatics Association, 2017.

[46] Mackey, Tim K., and Gaurvika Nayyar. "A review of existing and emerging digital technologies to combat the global trade in fake medicines." Expert opinion on drug safety 16, no. 5 (2017): 587-602.

[47] Clauson, Kevin A., Elizabeth A. Breeden, Cameron Davidson, and Timothy K. Mackey. "Leveraging blockchain technology to enhance supply chain management in healthcare." Blockchain in Healthcare Today (2018).

[48] Huckvale, Christopher, Josip Car, Masanori Akiyama, Safurah Jaafar, Tawfik Khoja, Ammar Bin Khalid, Aziz Sheikh, and Azeem Majeed. "Information technology for patient safety." BMJ Quality & Safety 19, no. Suppl 2 (2010): i25-i33.

[49] Porteous, Terry, Christine Bond, Roma Robertson, Philip Hannaford, and Ehud Reiter. "Electronic transfer of prescription-related information: comparing views of patients, general practitioners, and pharmacists." Br J Gen Pract 53, no. 488 (2003): 204-209.

[50] Romero, Alejandro. "Managing medicines in the hospital pharmacy: logistics inefficiencies." In Proceedings of the World Congress on Engineering and Computer Science, vol. 2, pp. 1-6. 2013.

[51] Raju, Saravanan, Vandita Rajesh, and Jitender S. Deogun. "The Case for a Data Bank: an Institution to Govern Healthcare and Education." In Proceedings of the 10th International Conference on Theory and Practice of Electronic Governance, pp. 538-539. ACM, 2017.

[52] Zhang, Peng, Jules White, Douglas C. Schmidt, and Gunther Lenz. "Design of blockchain-based apps using familiar software patterns with a healthcare focus." In Proceedings of the 24th Conference on Pattern Languages of Programs, p. 19. The Hillside Group, 2017.

[53] Kuo, Tsung-Ting, and Lucila Ohno-Machado. "Modelchain: Decentralized privacy-preserving healthcare predictive modeling framework on private blockchain networks." arXiv preprint arXiv:1802.01746 (2018).

[54] Lazar, Max A., Zihang Pan, Renee-Marie Ragguett, Yena Lee, Mehala Subramaniapillai, Rodrigo B. Mansur, Nelson Rodrigues, and Roger S. McIntyre. "Digital revolution in depression: a technologies update for clinicians." Personalized Medicine in Psychiatry 4 (2017): 1-6.

[55] Reyna, Ana, Cristian Martín, Jaime Chen, Enrique Soler, and Manuel Díaz. "On blockchain and its integration with IoT. Challenges and

opportunities." Future Generation Computer Systems 88 (2018): 173-190.

[56] Esposito, Christian, Alfredo De Santis, Genny Tortora, Henry Chang, and Kim-Kwang Raymond Choo. "Blockchain: A panacea for healthcare cloud-based data security and privacy?." IEEE Cloud Computing 5, no. 1 (2018): 31-37.

[57] Tasatanattakool, Pinyaphat, and Chian Techapanupreeda. "Blockchain: Challenges and applications." In 2018 International Conference on Information Networking (ICOIN), pp. 473-475. IEEE, 2018.

# Deep Gated Recurrent and Convolutional Network Hybrid Model for Univariate Time Series Classification

Nelly Elsayed[1], Anthony S Maida[2]
School of Computing and Informatics,
University of Louisiana at Lafayette,
Louisiana, USA

Magdy Bayoumi[3]
Dep. of Electrical and Computer Engineering,
University of Louisiana at Lafayette,
Louisiana, USA

*Abstract*—Hybrid LSTM-fully convolutional networks (LSTM-FCN) for time series classification have produced state-of-the-art classification results on univariate time series. We empirically show that replacing the LSTM with a gated recurrent unit (GRU) to create a GRU-fully convolutional network hybrid model (GRU-FCN) can offer even better performance on many time series datasets without further changes to the model. Our empirical study showed that the proposed GRU-FCN model also outperforms the state-of-the-art classification performance in many univariate time series datasets without additional supporting algorithms requirement. Furthermore, since the GRU uses simpler architecture than the LSTM, it has fewer training parameters, less training time, smaller memory storage requirements, and simpler hardware implementation, compared to the LSTM-based models.

*Keywords*—*GRU-FCN; LSTM; fully convolutional neural network; time series; classification*

## I. Introduction

A time series (TS) is a sequence of data points obtained at successive equally-spaced time points, ordinarily in a uniform interval time domain [1]. TSs are used in several research and industrial fields where temporal analysis measurements are involved such as in signal processing [2], pattern recognition [3], mathematics [1], psychological and physiological signals analysis [4], [5], earthquake prediction [6], weather readings [7], and statistics [1]. There are two types of time series: univariate and multivariate. In this paper, our objective is to study the univariate time series classification.

There are many approaches to time series classification. The distance-based classifier based on the k-nearest neighbor (KNN) algorithm is considered a baseline technique for time series classification. Mostly, a distance-based classifier uses Euclidean or Dynamic Time Warping (DTW) as a distance measure [8]. Feature-based time series classifiers are also widely used such as the bag-of-SFA-symbols (BOSS) [9] and the bag-of-features framework (TSBF) [10] classifiers. Ensemble-based classifiers combine separate classifiers into one model to reach a higher classification accuracy such as the elastic ensemble (PROP) [11], and the collective of transform-based ensemble (COTE) [12] classifiers.

Convolutional neural network (CNN) based classifiers have advantages over other classification methods because CNNs provide the classifier with a preprocessing mechanism within

TABLE I. COMPARISON OF GRU AND LSTM COMPUTATIONAL ELEMENTS.

| Comparison | LSTM | GRU |
|---|---|---|
| number of gates | 3 | 2 |
| number of activations | 2 | 1 |
| state memory cell | Yes | No |
| number of weight matrices | 8 | 6 |
| number of bias vectors | 3 | 4 |
| number of elementwise multiplies | 3 | 3 |
| number of matrix multiplies | 8 | 6 |

the model. Examples are the multi-channel CNN (MC-CNN) classifier [13], the multi-layered preceptron (MLP) [4], the fully convolutional network (FCN) [4] and, specifically, the residual network (ResNet) [4].

The present paper focuses on the recurrent neural network based classification approaches such as LSTM-FCN [5] and ALSTM-FCN [5] that are the first recurrent-based time series classification models. These models combine both temporal CNNs and long short-term memory (LSTM) models to provide the classifier with both feature extraction and time dependencies through the dataset during the classification process. These models use additional support algorithms such as attention and fine-tuning algorithms to enhance the LSTM learning due to its complex structure and data requirements.

This paper attempts to emerge the difference between the GRU and LSTM in univariate time series classification purpose. This paper studies whether the use of gated-recurrent units (GRUs) can improve the hybrid classifiers listed above with. We create the GRU-FCN by only replacing the LSTM with a GRU in the LSTM-FCN [5]. We intentionally kept the other components of the entire model without changes to make an empirical comparison between the LSTM and GRU in the same model structure to obtain a fair comparison between both architectures regarding the univariate time series classification task. Like the LSTM-FCN, our model does not require feature engineering or data preprocessing before the training or testing stages. The GRU is able to learn the temporal dependencies within the dataset. Moreover, the GRU has a smaller block architecture and shows comparable performance to the LSTM without a need for additional algorithms to support the model.

Although it is difficult to determine the best classifier for all time series types, the proposed model seeks to achieve equivalent accuracy to state-of-the-art classification models

Fig. 1. Block architecture for an unrolled GRU.

in univariate time series classification. Following [4] and [5], our tests use the UCR time series classification archive benchmark [14] to compare our model with other state-of-the-art univariate time series classification models. Our model achieved higher classification performance on several datasets compared to other state-of-the-art classification models.

## II. Model Components

### A. Gated Recurrent Unit (GRU)

The gated recurrent unit (GRU) was introduced in [15] as another type of gate-based recurrent unit which has a smaller architecture and comparable performance to the LSTM unit. The GRU consists of two gates: reset and update. The architecture of an unrolled GRU block is shown in Fig. 1. $r^{(t)}$ and $z^{(t)}$ denote the values of the reset and update gates at time step $t$, respectively. $x_i \in \mathbb{R}^n$ is a 1D input vector to the GRU block at time step $t$. $\tilde{h}^{(t)}$ is the output candidate of the GRU block. $h^{(t-1)}$ is the recurrent GRU block output of time step $t-1$ and the current output at time $t$ is $h^{(t)}$. Assuming a one-layer GRU, the reset gate, update gate, output candidate, and GRU output are calculated as follows [15]:

$$z^{(t)} = \sigma(W_{zx}x^{(t)} + U_{zh}h^{(t-1)} + b_z) \tag{1}$$

$$r^{(t)} = \sigma(W_{rx}x^{(t)} + U_{rh}h^{(t-1)} + b_r) \tag{2}$$

$$\tilde{h}^{(t)} = \tanh(W_x x^{(t)} + U_h(r^{(t)} \odot h^{(t-1)}) + b) \tag{3}$$

$$h^{(t)} = (1 - z^{(t)}) \odot h^{(t-1)} + z^{(t)} \odot \tilde{h}^{(t)} \tag{4}$$

where $W_{zx}$, $W_{rx}$, and $W_x$ are the feedforward weights and $U_{hz}$, $U_{hr}$, and $U_h$ are the recurrent weights of the update gate, reset gate, and output candidate activation respectively. $b_z$, $b_r$ and $b$ are the biases of the update gate, reset gate and the output candidate activation $\tilde{h}^{(t)}$, respectively. Fig. 3 shows the GRU architecture with weights and biases made explicit.

Like the RNN and LSTM, the GRU models temporal (sequential) datasets. The GRU uses its previous time step output and current input to calculate the next output. The GRU has the advantage of a smaller size over the LSTM. The GRU consists of two gates (reset and update), while the LSTM has three gates: input, output and forget. The GRU has one unit activation, but the LSTM has two unit activations:



Fig. 2. The proposed GRU-FCN model architecture diagram rendered using the Keras visualization tool and modified from [4], [5] architectures.

Fig. 3. The GRU architecture showing the weights of each component.

input-update and output activations. Also, the GRU does not contain the memory state cell which exists in the LSTM model. Thus, the GRU requires fewer trainable parameters, and shorter training time compared to the LSTM. Table I compares GRU and LSTM architecture components.

### B. Temporal Convolutional Neural Network

The Convolutional Neural Network (CNN), introduced in 1989 [16], utilizes weight sharing over grid-structured datasets such as images and time series [17], [18]. The convolutional layers within the CNN learn to extract complex feature representations from the data with little or no preprocessing. The temporal FCN consists of many layers of convolutional blocks that may have different or same kernel sizes, followed by a dense layer softmax classifier [18]. For time series problems, the values of each convolutional block in the FCN, are calculated as follows [4]:

$$y_i = W_i * x_i + b_i \qquad (5)$$
$$z_i = BN(y) \qquad (6)$$
$$out_i = ReLU(z) \qquad (7)$$

where $x_i \in \mathbb{R}^n$ is a 1D input vector which represents a time series segment, $W_i$ is the 1D convolutional kernel of weights, $b_i$ is the bias, and $y$ is the output vector of the convolutional block $i$. $z_i$ is the intermediate result after applying batch normalization [19] on the convolutional block which then is passed to the rectified linear unit $ReLU$ [20] to calculate the output of the convolutional layer $out_i$.

### III. MODEL ARCHITECTURE

As stated in the introduction, our model replaces the LSTM with a GRU in a hybrid gated-FCN. We intentionally did not change the other components of the entire model to attain a fair comparison between GRU and LSTM architectures in the same model structure for univariate time series classification. Our model is based on the framework introduced in [4], [5]. The proposed architecture actual implementation is shown in Fig. 2. The architecture has two parallel parts: a GRU and a temporal FCN. Our model uses three-layered FCN architecture proposed in [4]. The dimension adjustment aims to change the dimensions of the input to be compatible with the GRU recurrent

design [21]. We also used the global average pooling layer [22] to interpret the classes and to reduce the number of trainable parameters comparing to the fully connected layer, without any sacrifice in the accuracy. The FCN 1D kernel numbers are 128, 256, and 128 with kernel sizes 8, 5, and 3 in each convolutional layer, respectively. The weights were initialized using the He uniform variance scaling initializer [23]. In addition, we used the GRU instead of LSTMs that were used in [5] models to reduce the number of trainable parameters, memory, and training time. Moreover, we removed the masking and any extra supporting algorithms such as an attention mechanism, and fine-tuning that were used in the LSTM-FCN and ALSTM-FCN models [5]. The GRU is unfolded by eight unfolds as used in [5] for univariate time series. The hyperbolic tangent ($tanh$) function used as the unit activation and the hard-sigmoid ($hardSig$) function [24] is used as the recurrent activation (gate activation) of the GRU architecture. The weights were initialized using the $glorot\_uniform$ initializer [25], [26] and the biases were initialized to zero. The input was fitted using the concept used in [5] to fit an input to a recurrent unit. We used the Adam optimization function [27] with $\beta_1 = 0.9$, $\beta_2 = 0.999$ and initial learning rate $\alpha = 0.01$. The learning rate $\alpha$ was reduced by a factor of 0.8 every 100 training steps until it reached the minimum rate $\alpha = 0.0001$. The dense layer uses the softmax classifier [28] using the categorical crossentropy loss function [18]. In this paper, our goal is to make a fair comparison between the LSTM-based model and our GRU-based model. Thus, we used the same number of epochs that were assigned by the original LSTM-FCN model [5] for each univariate time series. The number of epochs that we assigned for each dataset used is shown in Table II.

The input to the model is the raw dataset without applying any normalizations or feature engineering prior to the training process. The FCN is responsible for feature extraction from the time series [4] and the GRU enables the model to learn temporal dependencies within the time series. Therefore the model learns both the features and temporal dependencies to predict the correct class for each training example.

### IV. METHOD AND RESULTS

We implemented our model by modifying the original LSTM-FCN [5]. We found that the fine-tuning algorithm has not been applied in the actual LSTM-FCN and ALSTM-FCN implementation source code on Github which shared by the authors [5] and mentioned in their literature. In addition, the LSTM-FCN [5] authors used a permutation algorithm for fitting the input to the FCN part which was not mentioned in their literature. Therefore, we generated the actual LSTM-FCN and ALSTM-FCN implementations to record the results based on their actual code implementation. In addition, to record their training time, memory requirement, the number of parameters and f1-score. The Keras API [26] with TensorFlow backend [29] were used in the implementation of the LSTM-FCN, ALSTM-FCN and GRU-FCN models. The source code of our GRU-FCN implementation can be found on Github: https://github.com/NellyElsayed/GRU-FCN-model-for-univariate-time-series-classification.

We tested our model on the UCR time series archive [14] as one of the standard benchmarks for time series classification.

TABLE II.     THE UCR DATASETS DESCRIPTIONS BASED ON [14] AND THEIR EXPERIMENTAL ADJUSTMENTS USED IN THE GRU-FCN IMPLEMENTATION.

| Dataset | Type | # Classes | Length | Train size | Test size | # epochs | Train Batch | Test Batch |
|---|---|---|---|---|---|---|---|---|
| Adiac | Image | 37 | 176 | 390 | 391 | 4000 | 128 | 128 |
| ArrowHead | Image | 3 | 251 | 36 | 175 | 4000 | 32 | 128 |
| Beef | Spectro | 5 | 470 | 30 | 30 | 8000 | 64 | 64 |
| BeetleFly | Image | 2 | 512 | 20 | 20 | 8000 | 64 | 64 |
| BirdChicken | Image | 2 | 512 | 20 | 20 | 8000 | 64 | 64 |
| Car | Sensor | 4 | 577 | 60 | 60 | 2000 | 128 | 128 |
| CBF | Simulated | 3 | 128 | 30 | 900 | 2000 | 32 | 128 |
| ChlorineConc | Sensor | 3 | 166 | 467 | 3840 | 2000 | 128 | 128 |
| CinCECGTorso | Sensor | 4 | 1639 | 40 | 1380 | 500 | 128 | 128 |
| Coffee | Spectro | 2 | 286 | 28 | 28 | 500 | 64 | 64 |
| Computers | Device | 2 | 720 | 250 | 250 | 2000 | 128 | 128 |
| CricketX | Motion | 12 | 300 | 390 | 390 | 2000 | 128 | 128 |
| CricketY | Motion | 12 | 300 | 390 | 390 | 2000 | 128 | 128 |
| CricketZ | Motion | 12 | 300 | 390 | 390 | 2000 | 64 | 128 |
| DiatomSizeR | Image | 4 | 345 | 16 | 306 | 2000 | 64 | 64 |
| DisPhOAgeGrp | Image | 3 | 80 | 400 | 139 | 2000 | 128 | 128 |
| DisPhOCorrect | Image | 2 | 80 | 600 | 276 | 2000 | 128 | 128 |
| DisPhTW | Image | 6 | 80 | 400 | 139 | 2000 | 128 | 128 |
| Earthquakes | Sensor | 2 | 512 | 322 | 139 | 2000 | 128 | 128 |
| ECG200 | ECG | 2 | 96 | 100 | 100 | 8000 | 64 | 64 |
| ECG5000 | ECG | 5 | 140 | 500 | 4500 | 2000 | 128 | 128 |
| ECGFiveDays | ECG | 2 | 136 | 23 | 861 | 2000 | 128 | 128 |
| ElectricDevices | Device | 7 | 96 | 8926 | 7711 | 2000 | 128 | 128 |
| FaceAll | Image | 14 | 131 | 560 | 1690 | 2000 | 128 | 128 |
| FaceFour | Image | 4 | 350 | 24 | 88 | 2000 | 128 | 128 |
| FacesUCR | Image | 14 | 131 | 200 | 2050 | 2000 | 128 | 128 |
| FiftyWords | Image | 50 | 270 | 450 | 455 | 2000 | 128 | 128 |
| Fish | Image | 7 | 463 | 175 | 175 | 2000 | 128 | 128 |
| FordA | Sensor | 2 | 500 | 3601 | 1320 | 2000 | 128 | 128 |
| FordB | Sensor | 2 | 500 | 3636 | 810 | 1600 | 128 | 128 |
| GunPoint | Motion | 2 | 150 | 50 | 150 | 2000 | 128 | 128 |
| Ham | Spectro | 2 | 431 | 109 | 105 | 2000 | 128 | 128 |
| HandOutlines | Image | 2 | 2709 | 1000 | 370 | 2000 | 64 | 128 |
| Haptics | Motion | 5 | 1092 | 155 | 308 | 2000 | 128 | 128 |
| Herring | Image | 2 | 512 | 64 | 64 | 2000 | 128 | 128 |
| InlineSkate | Motion | 7 | 1882 | 100 | 550 | 2000 | 128 | 128 |
| InsWingSound | Sensor | 11 | 256 | 220 | 1980 | 1000 | 128 | 128 |
| ItalyPowD | Sensor | 2 | 24 | 67 | 1029 | 2000 | 64 | 128 |
| LargeKApp | Device | 3 | 720 | 375 | 375 | 2000 | 128 | 128 |
| Lightning2 | Sensor | 2 | 637 | 60 | 61 | 4000 | 128 | 128 |
| Lightning7 | Sensor | 7 | 319 | 70 | 73 | 3000 | 32 | 32 |
| Mallat | Simulated | 8 | 1024 | 55 | 2345 | 2500 | 128 | 128 |
| Meat | Spectro | 3 | 448 | 60 | 60 | 2000 | 64 | 128 |
| MedicalImages | Image | 10 | 99 | 381 | 760 | 2000 | 64 | 128 |
| MidPhOAgeGrp | Image | 3 | 80 | 400 | 154 | 2000 | 128 | 128 |
| MidPhOCorrect | Image | 2 | 80 | 600 | 291 | 2000 | 128 | 128 |
| MidPhTW | Image | 6 | 80 | 399 | 154 | 2000 | 128 | 128 |
| MoteStrain | Sensor | 2 | 84 | 20 | 1252 | 2000 | 128 | 128 |
| NonInvECGTh1 | ECG | 42 | 750 | 1800 | 1965 | 2000 | 128 | 128 |
| NonInvECGTh2 | ECG | 42 | 750 | 1800 | 1965 | 2000 | 128 | 128 |
| OliveOil | Spectro | 4 | 570 | 30 | 30 | 6000 | 64 | 128 |
| OSULeaf | Image | 6 | 427 | 200 | 242 | 2000 | 64 | 128 |
| PhalOCorrect | Image | 2 | 80 | 1800 | 858 | 2000 | 64 | 128 |
| Phoneme | Sensor | 39 | 1024 | 214 | 1896 | 2000 | 64 | 128 |
| Plane | Sensor | 7 | 144 | 105 | 105 | 200 | 16 | 16 |
| ProxPhOAgeGrp | Image | 3 | 80 | 400 | 205 | 2000 | 128 | 128 |
| ProxPhOCorrect | Image | 2 | 80 | 600 | 291 | 2000 | 128 | 128 |
| ProxPhTW | Image | 6 | 80 | 400 | 205 | 2000 | 128 | 128 |
| RefDevices | Device | 3 | 720 | 375 | 375 | 2000 | 64 | 64 |
| ScreenType | Device | 3 | 720 | 375 | 375 | 2000 | 64 | 128 |
| ShapeletSim | Simulated | 2 | 500 | 20 | 180 | 2000 | 128 | 128 |
| ShapesAll | Image | 60 | 512 | 600 | 600 | 4000 | 64 | 64 |
| SmlKitApp | Device | 3 | 720 | 375 | 375 | 2000 | 128 | 64 |
| SonyAIBORI | Sensor | 2 | 70 | 20 | 601 | 2000 | 64 | 128 |
| SonyAIBORII | Sensor | 2 | 65 | 27 | 953 | 2000 | 64 | 128 |
| StarLightCurves | Sensor | 3 | 1024 | 1000 | 8236 | 2000 | 64 | 64 |
| Strawberry | Spectro | 2 | 235 | 613 | 370 | 8000 | 64 | 64 |
| SwedishLeaf | Image | 15 | 128 | 500 | 625 | 8000 | 64 | 64 |
| Symbols | Image | 6 | 398 | 25 | 995 | 2000 | 64 | 64 |
| SynControl | Simulated | 6 | 60 | 300 | 300 | 4000 | 16 | 128 |
| ToeSegI | Motion | 2 | 277 | 40 | 228 | 2000 | 128 | 64 |
| ToeSegII | Motion | 2 | 343 | 36 | 130 | 2000 | 128 | 32 |
| Trace | Sensor | 4 | 275 | 100 | 100 | 1000 | 64 | 128 |
| TwoLeadECG | ECG | 2 | 82 | 23 | 1139 | 2000 | 64 | 64 |
| TwoPatterns | Simulated | 4 | 128 | 1000 | 4000 | 2000 | 32 | 128 |
| UWaveAll | Motion | 8 | 945 | 896 | 3582 | 500 | 16 | 16 |
| UWaveX | Motion | 8 | 315 | 896 | 3582 | 2000 | 64 | 16 |
| UWaveY | Motion | 8 | 315 | 896 | 3582 | 2000 | 64 | 64 |
| UWaveZ | Motion | 8 | 315 | 896 | 3582 | 2000 | 64 | 64 |
| Wafer | Sensor | 2 | 152 | 1000 | 6164 | 1500 | 64 | 64 |
| Wine | Spectro | 2 | 234 | 57 | 54 | 8000 | 64 | 64 |
| WordSynonyms | Image | 25 | 270 | 267 | 638 | 1500 | 64 | 64 |
| Worms | Motion | 5 | 900 | 181 | 77 | 2000 | 64 | 64 |
| WormsTwoClass | Motion | 2 | 900 | 181 | 77 | 1000 | 16 | 16 |
| Yoga | Image | 2 | 426 | 300 | 3000 | 1000 | 128 | 128 |

TABLE III.  CLASSIFICATION TESTING ERROR AND RANK FOR 85 TIME SERIES DATASETS FROM THE UCR BENCHMARK.

| Dataset | GRU-FCN | FCN | LSTMFCN | ALSTMFCN | ResNet | MCNN | MLP | COTE | DTW | PROP | BOSS | TSBF | ED |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Adiac | **0.127** | 0.143 | 0.141 | 0.139 | 0.174 | 0.231 | 0.248 | 0.233 | 0.396 | 0.353 | 0.235 | 0.231 | 0.389 |
| ArrowHead | **0.085** | 0.120 | 0.102 | 0.119 | 0.183 | / | 0.292 | 0.138 | 0.297 | 0.103 | 1.66 | 0.246 | 0.200 |
| Beef | **0.100** | 0.250 | 0.167 | 0.233 | 0.233 | 0.367 | 0.167 | 0.133 | 0.367 | 0.367 | 0.200 | 0.434 | 0.333 |
| BeetleFly | **0.050** | **0.050** | **0.050** | **0.050** | 0.200 | / | 0.200 | 0.050 | 0.300 | 0.400 | 0.100 | 0.200 | 0.250 |
| BirdChicken | **0** | 0.050 | 0 | **0** | 0.100 | / | 0.400 | 0.150 | 0.250 | 0.350 | 0.050 | 0.100 | 0.450 |
| Car | **0.016** | 0.050 | 0.033 | 0.159 | 0.067 | / | 0.117 | / | 0.267 | / | 0.167 | 0.217 | 0.267 |
| CBF | **0** | 0.008 | 0.003 | 0.004 | 0.006 | 0.002 | 0.14 | 0.001 | 0.003 | 0.002 | 0.002 | 0.013 | 0.148 |
| ChloConc | **0.002** | 0.157 | 0.191 | 0.193 | 0.172 | 0.203 | 0.125 | 0.314 | 0.352 | 0.360 | 0.339 | 0.308 | 0.350 |
| CinCECGTorso | 0.124 | 0.187 | 0.191 | 0.193 | 0.172 | 0.058 | 0.158 | **0.064** | 0.349 | 0.062 | 0.125 | 0.288 | 0.103 |
| Coffee | **0** | **0** | **0** | **0** | **0** | 0.036 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Computers | 0.148 | 0.152 | 0.136 | 0.123 | 0.176 | / | 0.504 | 0.240 | 0.300 | **0.116** | 0.244 | 0.244 | 0.424 |
| CricketX | 0.156 | 0.185 | 0.193 | 0.203 | 0.179 | 0.182 | 0.431 | **0.154** | 0.246 | 0.203 | 0.259 | 0.295 | 0.423 |
| CricketY | 0.156 | 0.208 | 0.183 | 0.185 | 0.195 | **0.154** | 0.405 | 0.167 | 0.256 | 0.156 | 0.208 | 0.265 | 0.433 |
| Cricketz | 0.154 | 0.187 | 0.190 | 0.175 | 0.169 | 0.142 | 0.408 | **0.128** | 0.246 | 0.156 | 0.246 | 0.285 | 0.413 |
| DiatomSizeR | 0.036 | 0.069 | 0.046 | 0.063 | 0.069 | **0.023** | 0.036 | 0.082 | 0.033 | 0.059 | 0.046 | 0.102 | 0.065 |
| DisPhOAgeGr | 0.142 | 0.165 | 0.145 | **0.137** | 0.202 | / | 0.178 | 0.229 | 0.230 | 0.223 | 0.272 | 0.218 | 0.374 |
| DisPhOCorrect | 0.168 | 0.188 | 0.168 | **0.163** | 0.180 | / | 0.195 | 0.238 | 0.283 | 0.232 | 0.252 | 0.288 | 0.283 |
| DisPhalanxTW | **0.180** | 0.210 | 0.185 | 0.185 | 0.260 | / | 0.375 | 0.317 | 0.410 | 0.317 | 0.324 | 0.324 | 0.367 |
| Earthquakes | **0.171** | 0.199 | 0.177 | 0.173 | 0.214 | / | 10.208 | / | 0.281 | 0.281 | 0.186 | 0.252 | 0.288 |
| ECG200 | **0.080** | 0.100 | 0.100 | 0.090 | 0.130 | / | 0.210 | 0.150 | 0.230 | / | 0.130 | 0.160 | 0.120 |
| ECG5000 | **0.052** | 0.059 | 0.053 | **0.052** | 0.069 | / | 0.068 | 0.054 | 0.076 | 0.350 | 0.059 | 0.061 | 0.075 |
| ECG5Days | **0** | 0.010 | 0.011 | 0.009 | 0.045 | **0** | 0.030 | 0 | 0.232 | 0.178 | **0** | 0.124 | 0.203 |
| ElectricDevices | **0.037** | 0.277 | **0.037** | **0.037** | 0.272 | / | 0.360 | 0.230 | 0.399 | 0.277 | 0.201 | 0.298 | 0.449 |
| FaceAll | **0.040** | 0.071 | 0.060 | 0.045 | 0.166 | 0.235 | 0.115 | 0.105 | 0.192 | 0.115 | 0.210 | 0.256 | 0.286 |
| FaceFour | 0.136 | 0.068 | 0.057 | 0.057 | 0.068 | 0 | 0.167 | 0.091 | 0.171 | 0.091 | **0** | **0** | 0.216 |
| FourUCR | 0.050 | 0.052 | 0.071 | 0.057 | **0.042** | 0.063 | 0.185 | 0.057 | 0.095 | 0.063 | **0.042** | 0.134 | 0.231 |
| FiftyWords | **0.167** | 0.321 | 0.196 | 0.176 | 0.273 | 0.190 | 0.288 | 0.191 | 0.301 | 0.180 | 0.301 | 0.242 | 0.369 |
| Fish | **0.006** | 0.029 | 0.017 | 0.023 | 0.011 | 0.051 | 0.126 | 0.029 | 0.177 | 0.034 | 0.011 | 0.166 | 0.217 |
| FordA | 0.074 | 0.094 | **0.072** | 0.073 | **0.072** | / | 0.231 | / | 0.444 | 0.182 | 0.083 | 0.150 | 0.335 |
| FordB | 0.083 | 0.117 | 0.088 | **0.081** | 0.100 | / | 0.371 | / | 0.380 | 0.265 | 0.109 | 0.402 | 0.394 |
| GunPoint | **0** | 0 | 0 | **0** | 0.007 | **0** | 0.067 | 0.007 | 0.093 | 0.007 | **0** | 0.014 | 0.087 |
| Ham | 0.209 | 0.238 | 0.209 | 0.228 | 0.219 | / | **0.162** | 0.334 | 0.533 | / | 0.334 | 0.239 | 0.400 |
| HandOutlines | 0.112 | 0.224 | 0.113 | 0.358 | 0.139 | / | 0.117 | 0.068 | 0.119 | / | **0.098** | 0.146 | 0.138 |
| Haptics | 0.455 | 0.449 | **0.425** | 0.435 | 0.495 | 0.530 | 0.539 | 0.488 | 0.623 | 0.584 | 0.536 | 0.510 | 0.630 |
| Herring | 0.250 | 0.297 | 0.250 | 0.265 | 0.406 | / | 0.360 | 0.313 | 0.469 | **0.079** | 0.454 | 0.360 | 0.484 |
| InlineSkate | 0.625 | 0.589 | 0.534 | **0.507** | 0.635 | 0.618 | 0.649 | 0.551 | 0.616 | 0.567 | 0.511 | 0.615 | 0.658 |
| InsWSound | 0.446 | 0.598 | 0.342 | **0.329** | 0.469 | / | 0.369 | / | 0.643 | / | 0.479 | 0.376 | 0.438 |
| ItalyPower | **0.027** | 0.030 | 0.037 | 0.040 | 0.040 | 0.030 | 0.034 | 0.036 | 0.050 | 0.039 | 0.053 | 0.117 | 0.045 |
| LKitApp | 0.090 | 0.104 | 0.090 | **0.083** | 0.107 | / | 0.520 | 0.136 | 0.205 | 0.232 | 0.235 | 0.472 | 0.507 |
| Lightening2 | 0.197 | 0.197 | 0.197 | 0.213 | 0.246 | 0.164 | 0.279 | 0.164 | 0.131 | **0.115** | 0.148 | 0.263 | 0.246 |
| Lightening7 | **0.137** | **0.137** | 0.164 | 0.178 | 0.164 | 0.219 | 0.356 | 0.247 | 0.274 | 0.233 | 0.342 | 0.274 | 0.427 |
| MALLAT | 0.048 | 0.020 | 0.019 | **0.016** | 0.021 | 0.057 | 0.064 | 0.036 | 0.066 | 0.050 | 0.058 | 0.040 | 0.086 |
| Meat | 0.066 | 0.033 | 0.116 | 0.033 | 0 | / | 0 | 0.067 | 0.067 | / | 0.100 | 0.067 | 0.067 |
| MedicalImages | **0.199** | 0.208 | **0.199** | 0.204 | 0.228 | 0.260 | 0.271 | 0.258 | 0.263 | 0.245 | 0.288 | 0.295 | 0.316 |
| MidPhOAgeGrp | 0.187 | 0.232 | 0.188 | 0.189 | 0.240 | / | 0.193 | **0.169** | 0.500 | 0.474 | 0.220 | 0.186 | 0.481 |
| MidPhOCorrect | 0.160 | 0.205 | 0.160 | 0.163 | 0.207 | / | 0.442 | 0.403 | 0.302 | 0.210 | 0.455 | 0.423 | 0.234 |
| MidPhTW | **0.363** | 0.388 | 0.383 | 0.373 | 0.393 | / | 0.429 | 0.429 | 0.494 | 0.630 | 0.455 | 0.403 | 0.487 |
| MoteStrain | 0.076 | **0.050** | 0.061 | 0.064 | 0.105 | 0.079 | 0.131 | 0.085 | 0.165 | 0.114 | 0.073 | 0.097 | 0.121 |
| NonInvECGTh1 | **0.034** | 0.039 | 0.035 | 0.025 | 0.052 | 0.064 | 0.058 | 0.093 | 0.210 | 0.178 | 0.161 | 0.158 | 0.171 |
| NonInvECGTh2 | **0.035** | 0.045 | 0.038 | 0.034 | 0.049 | 0.060 | 0.057 | 0.073 | 0.135 | 0.112 | 0.101 | 0.139 | 0.120 |
| OliveOil | **0.012** | 0.167 | 0.133 | 0.067 | 0.133 | 0.133 | 0.600 | 0.100 | 0.167 | 0.133 | 0.100 | 0.167 | 0.133 |
| OSULeaf | **0** | 0.012 | 0.004 | 0.004 | 0.021 | 0.271 | 0.430 | 0.145 | 0.409 | 0.194 | 0.012 | 0.240 | 0.479 |
| PhalOCorrect | **0.165** | 0.174 | 0.177 | 0.170 | 0.175 | / | 0.164 | 0.194 | 0.272 | / | 0.229 | 0.171 | 0.239 |
| Phoneme | 0.644 | 0.655 | 0.650 | **0.640** | 0.676 | / | 0.902 | / | 0.772 | / | 0.733 | 0.724 | 0.891 |
| Plane | **0** | 0 | 0 | **0** | **0** | / | 0.019 | / | **0** | / | **0** | 0 | 0.038 |
| ProxPhOeAgeGrp | 0.117 | 0.151 | 0.117 | **0.107** | 0.151 | / | 0.135 | 0.121 | 0.195 | 0.117 | 0.152 | 0.128 | 0.215 |
| ProxPhOCorrect | 0.079 | 0.100 | **0.065** | 0.075 | 0.082 | / | 0.200 | 0.142 | 0.217 | 0.172 | 0.166 | 0.152 | 0.192 |
| ProxPhTW | **0.167** | 0.190 | **0.167** | 0.173 | 0.193 | / | 0.210 | 0.186 | 0.244 | 0.244 | 0.200 | 0.191 | 0.293 |
| RefDevices | **0.407** | 0.467 | 0.421 | 0.429 | 0.472 | / | 0.632 | 0.443 | 0.536 | 0.424 | 0.498 | 0.528 | 0.605 |
| ScreenType | 0.297 | 0.333 | 0.351 | 0.341 | **0.293** | / | 0.614 | 0.411 | 0.603 | 0.440 | 0.536 | 0.491 | 0.640 |
| ShapeletSim | 0.011 | 0.133 | 0.011 | 0.011 | **0** | / | 0.528 | 0 | 0.350 | / | **0** | 0.039 | 0.461 |
| ShapesAll | 0.097 | 0.102 | 0.098 | 0.100 | **0.088** | / | 0.350 | 0.095 | 0.232 | 0.187 | 0.092 | 0.815 | 0.248 |
| SmlKitApp | 0.186 | 0.197 | 0.184 | 0.203 | 0.203 | / | 0.667 | **0.147** | 0.357 | 0.187 | 0.275 | 0.328 | 0.659 |
| SonyAIBORI | 0.017 | 0.032 | 0.018 | 0.030 | **0.015** | 0.230 | 0.273 | 0.146 | 0.275 | 0.293 | 0.321 | 0.205 | 0.305 |
| SonyAIBORII | 0.018 | 0.038 | 0.022 | 0.025 | 0.038 | 0.070 | 0.161 | 0.076 | 0.169 | 0.124 | 0.098 | 0.223 | 0.141 |
| StarLightCurves | 0.025 | 0.033 | 0.024 | 0.023 | 0.029 | 0.023 | 0.043 | 0.031 | 0.093 | 0.079 | **0.021** | 0.023 | 0.151 |
| Strawberry | **0.013** | 0.031 | **0.013** | **0.013** | 0.042 | / | 0.038 | 0.030 | 0.059 | / | 0.025 | 0.046 | 0.054 |
| SwedishLeaf | 0.016 | 0.034 | 0.021 | **0.014** | 0.042 | 0.066 | 0.107 | 0.046 | 0.208 | 0.085 | 0.272 | 0.085 | 0.211 |
| Symbols | 0.024 | 0.038 | 0.016 | **0.013** | 0.128 | 0.049 | 0.147 | 0.046 | 0.050 | 0.049 | 0.032 | 0.055 | 0.101 |
| SynControl | **0** | 0.010 | 0.003 | 0.006 | **0** | 0.003 | 0.050 | **0** | 0.007 | 0.010 | 0.030 | 0.007 | 0.120 |
| ToeSeg1 | 0.021 | 0.031 | **0.013** | **0.013** | 0.035 | / | 0.500 | 0.018 | 0.228 | 0.079 | 0.062 | 0.220 | 0.320 |
| ToeSeg2 | 0.076 | 0.085 | 0.084 | 0.077 | 0.138 | / | 0.408 | 0.047 | 0.162 | 0.085 | **0.039** | 0.200 | 0.192 |
| Trace | **0** | 0 | 0 | **0** | 0 | **0** | 0.180 | 0.010 | **0** | 0.010 | **0** | 0.020 | 0.240 |
| TwoLeadECG | **0** | 0 | 0.001 | 0.001 | 0 | 0.001 | 0.147 | 0.015 | 0.096 | **0** | 0.004 | 0.135 | 0.253 |
| TwoPatterns | 0.009 | 0.103 | 0.003 | 0.003 | **0** | 0.002 | 0.114 | 0 | **0** | 0.067 | 0.016 | 0.024 | 0.093 |
| UWaveAll | 0.078 | 0.174 | 0.096 | 0.107 | 0.132 | / | 0.253 | 0.161 | 0.108 | 0.199 | 0.238 | 0.170 | **0.052** |
| UWaveX | 0.171 | 0.246 | **0.151** | 0.152 | 0.213 | 0.180 | 0.232 | 0.196 | 0.273 | 0.199 | 0.241 | 0.264 | 0.261 |
| UWaveY | 0.240 | 0.275 | **0.233** | 0.234 | 0.332 | 0.268 | 0.297 | 0.267 | 0.366 | 0.283 | 0.313 | 0.228 | 0.338 |
| UWaveZ | 0.237 | 0.271 | 0.203 | **0.202** | 0.245 | 0.232 | 0.295 | 0.265 | 0.342 | 0.290 | 0.312 | 0.074 | 0.350 |
| Wafer | **0.001** | 0.003 | **0.001** | 0.002 | 0.003 | 0.002 | 0.004 | **0.001** | 0.020 | 0.003 | **0.001** | 0.005 | 0.005 |
| Wine | **0.111** | **0.111** | **0.111** | **0.111** | 0.204 | / | 0.056 | 0.223 | 0.426 | / | 0.260 | 0.389 | 0.389 |
| WordSynonyms | 0.262 | 0.420 | 0.329 | 0.332 | 0.368 | 0.276 | 0.406 | 0.266 | 0.351 | **0.226** | 0.345 | 0.312 | 0.382 |
| Worms | 0.325 | 0.331 | 0.325 | **0.320** | 0.381 | / | 0.585 | 0.442 | 0.416 | / | 0.442 | 0.312 | 0.545 |
| WormsTwoClass | 0.209 | 0.271 | 0.226 | 0.198 | 0.265 | / | 0.403 | 0.221 | 0.377 | / | **0.169** | 0.247 | 0.390 |
| Yoga | 0.090 | 0.098 | 0.082 | **0.081** | 0.142 | 0.112 | 0.145 | 0.113 | 0.164 | 0.121 | **0.081** | 0.181 | 0.170 |
| **no. best** | **39** | 9 | 19 | 25 | 13 | 5 | 3 | 11 | 4 | 5 | 13 | 3 | 2 |
| **Arith AVG Rank** | **2.947** | 5.841 | 3.818 | 3.729 | 6.035 | 9.118 | 9.100 | 6.071 | 9.882 | 8.253 | 7.071 | 8.459 | 10.676 |
| **MPCE** | **0.0308** | 0.0387 | 0.0327 | 0.0342 | 0.0415 | 0.1853 | 0.0725 | 0.0629 | 0.0734 | 0.1018 | 0.0558 | 0.0599 | 0.0807 |

Each dataset is divided into training and testing sets. The number of classes in each time series, the length and the size of both the training and test sets are shown in Table II based on the datasets description in [14]. The UCR benchmark datasets have different types of collected sources: 29 datasets of image source, 6 spectro source, 5 simulated source, 19 sensor source, 6 device source, 12 motion source, and 6 electrocardiogram (ECG) source. In addition, as we mentioned in the previous Section, Table II also shows the number of epochs through training, and the batch sizes of the training and testing stages based on our experiments.

We compared our GRU-FCN with several state-of-the-art time series methods that also were studied in [4] and [5]. These included FCN [4] which is based on a fully convolutional network, LSTM-FCN [5], ALSTM-FCN [5], that are based on long short-term memory and fully convolutional networks, ResNet [4] which based on convolutional residual networks, multi-scale convolution neural networks model (MCNN) [13], multi-layered perceptrons model (MLP) [4], collective of transformation-based ensembles model (COTE) [12] which based on transformation ensembles, dynamic time warping model (DTW) [30] that is based on a weighted dynamic time warping mechanism, PROP model [11] which is based on elastic distance measures, BOSS model [9] that based on noise reduction in the time series representation, time series based on a bag-of-features representation (TSBF) model [10], and Euclidean distance (ED) model [14]. Our model shows the overall highest number of being the best classifier for 39 time series out of 85. Our model also shows the overall smallest classification error, arithmetic average rank, and mean per-class classification error (MPCE) compared to the other models as shown in Table III.

Table IV shows a comparison between the number of parameters, training time and memory required to save the trainable weights of the GRU-FCN and both LSTM-FCN and ALSTM-FCN models as the existing LSTM-based to-date univariate classification models over the UCR 85 datasets. The GRU-FCN has a smaller number of parameters for all the datasets. The GRU-FCN saves overall 1207KB, and 5719KB memory requirements to save the trained model's weight; and 106.065, and 62.271 minutes to train the models over the UCR datasets comparing to the LSTM-FCN and ALSTM-FCN, respectively. Therefore, the GRU-FCN is preferable as a low budget classification model with high accuracy performance.

We evaluated our model using the Mean Per-Class Error (MPCE) used in [4] to evaluate the performance of a classification method over multiple datasets. The MPCE for a given model is calculated based on the per-class error (PCE) as follows:

$$\text{PCE}_m = \frac{e_m}{c_m} \tag{8}$$

$$\text{MPCE} = \frac{1}{M} \sum_{m=1}^{M} \text{PCE}_m \tag{9}$$

where $e_m$ is the error rate for dataset $m$ consisting of $c_m$ classes. $M$ is the number of tested datasets.

Table III shows the MPCE value for our GRU-FCN and other state-of-the-art models on the UCR benchmark datasets [14]. The results obtained by implementing GRU-FCN

and generating LSTM-FCN, and ALSTM models based on their actual implementation on Github. For the other models, we obtained the results from their own publications. Our GRU-FCN has the smallest MPCE value compared to the other state-of-the-art classification models. This means that generally, our GRU-FCN model performance across the different datasets is higher than the other state-of-the-art models.

Fig. 4, 5, 6, 7 are showing the loss value of both the training and validation processed of datasets. Each of these figures represents the loss process over image, motion, simulated, and source-obtained datasets from the UCR benchmark datasets respectively. These figures show that the average difference between the training and validation loss for the GRU-FCN is smaller than the LSTM-FCN and ALSTM-FCN models.

Table V shows the f1-score (also known as F-score or F-measure) [31], [32] for GRU-FCN, LSTM-FCN, and ALSTM-FCN classifiers. The f1-score shows the overall measure of a model's accuracy over each dataset used. The f1-score measuring based on both the precision and recall values of the classification model [31], [32]. The f1-score is calculated as follows [31], [32]:

$$precision = \frac{TP}{TP + FP} \tag{10}$$

$$recall = \frac{TP}{TP + FN} \tag{11}$$

$$f1\text{-}score = 2 \times \frac{precision \times recall}{precision + recall} \tag{12}$$

where TP, FP, FN stands for true-positive, false-positive and false-negative respectively. The GRU-FCN shows the highest f1-score for 53 out of 85 datasets comparing to the LSTM-FCN and ALSTM-FCN that both of these models have the highest f1-score for only 29 out of 85 datasets.

Fig. 8 shows the critical difference diagram [33] for Nemenyi or Bonferroni-Dunn test [34] with $\alpha = 0.05$ on our GRU-FCN and the state-of-the-art models based on the ranks arithmetic mean on the UCR benchmark datasets. This graph shows the significant classification accuracy improvement of our GRU-FCN compared to the other state-of-the-art models.

The Wilcoxon signed-rank test is one of the substantial tests to provide the classification method efficiency [35], [36]. Table VI shows the Wilcoxon signed-rank test [35], [37] among the twelve state-of-the-art classification models. This provides the overall accuracy evidence of each of the twelve classification methods.

## V. Conclusion

The proposed GRU-FCN classification model shows that replacing the LSTM by a GRU enhances the classification accuracy without requiring extra algorithm enhancements such as fine-tuning or attention algorithms. This The GRU also has a smaller architecture that requires fewer computations than the LSTM. Moreover, the GRU-based model requires a smaller number of trainable parameters, memory, and training time compared to the LSTM-based models. Furthermore, the proposed GRU-FCN classification model achieves the performance of state-of-the-art models and has the highest

TABLE IV. A COMPARISON BETWEEN THE GRU-FCN AND LSTM-BASED CLASSIFICATION MODELS FOR THE NUMBER OF PARAMETERS, TRAINING TIME (MINUTES), AND MEMORY (KB) REQUIRED TO SAVE THE MODEL WEIGHTS ON THE UCR 85 DATASETS [14].

| Dataset | Number of Parameters | | | Training Time (Minutes) | | | Memory (KB) | | |
|---|---|---|---|---|---|---|---|---|---|
| | GRU-FCN | LSTM-FCN | ALSTM-FCN | GRU-FCN | LSTM-FCN | ALSTM | GRU-FCN | LSTM-FCN | LSTM-FCN |
| Adiac | 275,237 | 276,717 | 283,837 | 9.597 | 9.560 | 10.056 | 1,114 | 1,119 | 1,150 |
| ArrowHead | 272,379 | 274,459 | 284,579 | 4.134 | 4.303 | 4.692 | 1,103 | 1,111 | 1,151 |
| Beef | 277,909 | 281,741 | 300,621 | 3.896 | 4.804 | 4.889 | 1,124 | 1,139 | 1,215 |
| BeetleFly | 278,506 | 282,674 | 303,234 | 3.937 | 4.208 | 4.545 | 1,126 | 1,144 | 1,225 |
| BirdChicken | 278,506 | 282,674 | 303,234 | 3.437 | 3.760 | 4.131 | 1,126 | 1,144 | 1,225 |
| Car | 280,340 | 285,028 | 308,188 | 1.899 | 1.972 | 2.045 | 1,134 | 1,152 | 1,245 |
| CBF | 269,427 | 270,523 | 275,723 | 5.243 | 5.248 | 5.339 | 1,092 | 1,096 | 1,117 |
| ChloConc | 270,339 | 271,739 | 278,459 | 13.324 | 14.601 | 14.813 | 1,095 | 1,110 | 1,127 |
| CinCECGTorso | 305,828 | 319,012 | 384,652 | 6.087 | 6.594 | 7.003 | 1,233 | 1,285 | 1,544 |
| Coffee | 273,082 | 275,442 | 286,962 | 0.504 | 0.524 | 0.540 | 1,104 | 1,115 | 1,161 |
| Computers | 283,498 | 289,330 | 318,210 | 7.722 | 8.049 | 8.436 | 1,145 | 1,170 | 1,283 |
| CricketX | 274,788 | 277,260 | 289,340 | 6.850 | 7.124 | 7.292 | 1,112 | 1,122 | 1,171 |
| CricketY | 274,788 | 277,260 | 289,340 | 6.673 | 6.978 | 7.224 | 1,112 | 1,122 | 1,171 |
| Cricketz | 274,788 | 277,260 | 289,340 | 8.601 | 8.933 | 9.539 | 1,112 | 1,122 | 1,171 |
| DiatomSizeR | 274,772 | 277,604 | 291,484 | 2.886 | 3.016 | 3.066 | 1,112 | 1,123 | 1,180 |
| DisPhOAgeGrp | 268,275 | 268,987 | 272,267 | 2.346 | 2.439 | 5.056 | 1,087 | 1,090 | 1,103 |
| DisPhOCorrect | 268,138 | 268,850 | 272,130 | 3.554 | 3.791 | 3.980 | 1,085 | 1,090 | 1,103 |
| DisPhTW | 268,686 | 269,398 | 272,678 | 2.611 | 2.723 | 2.876 | 1,088 | 1,091 | 1,106 |
| Earthquakes | 278,506 | 282,674 | 303,234 | 4.998 | 5.507 | 5.547 | 1,126 | 1,144 | 1,225 |
| ECG200 | 268,522 | 269,362 | 273,282 | 5.305 | 5.599 | 6.125 | 1,087 | 1,092 | 1,108 |
| ECG5000 | 269,989 | 271,181 | 276,861 | 13.223 | 13.797 | 14.162 | 1,093 | 1,098 | 1,123 |
| ECG5Days | 269,482 | 270,642 | 276,162 | 2.433 | 2.481 | 2.494 | 1,090 | 1,097 | 1,119 |
| ElectricDevices | 269,207 | 270,047 | 273,967 | 67.350 | 75.44 | 65.879 | 1,090 | 1,093 | 1,111 |
| FaceAll | 271,006 | 272,126 | 277,446 | 7.465 | 7.753 | 7.812 | 1,097 | 1,101 | 1,125 |
| FaceFour | 274,892 | 277,764 | 291,844 | 1.072 | 1.101 | 1.197 | 1,112 | 1,123 | 1,181 |
| FourUCR | 271,006 | 272,126 | 277,446 | 7.609 | 7.722 | 8.241 | 1,097 | 1,101 | 1,125 |
| FiftyWords | 279,274 | 281,506 | 292,386 | 6.052 | 6.353 | 6.428 | 1,129 | 1,138 | 1,183 |
| Fish | 278,015 | 281,791 | 300,391 | 3.770 | 3.850 | 3.912 | 1,125 | 1,139 | 1,214 |
| FordA | 278,218 | 282,290 | 302,370 | 43.135 | 44.861 | 47.525 | 1,124 | 1,142 | 1,221 |
| FordB | 278,218 | 282,290 | 302,370 | 26.781 | 27.341 | 27.890 | 1,124 | 1,142 | 1,221 |
| GunPoint | 269,818 | 271,090 | 277,170 | 1.003 | 1.046 | 1.138 | 1,092 | 1,098 | 1,123 |
| Ham | 276,562 | 280,082 | 297,402 | 2.048 | 2.127 | 2.160 | 1,118 | 1,133 | 1,202 |
| HandOutlines | 331,234 | 352,978 | 461,418 | 61.902 | 62.375 | 63.393 | 1,332 | 1,418 | 1,842 |
| Haptics | 292,837 | 301,645 | 345,405 | 9.787 | 10.023 | 10.631 | 1,183 | 1,217 | 1,390 |
| Herring | 278,506 | 282,674 | 303,234 | 1.633 | 1.668 | 1.706 | 1,126 | 1,144 | 1,225 |
| InlineSkate | 312,071 | 327,199 | 402,559 | 16.439 | 16.772 | 17.853 | 1,258 | 1,317 | 1,614 |
| InsWingSound | 273,595 | 275,715 | 286,035 | 4.332 | 4.510 | 4.599 | 1,107 | 1,115 | 1,158 |
| ItalyPowD | 266,794 | 267,058 | 268,098 | 2.719 | 3.015 | 3.048 | 1,080 | 1,083 | 1,087 |
| LargeKApp | 283,635 | 289,467 | 318,347 | 10.786 | 12.008 | 11.640 | 1,147 | 1,170 | 1,283 |
| Lightening2 | 281,506 | 286,674 | 312,234 | 3.887 | 3.940 | 4.065 | 1,137 | 1,159 | 1,260 |
| Lightening7 | 274,559 | 277,183 | 290,023 | 4.091 | 4.811 | 4.477 | 1,111 | 1,121 | 1,174 |
| MALLAT | 291,616 | 299,880 | 340,920 | 34.911 | 37.448 | 38.080 | 1,178 | 1,210 | 1,373 |
| Meat | 277,107 | 280,763 | 298,763 | 1.698 | 1.737 | 1.832 | 1,122 | 1,136 | 1,207 |
| MedicalImages | 269,690 | 270,554 | 274,594 | 5.361 | 5.456 | 6.498 | 1,092 | 1,095 | 1,114 |
| MidPhOAgeGrp | 268,275 | 268,987 | 272,267 | 1.802 | 2.138 | 2.182 | 1,087 | 1,090 | 1,103 |
| MidPhOCorrect | 268,138 | 268,850 | 272,130 | 3.219 | 3.374 | 3.528 | 1,085 | 1,090 | 1,103 |
| MidPhTW | 268,686 | 269,398 | 272,678 | 2.271 | 2.340 | 2.321 | 1,088 | 1,091 | 1,106 |
| MoteStrain | 268,234 | 268,978 | 272,418 | 2.398 | 2.423 | 2.481 | 1,085 | 1,090 | 1,104 |
| NonInvECGTh1 | 289,698 | 295,770 | 325,850 | 61.809 | 61.853 | 71.308 | 1,170 | 1,194 | 1,314 |
| NonInvECGTh2 | 289,698 | 295,770 | 325,850 | 59.212 | 60.554 | 60.754 | 1,170 | 1,194 | 1,314 |
| OliveOil | 280,172 | 284,804 | 307,684 | 3.267 | 3.670 | 4.073 | 1,133 | 1,151 | 1,243 |
| OSULeaf | 277,014 | 280,502 | 297,662 | 4.962 | 5.096 | 5.409 | 1,121 | 1,134 | 1,204 |
| PhalOCorrect | 268,138 | 268,850 | 272,130 | 16.319 | 19.269 | 21.159 | 1,085 | 1,090 | 1,103 |
| Phoneme | 295,863 | 304,127 | 345,167 | 29.778 | 31.34 | 37.147 | 1,194 | 1,226 | 1,389 |
| Plane | 270,359 | 271,583 | 277,423 | 0.497 | 0.502 | 0.575 | 1,095 | 1,099 | 1,125 |
| ProxPhOAgeGrp | 268,275 | 268,987 | 272,267 | 3.550 | 3.601 | 3.605 | 1,087 | 1,090 | 1,103 |
| ProxPhOCorrect | 268,138 | 268,850 | 272,130 | 4.142 | 4.538 | 4.678 | 1,085 | 1,090 | 1,103 |
| ProxPhTW | 268,686 | 269,398 | 272,678 | 2.050 | 2.201 | 2.126 | 1,088 | 1,091 | 1,106 |
| RefDevices | 283,635 | 289,467 | 318,347 | 12.878 | 14.160 | 14.460 | 1,147 | 1,170 | 1,283 |
| ScreenType | 283,635 | 289,467 | 318,347 | 13.327 | 13.890 | 14.283 | 1,147 | 1,170 | 1,283 |
| ShapeletSim | 278,218 | 282,290 | 302,370 | 1.596 | 1.628 | 2.004 | 1,124 | 1,142 | 1,221 |
| ShapesAll | 286,452 | 290,620 | 311,180 | 34.243 | 36.523 | 37.627 | 1,157 | 1,173 | 1,256 |
| SmlKitApp | 283,635 | 289,467 | 318,347 | 12.417 | 12.92 | 14.248 | 1,147 | 1,170 | 1,283 |
| SonyAIBORI | 267,898 | 268,530 | 271,410 | 0.982 | 1.931 | 2.042 | 1,084 | 1,088 | 1,100 |
| SonyAIBORII | 267,778 | 268,370 | 271,050 | 2.492 | 2.496 | 2.873 | 1,084 | 1,088 | 1,099 |
| StarLightCurves | 290,931 | 299,195 | 340,235 | 151.538 | 157.143 | 161.447 | 1,176 | 1,208 | 1,369 |
| Strawberry | 271,858 | 273,810 | 283,290 | 39.138 | 40.408 | 42.769 | 1,100 | 1,109 | 1,147 |
| SwedishLeaf | 271,071 | 272,167 | 277,367 | 6.931 | 7.572 | 7.891 | 1,098 | 1,102 | 1,125 |
| Symbols | 276,318 | 279,574 | 295,574 | 6.176 | 6.543 | 6.736 | 1,118 | 1,131 | 1,196 |
| SynControl | 268,206 | 268,758 | 271,238 | 20.562 | 21.735 | 23.209 | 1,086 | 1,088 | 1,101 |
| ToeSeg1 | 272,866 | 275,154 | 286,314 | 1.824 | 1.846 | 1.900 | 1,104 | 1,114 | 1,158 |
| ToeSeg2 | 274,450 | 277,266 | 291,066 | 1.415 | 1.549 | 1.629 | 1,110 | 1,122 | 1,177 |
| Trace | 273,092 | 275,364 | 286,444 | 0.977 | 1.021 | 1.093 | 1,105 | 1,114 | 1,160 |
| TwoLeadECG | 268,186 | 268,914 | 272,274 | 3.053 | 3.535 | 3.498 | 1,085 | 1,090 | 1,104 |
| TwoPatterns | 269,564 | 270,660 | 275,860 | 33.994 | 37.673 | 38.303 | 1,092 | 1,096 | 1,119 |
| UWaveAll | 289,720 | 297,352 | 335,232 | 24.983 | 28.702 | 28.874 | 1,170 | 1,200 | 1,351 |
| UWaveX | 274,600 | 277,192 | 289,872 | 30.214 | 32.095 | 33.573 | 1,111 | 1,121 | 1,173 |
| UWaveY | 274,600 | 277,192 | 289,872 | 30.214 | 31.526 | 32.526 | 1,111 | 1,121 | 1,173 |
| UWaveZ | 274,600 | 277,192 | 289,872 | 30.214 | 31.881 | 33.573 | 1,111 | 1,121 | 1,173 |
| Wafer | 269,866 | 271,154 | 277,314 | 20.438 | 21.835 | 22.018 | 1,092 | 1,099 | 1,123 |
| Wine | 271,834 | 273,778 | 283,218 | 3.771 | 4.021 | 4.530 | 1,099 | 1,109 | 1,146 |
| WordSynonyms | 275,849 | 278,081 | 288,961 | 4.911 | 5.155 | 5.498 | 1,116 | 1,125 | 1,170 |
| Worms | 288,229 | 295,501 | 331,581 | 4.484 | 4.669 | 5.019 | 1,165 | 1,193 | 1,336 |
| WormsTwoClass | 287,818 | 295,090 | 331,170 | 3.536 | 3.586 | 4.134 | 1,162 | 1,192 | 1,334 |
| Yoga | 276,442 | 279,922 | 297,042 | 10.970 | 11.606 | 10.753 | 1,118 | 1,133 | 1,200 |
| Total | **23,555,876** | 23,849,100 | 25,291,420 | **1145.645** | 1207.916 | 1251.71 | **95,273** | 96,480 | 100,992 |

TABLE V.     THE f1-SCORE VALUE OF THE PROPOSED GRU-FCN MODEL AND THE LSTM-BASED ARCHITECTURES OVER THE UCR BENCHMARK
DATASETS [14].

| Dataset | f1-Score | | |
|---|---|---|---|
| | **GRU-FCN** | **LSTM-FCN** | **ALSTM-FCN** |
| Adiac | **0.795** | 0.770 | 0.780 |
| ArrowHead | **0.711** | 0.694 | 0.695 |
| Beef | 0.819 | **0.873** | 0.765 |
| BeetleFly | **1.0** | **1.0** | 0.949 |
| BirdChicken | **1.0** | **1.0** | **1.0** |
| Car | **0.954** | 0.952 | 0.947 |
| CBF | **0.995** | 0.994 | 0.989 |
| ChlorineCon | 0.766 | **0.791** | 0.767 |
| CinCECGTorso | **0.379** | 0.321 | 0.375 |
| Coffee | **1.0** | **1.0** | **1.0** |
| Computers | **0.916** | 0.914 | 0.913 |
| CricketX | **0.786** | 0.782 | 0.784 |
| CricketY | 0.756 | **0.786** | 0.776 |
| CricketZ | **0.779** | 0.778 | 0.761 |
| DiatomSizeR | 0.926 | 0.926 | **0.935** |
| DisPhOAgeGrp | **0.645** | 0.614 | 0.636 |
| DisPhOCorrect | **0.813** | 0.804 | **0.813** |
| DisPhTW | 0.477 | 0.469 | **0.479** |
| Earthquakes | **0.483** | 0.466 | 0.466 |
| ECG200 | **0.910** | 0.900 | 0.909 |
| ECG5000 | 0.253 | 0.251 | **0.263** |
| ECGFiveDays | **0.991** | **0.991** | **0.991** |
| ElectricDevices | 0.195 | 0.196 | **0.197** |
| FaceAll | **0.137** | 0.134 | 0.136 |
| FaceFour | **0.960** | 0.949 | 0.949 |
| FacesUCR | 0.892 | **0.898** | 0.896 |
| 50words | **0.353** | 0.330 | **0.353** |
| Fish | 0.962 | **0.964** | 0.957 |
| FordA | 0.926 | **0.928** | **0.928** |
| FordB | 0.928 | **0.930** | 0.929 |
| GunPoint | **1.0** | **1.0** | **1.0** |
| Ham | **0.788** | **0.788** | 0.770 |
| HandOutlines | **0.875** | 0.873 | 0.866 |
| Haptics | **0.528** | 0.523 | 0.515 |
| Herring | 0.717 | **0.722** | 0.694 |
| InlineSkate | 0.454 | **0.474** | 0.446 |
| InWingSound | **0.477** | 0.432 | 0.410 |
| ItalyPower | 0.970 | 0.970 | **0.972** |
| LargeKApp | 0.406 | 0.407 | **0.410** |
| Lightning2 | 0.765 | **0.767** | **0.767** |
| Lightning7 | **0.872** | 0.833 | 0.858 |
| MALLAT | 0.971 | 0.970 | 0.971 |
| Meat | 0.925 | 0.870 | **0.973** |
| MedicalImages | **0.714** | 0.686 | 0.701 |
| MidPhOutlineAgeGrp | **0.507** | 0.347 | 0.445 |
| MidPhOCorrect | **0.823** | 0.821 | 0.819 |
| MidPhTW | **0.329** | 0.314 | 0.320 |
| MoteStrain | 0.925 | 0.920 | 0.915 |
| NonInvECGTh1 | **0.911** | 0.908 | 0.905 |
| NonInvECGTh2 | **0.899** | 0.896 | 0.894 |
| OliveOil | 0.853 | 0.611 | **0.885** |
| OSULeaf | **0.988** | 0.979 | **0.988** |
| PhalOCorrect | **0.812** | 0.803 | 0.809 |
| Phoneme | 0.025 | **0.026** | **0.026** |
| Plane | **0.888** | **0.888** | 0.882 |
| ProxPhOeAgeGrp | **0.600** | 0.594 | 0.436 |
| ProxPhOCorrect | 0.896 | **0.904** | 0.896 |
| ProxPhTW | 0.545 | 0.504 | 0.469 |
| RefDevices | **0.277** | 0.241 | 0.241 |
| ScreenType | 0.297 | 0.302 | **0.308** |
| ShapeletSim | **0.842** | **0.842** | **0.842** |
| ShapesAll | **0.108** | **0.108** | 0.107 |
| SmlKitApp | 0.345 | 0.361 | **0.370** |
| SonyAIBORI | **0.984** | 0.974 | 0.983 |
| SonyAIBORII | **0.980** | 0.978 | 0.977 |
| StarLightCurves | **0.975** | 0.961 | 0.962 |
| Strawberry | **0.818** | **0.818** | **0.818** |
| SwedishLeaf | 0.807 | 0.801 | **0.811** |
| Symbols | 0.980 | **0.982** | 0.974 |
| SynControl | **0.522** | 0.516 | 0.511 |
| ToeSeg1 | 0.708 | **0.746** | **0.746** |
| ToeSeg2 | **0.582** | 0.563 | 0.577 |
| Trace | **1.0** | 0.986 | 0.983 |
| TwoLeadECG | **0.999** | **0.999** | **0.999** |
| TwoPatterns | 0.986 | **0.989** | 0.971 |
| UWaveAll | **0.782** | 0.766 | 0.754 |
| UWaveX | **0.665** | 0.654 | 0.659 |
| UWaveY | **0.698** | 0.695 | 0.686 |
| UWaveZ | 0.736 | **0.739** | 0.743 |
| Wafer | **0.996** | **0.996** | **0.996** |
| Wine | **0.887** | **0.887** | **0.887** |
| WordSynonyms | **0.380** | 0.327 | 0.345 |
| Worms | **0.448** | 0.423 | 0.425 |
| WormsTwoClass | 0.530 | 0.525 | **0.542** |
| Yoga | 0.882 | 0.906 | **0.914** |

Fig. 4.   The loss value of GRU-FCN, LSTM-FCN, and ALSTM-FCN models over the image-source obtained (DiatomSizeR dataset) training and validation processes.



Fig. 5.   The loss value of GRU-FCN, LSTM-FCN, and ALSTM-FCN models over the motion-source obtained (CricketX dataset) training and validation processes.



Fig. 6.   The loss value of GRU-FCN, LSTM-FCN, and ALSTM-FCN models over the simulated-source obtained (CDF dataset) training and validation processes.


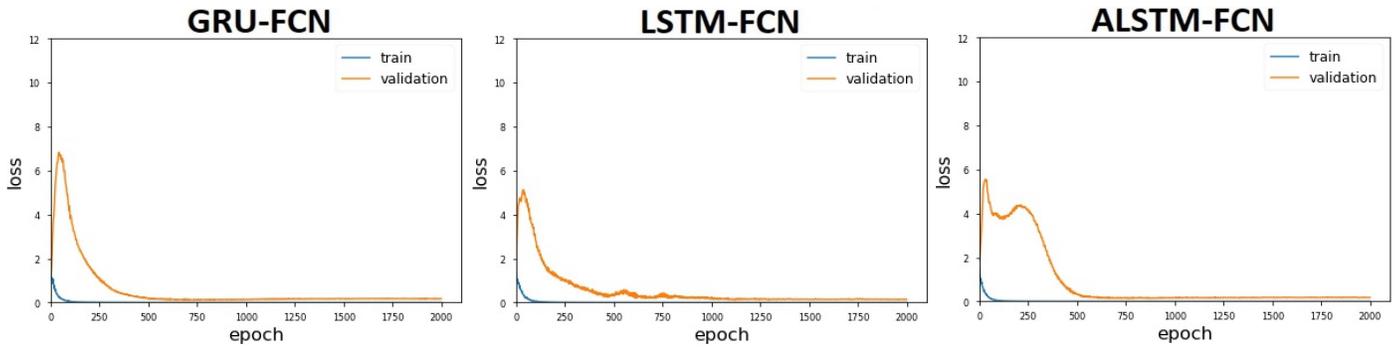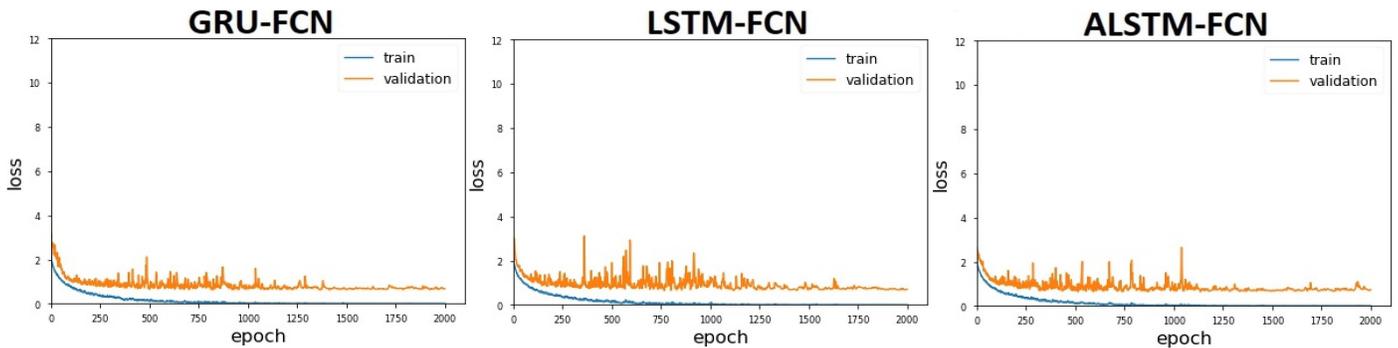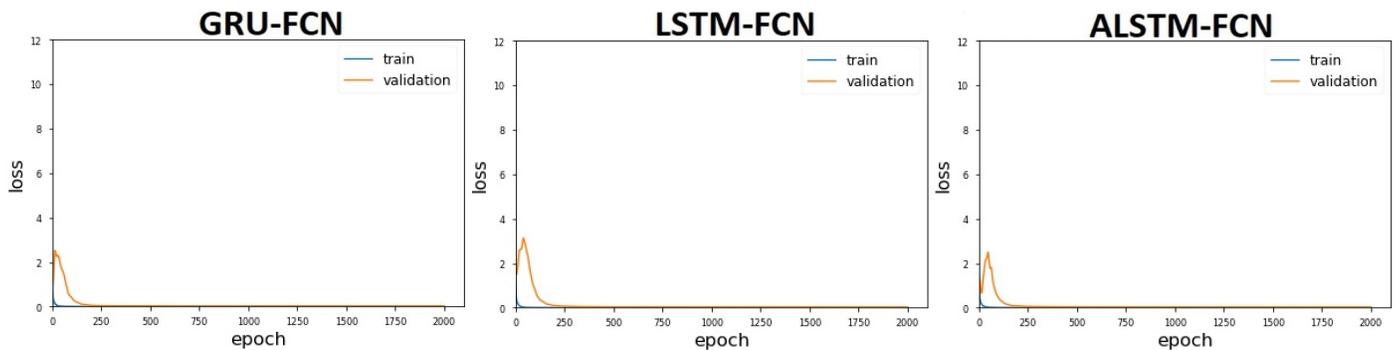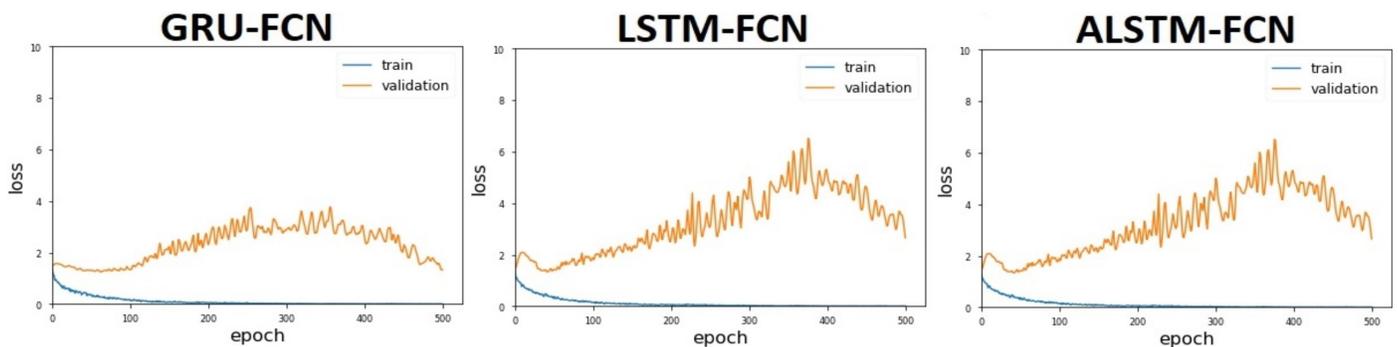
Fig. 7.   The loss value of GRU-FCN, LSTM-FCN, and ALSTM-FCN models over the sensor-source obtained (ChlorineCon dataset) training and validation processes.
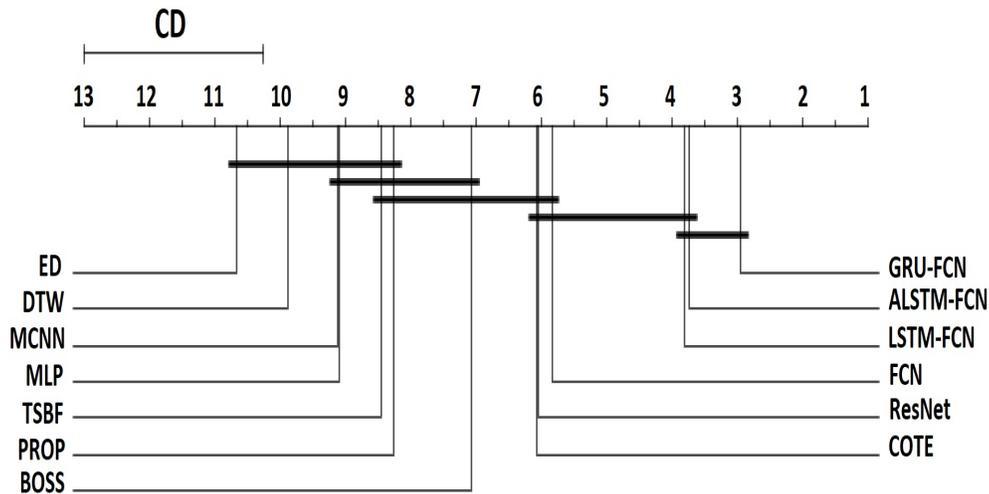
Fig. 8. Critical difference diagram based on the arithmetic mean of model ranks.

TABLE VI. WILCOXON SIGNED-RANK TEST ON GRU-FCN AND 10 BENCHMARK MODEL ON THE 85 DATASETS FROM UCR BENCHMARK [14].

| | FCN | LSTM-FCN | ALSTM-FCN | ResNet | MCNN | MLP | COTE | DTW | PROP | BOSS | TSBF | ED |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **GRU-FCN** | 3.44E-10 | 4.95E-03 | **4.00E-02** | 2.53E-11 | 1.05E-12 | 1.43E-13 | 1.25E-08 | 1.23E-14 | 1.58E-11 | 4.37E-10 | 2.77E-12 | 2.93E-15 |
| **FCN** | | 4.37E-09 | 8.58E-08 | **1.68E-01** | 9.31E-10 | 1.12E-09 | **1.85E-02** | 3.49E-12 | 1.31E-07 | 8.02E-04 | 1.10E-07 | 7.07E-13 |
| **LSTM-FCN** | | | **7.45E-01** | 2.24E-09 | 1.40E-11 | 6.09E-13 | 1.03E-06 | 2.35E-14 | 5.72E-11 | 2.85E-9 | 6.40E-13 | 1.08E-14 |
| **ALSTM-FCN** | | | | **1.40E-07** | 1.02E-11 | 8.35E-12 | **2.33E-07** | 1.55E-13 | 7.95E-11 | 4.71E-09 | 3.30E-12 | 4.73E-14 |
| **ResNet** | | | | | **6.28E-09** | 1.79E-08 | **2.46E-01** | 9.32E-13 | 1.76E-06 | **1.28E-03** | 1.56E-07 | 1.11E-13 |
| **MCNN** | | | | | | **4.35E-05** | 4.77E-08 | 5.76E-05 | **6.10E-04** | 1.20E-06 | 7.72E-06 | **2.18E-04** |
| **MLP** | | | | | | | 7.04E-05 | **7.28E-01** | **7.13E-01** | **1.08E-03** | 5.70E-03 | 3.25E-04 |
| **COTE** | | | | | | | | 1.62E-06 | 2.28E-05 | **7.74E-03** | 3.59E-04 | 3.22E-07 |
| **DTW** | | | | | | | | | **2.05E-01** | 2.37E-07 | 1.80E-04 | 2.13E-03 |
| **PROP** | | | | | | | | | | **8.82E-03** | **5.13E-01** | **3.14E-02** |
| **BOSS** | | | | | | | | | | | **3.18E-02** | 7.02E-10 |
| **TSBF** | | | | | | | | | | | | **6.65E-08** |

average arithmetic ranking and the lowest mean per-class error (MPCE) through time series datasets classification of the UCR benchmark compared to the state-of-the-art models. Therefore, replacing the LSTM by GRU in the LSTM-FCN for univariate time series classification can improve the classification with smaller model architecture.

## VI. ACKNOWLEDGMENT

## REFERENCES

[1] J. D. Hamilton, *Time series analysis.* Princeton University Press, Princeton, NJ, 1994, vol. 2.

[2] H. Sohn and C. R. Farrar, "Damage diagnosis using time series analysis of vibration signals," *Smart Materials and Structures*, vol. 10, no. 3, p. 446, 2001.

[3] M. Gul and F. N. Catbas, "Statistical pattern recognition for structural health monitoring using time series modeling: theory and experimental verifications," *Mechanical Systems and Signal Processing*, vol. 23, no. 7, pp. 2192–2204, 2009.

[4] Z. Wang, W. Yan, and T. Oates, "Time series classification from scratch with deep neural networks: a strong baseline," in *Neural Networks (IJCNN), 2017 International Joint Conference on.* IEEE, 2017, pp. 1578–1585.

[5] F. Karim, S. Majumdar, H. Darabi, and S. Chen, "LSTM fully convolutional networks for time series classification," *IEEE Access*, vol. 6, pp. 1662–1669, 2018.

[6] A. Amei, W. Fu, and C.-H. Ho, "Time series analysis for predicting the occurrences of large scale earthquakes," *International Journal of Applied Science and Technology*, vol. 2, no. 7, 2012.

[7] J. Rotton and J. Frey, "Air pollution, weather, and violent crimes: concomitant time-series analysis of archival data." *Journal of Personality and Social Psychology*, vol. 49, no. 5, p. 1207, 1985.

[8] E. Keogh and C. A. Ratanamahatana, "Exact indexing of dynamic time warping," *Knowledge and Information Systems*, vol. 7, no. 3, pp. 358–386, 2005.

[9] P. Schäfer, "The BOSS is concerned with time series classification in the presence of noise," *Data Mining and Knowledge Discovery*, vol. 29, no. 6, pp. 1505–1530, 2015.

[10] M. G. Baydogan, G. Runger, and E. Tuv, "A bag-of-features framework to classify time series," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 11, pp. 2796–2802, 2013.

[11] J. Lines and A. Bagnall, "Time series classification with ensembles of elastic distance measures," *Data Mining and Knowledge Discovery*, vol. 29, no. 3, pp. 565–592, 2015.

[12] A. Bagnall, J. Lines, J. Hills, and A. Bostrom, "Time-series classification with COTE: the collective of transformation-based ensembles," *IEEE Transactions on Knowledge and Data Engineering*, vol. 27, no. 9, pp. 2522–2535, 2015.

[13] Z. Cui, W. Chen, and Y. Chen, "Multi-scale convolutional neural networks for time series classification," *arXiv preprint arXiv:1603.06995*, 2016.

[14] Y. Chen, E. Keogh, B. Hu, N. Begum, A. Bagnall, A. Mueen, and G. Batista. (2015, July) The UCR time series classification archive. [Online]. Available: http://www.cs.ucr.edu/~eamonn/time_series_data/

[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] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel, "Backpropagation applied to handwritten zip code recognition," *Neural Computation*, vol. 1, no. 4, pp. 541–551, 1989.

[17] Y. LeCun and Y. Bengio, "Convolutional networks for images, speech, and time series," in *The Handbook of Brain Theory and Neural Networks*. MIT Press, 1995, pp. 255–258.

[18] I. Goodfellow, Y. Bengio, A. Courville, and Y. Bengio, *Deep Learning*. MIT press Cambridge, 2016, vol. 1.

[19] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," *arXiv preprint arXiv:1502.03167*, 2015.

[20] V. Nair and G. E. Hinton, "Rectified linear units improve restricted Boltzmann machines," in *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*, 2010, pp. 807–814.

[21] (2019) Keras recurrent layers documentation. [Online]. Available: https://keras.io/layers/recurrent/

[22] Y.-L. Boureau, J. Ponce, and Y. LeCun, "A theoretical analysis of feature pooling in visual recognition," in *Proceedings of the 27th international conference on machine learning (ICML-10)*, 2010, pp. 111–118.

[23] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: surpassing human-level performance on imagenet classification," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 1026–1034.

[24] C. Gulcehre, M. Moczulski, M. Denil, and Y. Bengio, "Noisy activation functions," in *International Conference on Machine Learning*, 2016, pp. 3059–3068.

[25] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks," in *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, 2010, pp. 249–256.

[26] F. Chollet *et al.*, "Keras," https://keras.io, 2015.

[27] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[28] N. M. Nasrabadi, "Pattern recognition and machine learning," *Journal of Electronic Imaging*, vol. 16, no. 4, p. 049901, 2007.

[29] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng, "TensorFlow: Large-scale machine learning on heterogeneous systems," 2015, software available from tensorflow.org. [Online]. Available: https://www.tensorflow.org/

[30] Y.-S. Jeong, M. K. Jeong, and O. A. Omitaomu, "Weighted dynamic time warping for time series classification," *Pattern Recognition*, vol. 44, no. 9, pp. 2231–2240, 2011.

[31] Y. Sasaki *et al.*, "The truth of the f-measure," *Teach Tutor mater*, vol. 1, no. 5, pp. 1–5, 2007.

[32] D. M. Powers, "Evaluation: from precision, recall and f-measure to roc, informedness, markedness and correlation," 2011.

[33] J. Demšar, "Statistical comparisons of classifiers over multiple data sets," *Journal of Machine Learning Research*, vol. 7, no. Jan, pp. 1–30, 2006.

[34] T. Pohlert, "The pairwise multiple comparison of mean ranks package (PMCMR)," *R Package*, vol. 27, 2014.

[35] R. Woolson, "Wilcoxon signed-rank test," *Wiley encyclopedia of clinical trials*, pp. 1–3, 2007.

[36] D. Rey and M. Neuhäuser, "Wilcoxon-signed-rank test," *International encyclopedia of statistical science*, pp. 1658–1659, 2011.

[37] R. Lowry, "Concepts and applications of inferential statistics," 2014.

# Impact Study and Evaluation of Higher Modulation Schemes on Physical Layer of Upcoming Wireless Mobile Networks

Heba Haboobi[1], Mohammad R Kadhum[2]
Faculty of Arts, Science & Technology
University of Northampton
Northampton, U.K.

*Abstract*—In this paper, the higher modulation formats (128 and 256) Quadrature Amplitude Modulation (QAM), for modulation/demodulation the digital signal of the currently used Orthogonal Frequency Division Multiplexing (OFDM) system, is proposed, explored and evaluated at a wireless transmission system. The proposed modulation schemes are utilized to study the impact of adding extra bits for each transmitted sample on system performance in terms of the channel capacity, Bit Error Rate (BER) and Signal to Noise Ratio (SNR). As such, the key purpose of this research is to identify the advantages and disadvantages of using higher modulation schemes on the physical layer (PHY) of future mobile networks. In addition, the trade-off relation between the achieved bit rate and the required power of the receiver is examined in the presence of the Additive White Gaussian Noise (AWGN) and Rayleigh noise channels. Besides, the currently employed waveform (OFDM) is considered herein as an essential environment to test the effect of receiving additional complex numbers on the constellations table. Thus, investigate the ability to recognize both the phase and amplitude of intended constellations for the upcoming design of wireless transceivers. Moreover, a MATLAB simulation is employed to evaluate the proposed system mathematically and physically in an electrical back-to-back transmission system.

*Keywords*—*Orthogonal Frequency Division Multiplexing (OFDM); Quadrature Amplitude Modulation (QAM); Bit Error Rate (BER); Signal to Noise Ratio (SNR); Bandwidth (BW); Additive White Gaussian Noise (AWGN); Rayleigh noise; physical layer (PHY)*

## I. Introduction

Enlarged growth of propagated data over wireless mobile networks is associated directly with an expanded range of emerging data-hungry applications. Hence, nowadays, using smart phones, users are looking for downloading and uploading huge amounts of data rather than the conventional operations. Thus, the predicted scenarios for the next generation of mobile are considered as an essential driving force for introducing higher modulation schemes [1].

Recently, many research studies have proposed new solutions to meet with the future market demand for mobile networks [2]. Consequently, achieve developed rates of channel capacity at acceptable limits of errors get overwhelming interest in terms of R&D. As such, high efforts have been made to promote both bit rate and bandwidth (BW) efficiency by employing higher modulation formats like 64 Quadrature Amplitude Modulation (QAM).

Nevertheless, so far, investigation in more advanced formats of modulation like the 128 and 256 QAM is yet not introduced for future wireless networks of mobile [3].

It's worth noting that, comparing to the lower modulation formats, the higher order schemes can significantly improve both BW efficiency and maximum bit rate with raising the minimum limits of the Signal to Noise Ratio (SNR) to keep the required signal at the acceptable level of the Bit Error Rate (BER). As such, for the wireless transmission system, the modulation and demodulation process can play a key role in regulating the BER levels and maximum rates of transmitted data [4].

The big challenge for future wireless networks is increasing transmission and reception rates for mobile communications and decreasing required levels of power and BER [5]. That's mean, the required modulation format is essentially selected in accordance with the future system performance, in terms of the bit rate, and the BER.

Utilizing an appropriate type of a modulation format for the advanced mobile transmission is critical due to key technical limitations. Thus, factors like an optimal usage of the BW, limitation of received power and probability of constellation table noise are major restrictions for developers in the modulation format field of wireless communication.

From digital signal processing (DSP) perspective, advanced digital modulation schemes can play a big role in expanding channel capacity by transmitting high amounts of data per each slice of time and with an acceptable immunity of noise particularly when the required level of power is afforded [6]. Thus, the performance in terms of maximum data bit rate and minimum error rate of a modulation scheme is decided depending on the efficient usage of the BW and power.

The key question of this study is, can the higher modulation formats be a good solution for a developed performance of future wireless mobile technology?

In this context, this study discusses how the developed modulation formats can impact the performance of transmission in a thorough manner by comparing not only the major principles of the conventional modulation formats but also including the most recent one. As such, investigating a higher scheme like 128 and 256 QAM to get a promising progress in this field of research. Hence, emphasis on the trade-off relation

between the improved BW efficiency and required level of received power for future mobile networks.

Depending on Shannon, the channel capacity can be decided by the maximum number of transmitted bits throw an obtainable BW and at an acceptable level of errors. Thus, in digital wireless communication, the quality of the transmission service is typically measured by both the BER and the maximum achieved bit rate [7].

The main principle beyond designing an optimal digital communication system is to obtain lower levels of error and get more efficient usage of channel BW simultaneously. Hence, an efficient modulation system can basically assistant a better spectral efficiency (SE), where the SE represents a number of bits transmitted per second for each Hertz of the BW [8].

To accommodate the developed schemes of modulation format in the physical layer (PHY) of wireless mobile networks, the current waveform Orthogonal Frequency Division Multiplexing (OFDM) is utilized due to its high resistance against the Inter-Symbol Interference (ISI), better BW efficiency, and improved bit rate of transmission [9]. Thus, the newly examined formats of modulation (128 QAM & 256 QAM) are explored employing the most popular transmission air interface OFDM.

Since the orthogonality is the key feature of the currently applied waveform, a better level of transmission can be delivered by the OFDM in comparison with the old-style Frequency Division Multiplexing (FDM) [10].

It's worth noting that, the impact of unwanted signals (noise) is considered herein for both the line of sight (LoS) and non-line of sight (NLoS) transmission. Thus, different models of channels like the Additive white Gaussian noise (AWGN) and Rayleigh noise are adopted to compare the performance of the transmission system for the most popular modulation schemes.

The rest of the paper is structured as follows: Section 2 demonstrates a literature review of the prepared study. Section 3 deliberates theoretically the main concepts of the proposed system highlighting physically and mathematically the fundamentals beyond it. Section 4 simulates numerically system performance (channel capacity and BER). Section 5 concludes the outlines of the paper.

## II. LITERATURE REVIEW

Since the main target of a future wireless system designer is to transmit bigger data rate within the offered BW and at an affordable expense [11], digital Modulation still has a high impact relation on the world of the developed throughput communication system. In this context, various studies about the relation between the most common modulation formats and channel models of the wireless communication system are investigated.

Regarding the Binary Phase Shift Keying (BPSK), it's the simplest form of modulation formats which is used for representing digital data employing two changed phases of each available subcarrier. The two utilized phases are detached by 180 degree, and a regular angular interval around a circle is ordinarily applied to place the selected points at the constellation table.

The best immunity to distortion is achieved due to the maximum separation in phase between adjacent constellation points. Thus, two phases are assigned arbitrarily, (0 and $\pi$), representing the binary "0" and "1".

Besides, their circular position makes them able to get an equal energy of transmission. In addition, the complex numbers represented by this way will have a similar modulus due to apply the same real and imaginary part for each point [12].

Regarding the Quadrature Phase Shift Keying (QPSK), it is another form of Phase Shift Keying where couple bits are utilized for each used subcarrier. Hence, employing the same BW, one of four phase shifts (0, 90, 180, or 270) degree are possibly selected to represent the generated points.

In comparative to the BPSK scheme, this modulation format is used to get a dual bit rate for a similar maintaining boundary of BW. Nevertheless, to achieve the same BPSK level of the BER, the QPSK need to twice the received power due to transmitting two bits simultaneously.

The investigations about employing the digital modulation with the OFDM particularly low modulation schemes like BPSK, QPSK were started in last three decades. The performance of the OFDM based BPSK, QPSK, etc., was explored in terms of the BER and channel capacity. Hence, a number of studies discussed in detail the low bit rate transmission and showed the trade-off relation between the size of achieved capacity and the SNR, particularly, a power of the signal and how it improved to recover the overall signal in the receiver.

In addition, the effect of transmission through noisy channels has been shown, thus, the performance (BER, bit rate) for a transmitted signal has been investigated under the AWGN channel response showing how the added noise can eliminate the performance of transmission for gained capacities at intended limits of errors.

Moreover, the influence of the AWGN is experimented for different received power at the receiver side demonstrating the impact of the applied power on the received signal. Furthermore, some researchers investigated the effect of expanding the available BW and how that can improve the bit rate for low modulation format [13].

However, those studies discussed the transmission performance of both BER and channel capacity for only low bit rate modulation format, thus, explore the activity for low capacity systems.

It's worth noting that, whenever the wanted bit rate is increased (16 PSK or more), it is strongly advisable to move to the Quadrature Amplitude Modulation (QAM) due to larger distances among the adjacent points in the constellation table. On the other hand, detect the points with both the phase and modulus due to having different amplitudes rises the involved complexity of the demodulation system.

From a general perspective, the QAM is a modulation mechanism, in which, both the amplitude and phase are utilized to express the modulated points for each frequency subcarrier.

The combination of shifted amplitude and phase produces higher modulation system with enhanced data representation. Hence, the mixture of changing phases and amplitudes can improve data transmission efficiency.

In addition, such a kind of modulation can supply a higher bit rate than formerly mentioned modulation scheme (PSK) for digital communications of mobile. Thus, increasing the number of the indicated points at the constellation table.

Besides, a square grid arrangement with equivalent horizontal and vertical spacing is applied for representing the QAM points, where the amplitude can vary with the phase at the In-phase (I) and Quadrate-phase (Q) plane.

The possibility of transmitting further bits per symbol is achieved employing orderly higher modulation schemes. Hence, supply room for added points inside the constellation table, to raise the ability of transmitting a similar amount of data in a smaller BW.

However, this is come up with raising the level of noise since constellation points number is enlarged and the spaces among them are narrowed resulting in higher BER at the constellation table.

This, nevertheless, can improve the BW efficiency but with less reliability in comparative with the lower order of modulation formats. However, this more liable to noise system is treated by promoting the needed power of the signal at the receiver side.

Regarding the first most popular one of the QAM family, which is termed as the 16 QAM, number of studies started to appear exploring an alternative approach for the digital modulation depends on both phase and amplitude of the signal instead of the phase only.

The new style of the modulation represented a big movement in a wireless digital modulation due to its ability to improve BW efficiency. Hence, employ an extra number of bits per each sample (4 bits) to generate each complex number. The performance in terms of the BER and channel capacity for the OFDM based 16 QAM, were discussed extremely by a number of researchers.

In addition, the close relation between the SNR and BER investigated to show the effect of adding extra bit for each sample in the constellation table. Hence, the received electrical power was also another determination for the improved signal.

As a result, the received signal is not recovered perfectly unless increase the power at the receiver side. The utilized BW was increased side by side by upgrading the modulation format giving an extra improvement for the performance in terms of the maximum bit rate and the BER of system.

Moreover, the effect of testing the AWGN channel has been demonstrated showing how the performance of channel capacity for transmitted signal was decreased due to increased BER [14]. However, this study demonstrated the system performance for only 16 QAM and future market demands requires more and more capacity like 32, 64, 128, 256.

Recently, researchers investigated another developed modulation scheme for upgrading the digital wireless transmission systems. This investigation showed that by increasing the number of bits for each transmitted sample (5 bits), the performance in terms of capacity and BER for the OFDM based 32 QAM was raised.

Nevertheless, this improvement caused in increasing the required SNR to mitigate the internal constellation noise.

In addition, the influence of the unwanted signal (AWGN) was explored demonstrating the relationship between increased SNR and reduced BER due to a promoted power of the signal. Hence, the utilized power for the received signal was strengthened against the unwanted power of the noise.

Moreover, a number of the exercised received powers was examined clarifying how the SNR is improved side by side with the raised power at the receiver side.

Improving the level of modulation format was accompanied by a number of trials for expansion the utilized BW and get some extra enhancements for the performance of wireless transmission systems with indeed an added cost [15].

However, the hangry data applications still looking for employing further developed schemes to sustain high capacities at a good level of errors.

Very recent, new studies have demonstrated the influence of utilizing a modern developed modulation format (64 QAM). This advanced modulation which is considered lately by the Long-Term Evolution (LTE) can accommodate 6 bits per each sample. Thus, improve the spectrum efficiency of the employed BW.

In addition, researchers have focused their attention on the trade-off relation between the enhanced performance of both maximum bit rate and the BER, and the expanded limits of power for a received signal. Hence, clarify the ability of the developed system to address the generated noise due to raised interference among the intended values of the signal based the new modulation format.

On the other hand, the impact of the uniform AWGN in a wireless channel showed how the developed system needs an extra power to mitigate the signal at the receiver side.

Besides, different values of received powers have been tested exploring how the optimal value of the SNR was achieved [16].

The promoted modulation schemes synchronized with indicated efforts to invest in the obtainable BW. Hence, seeking further solutions to enhance the efficiency of spectrum for acceptable boundaries of errors.

Nevertheless, the need for a higher data rate of greedy data services is still required by the future generation of wireless networks.

In a wireless mobile radio channel, the type of attenuation for a transmitted signal is decided by the nature of the propagated signal and the features of the utilized channel. Hence, various transmitted signals will suffer from various kinds of fading effects due to the relation between the parameters of both the channel like delay/doppler and the signal like symbol interval.

To explore the underlying behaviour of a wireless communications system in the absence of nonlinearity factors of noise, the AWGN model is employed. Hence, diving the internal behaviour of the spread signal before considering the other complicated phenomena like interference, dispersion, etc.

This basic, tractable and organized mathematical model added a steady spectral density linearly which is represented as watts per Hertz of the available BW over a specific transmission media.

In addition, the white noise with a Gaussian distribution of modulus is assumed for a flat channel with a fixed transmitter and receiver [17].

To illustrate the NLoS effect of radio signal propagation between the transmitter and receiver, the Rayleigh channel is discussed.

Such statistical model is important for signals which basically suffer from a similar attenuation, but different phases of arrival. Hence, for a wireless signal passed through the communication channel, the changed amplitudes of a signal for various coming times are assumed according to the Rayleigh distribution. Thus, the Rayleigh model is mainly applied due to the multipath phenomenon at the receiver [17].

To explain more about the multipath propagation, in mobile communications, the multipath phenomenon occurs when a radio signal is received at antenna by more than one path. This, as a result, influences the quality of telecommunications due to shifted phases of the same signal.

It's worth noting that, a maximum bit rate is directly related to the intended level of the BER. Hence, the BER considers as a significant parameter in assessing the performance of the digital wireless transmission system. Thus, the BER is utilized to examine the overall performance of the electrical back-to-back transmission system (transmitter, receiver and wireless channel) referring to the rate of error occurrence for a delivered signal.

## III. System Model

As is seen in Fig. 1, utilizing the OFDM environment, the digital signal is converted as the first step to its frequency domain using the proposed system of higher modulation formats (128 and 256 QAM), thereafter, the signal is transformed to the time domain using Inverse Fast Fourier Transform (IFFT).

The time guard interval, Cyclic Prefix (CP) is added at this stage of transmission to supply an appropriate level of protection for the transmitted symbols. Hence, these offered intervals of time play a big role in preventing any probable interference between transmitted symbols. Thus, sustain the
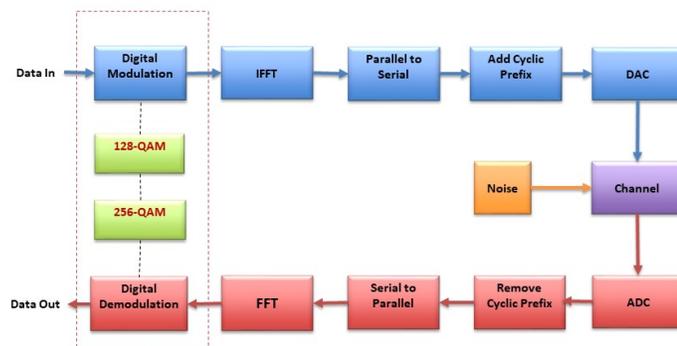


Fig. 1. Proposed OFDM transceiver including the higher modulation formats.

robustness of the transmission system against a probably happened Inter-Symbol Interference (ISI).

Utilizing a proper sampling frequency digital to analog converter ($F_{DAC}$), the prepared signal is converted to analog domain to be broadcasting later by an antenna.

In the receiver side, inverse processing is applied. The wireless received signal is converted back to its digital domain using a corresponding sampler analog to digital converter ($F_{ADC}$). The CP's that followed symbols of the OFDM are removed periodically to conclude the main signal in the time domain.

Employing Fast Fourier Transform (FFT), the received samples are turned back into the frequency domain. Finally, the digital values are recovered from their complex numbers using the demodulation operation.

The main concept beyond the OFDM is to divide a signal stream with high data rate into a group of low data rate which are simultaneously transmitted through a group of subcarriers.

As the modulation / demodulation process represents a key stage of the OFDM system, it is quite important to discuss, from a mathematical perspective, some related concepts that can affect directly the transmission operation.

In the transmitter side, particularly, in the modulation process, each token of the binary digits is converted to a complex number which is expressed in Cartesian forms as follow [18]:

$$C_k = I_k + jQ_k \quad (1)$$

where $j = \sqrt{-1}$, and $I$, $Q$ represents the real and imaginary parts for $k^{th}$ complex number respectively.

In this context, the mathematical relationship between the amplitude (A) and a produced complex number is clarified as follow [19]:

$$A_k = \sqrt{I_k^2 + Q_k^2} \quad (2)$$

In addition, the relationship between the phase (Θ) and a complex number is illustrated as follows [18]:

$$\Theta_k = \arctan(\frac{Q_k}{I_k}) \quad (3)$$

Thus, both the amplitude and phase are concluded based on the complex number formula in the Cartesian domain. This is, however, not considered for all mathematical operation of the complex number. Hence, the Cartesian form can fit for addition and subtraction operations more than multiplication and division.

To address this issue, the Polar form is accounted to state both multiplication and division in an easier way. The complex number formula in the Polar domain is demonstrated as follows [18]:

$$C = Ae^{j\Theta} \quad (4)$$

Then, according to Euler's formula [18]:

$$e^{j\Theta} = \cos\Theta + j\sin\Theta \qquad (5)$$

This, as result, leads to clarify a complex number in a sinusoidal form as follows [18]:

$$C = A\cos\Theta + Aj\sin\Theta \qquad (6)$$

Consequently, $I$ (real part) = $A\cos\Theta$ and $Q$ (imaginary part) = $A\sin\Theta$.

In the receiver side, particularly, when the channel response is counted, it is very necessary to make use of channel estimation concepts due to probably higher changes in both phase and amplitude of the transmitted signal over a wireless channel.

Thus, there is no benefit of applying the demodulation process, before achieving such a kind of corrections. According to this principle, the pilot-aided scheme with two key steps is applied to perform the equalization for the received signal as follows:

The first step is calculating the inverse of the channel response to mitigate the bad effect of the channel. Then, averaging the received pilots over $m$ transmitted samples to reduce the noise of signal and then to recover the signal itself.

To perform the equalization in a simple way, the Channel Transfer Function (CTF) is estimated. This, however, needs to make both the amplitude and phase of each pilot sample are known. Hence, the training samples $T_k$ are transmitted periodically with recognized magnitude and phase as the following [19]:

$$T_k = X_k * e^{j\phi_K} \qquad (7)$$

where $X_k$, and $(\phi_K)$ represent the amplitude and phase of the $K_{th}$ transmitted sample, respectively.

The equivalent received training samples $R_k$ is [16]:

$$R_k = Y_k * e^{j\varphi_K} + N_k \qquad (8)$$

where $Y_k$, and $(\varphi_K)$ represent the amplitude and phase of the $K_{th}$ received sample, respectively and $N_k$ is the $K_{th}$ sample noise.

By making use of the identified transmitted and received training samples, the CTF in frequency domain $E_k$ , is determined as the following [19]:

$$E_k = \frac{R_k - N_k}{T_k} \qquad (9)$$

Thus, the estimated CTF can be explained as the following:

$$\hat{E}_k = E_k + \frac{N_k}{T_k} \qquad (10)$$

Finally, the other received complex values in frequency domain $\check{V}_i$ is equalized by multiplying the inverse of the estimated CTF, $\hat{E}_k^{-1}$ as follow [19]:

$$\check{V} = \hat{E}_k^{-1} * \acute{V}_i \qquad (11)$$

Regarding the probable shapes of complex numbers on the constellation table, the following equation achieves this purpose [18]:

$$Y = 2^X \qquad (12)$$

where $Y$ represents the total number of probabilities which might be assigned for each sample employing $X$ of input bits.

Despite each spectrum of subcarrier can coincide with the others, the ability to extract each subcarrier is achieved over the digital signal processing. Hence, this overlapped property of subcarriers increases the spectrum efficiency of the current OFDM in comparative with the previous multicarrier design of waveform.

Thus, the OFDM technique splits a wireless channel into smaller subcarriers each one is modulated with an amount of data according to the applied modulation format. The improved efficiency of the OFDM spectrum is gained due to applying the orthogonality between adjacent subcarriers. Consequently, obtain a larger benefit for the same offered BW.

According to Shannon's theorem, the channel capacity represents the maximum achieved bit rate with a vanished amount of errors as follows [20]:

$$Capacity = BW.\log_2(1 + SNR) \qquad (13)$$

Besides, the SNR is gained as follow:

$$SNR = P_{receiver}/P_{noise} \qquad (14)$$

As such, to improve the obtained channel capacity, it's better to make an extra investigation in the field of BW efficiency than going with expanding the offered BW itself since such expanded resources require a highly increased cost.

In this study, the modulation operation is explored in the presence of AWGN and Rayleigh channels. Hence, the received signal $R$ is composed of the transmitted signal $T$ multiplied by the response of the channel $E$, where $E = 1$ with the AWGN and $E \neq 1$ with the Rayleigh. Furthermore, the composition is combined with the AWGN, which is represented here as $N$, as follows [21]:

$$R = ET + N \qquad (15)$$

It's worth noting that, the assigned length $k$ of the utilized IFFT/FFT represents the total number of subcarriers presented in the system. This is, nevertheless, not actually account for a real number of the employed subcarriers. Thus, when the system is supported with $k$ subcarriers, the utilized subcarriers for data transmission is about half of them. This is due to apply an equivalent number of conjugates which are employed in converting the signal from the frequency domain to the time domain.

TABLE I. KEY PARAMETERS OF OFDM SYSTEM

| Parameter | Value |
|---|---|
| FFT size | 40 |
| Sampling frequency | 40 MHz |
| Number of subcarriers | 15 |
| Cyclic prefix | 0.25 |

In addition, a specific duration is assigned for the CP to ensure more reliable transmission system.

## IV. EXPERIMENTAL RESULTS

In this part, numerical simulation with the MATLAB is introduced for the promoted transmission system demonstrating the advantages and disadvantages of adding higher modulation schemes for future mobile networks. Thus, examining experimentally the performance in terms of a maximum bit rate and the BER for the OFDM based 128 QAM and 256 QAM to explore the behaviour of wireless channel responses over variant modulation schemes (low and high).

The experiment is set up herein for 15 frequencies of subcarriers where the advanced modulation formats are considered side by side with the conventional modulation schemes. Besides, as the modulation system is the core of this proposed system, the number of stated bits for each used subcarrier is accurately fixed in accordance with the corresponding level of employed modulation schemes. The optimal parameters of the configured model that utilized for a wireless electrical back-to-back system is basically achieved under the conditions shown in Table I.

As is shown in Fig. 2, the newly proposed modulation formats, 128, and 256 QAM can improve the transmission bit rate compared to currently applied, 64 QAM, to about 16% and 33%, respectively.

In addition, in comparison with the 32 QAM and 16 QAM, the higher modulation 256 QAM increase the overall channel capacity by about 60% and 100% sequentially.
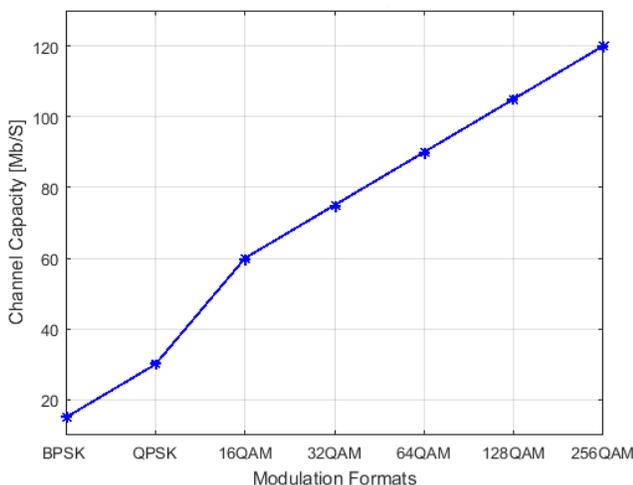
Moreover, 3- and 7-times enhancement are registered for the higher modulation (256 QAM) in comparison to both QPSK and BPSK.

This, however, comes up with raising the level of errors for transmitted samples due to enlarged interference at the constellation map. As such, a higher BER is presented for each utilized subcarrier because of the need for higher received power to sustain the required limits of errors at $10^{-3}$.

The experiment illustrates also, how the BER of the employed subcarriers is varied for diverse schemes of modulation under a similar level of the SNR and over the AWGN channel.

As is clear in Fig. 3, the BERs of the applied subcarriers are calculated for the currently employed modulation format (64 QAM). The measured BER is essentially achieved with an appropriate level of the SNR, which is equivalent herein to 23 dB. Accordingly, the acceptable limits of errors are gained due to supplying a suitable received power.

This scenario, however, is not typically fit for higher formats of modulation like 128 and 256 QAM.

According to Fig. 4, utilizing the same level of the SNR, the BER gets worse with increasing the modulation technique to the 128 QAM.

Furthermore, as is shown in Fig. 5, an extra rise of the BER level is recorded with moving to the higher modulation format (256 QAM) and keep the SNR at 23 dB. This fundamentally, due to reduce the distances between the adjacent samples on the constellation table resulting in inability in recognizing the received signal of the enlarged modulation schemes.

To explain the impact of increasing the supplied levels of received power on the signal strength, the SNR is raised, firstly, by about 3 dB to be more suitable for the 128 QAM.

As is clarified in Fig. 6, the new drawn map of the BERs refers to a feasible enhancement at the gained averages of errors with 128 QAM due to reduce the influence of interference among the constellations.



Fig. 2. Maximum achieved bit rate of different modulation formats including the 128 and 256 QAM.
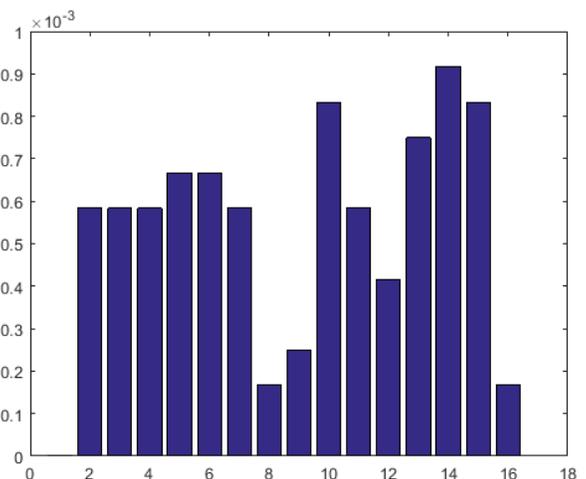


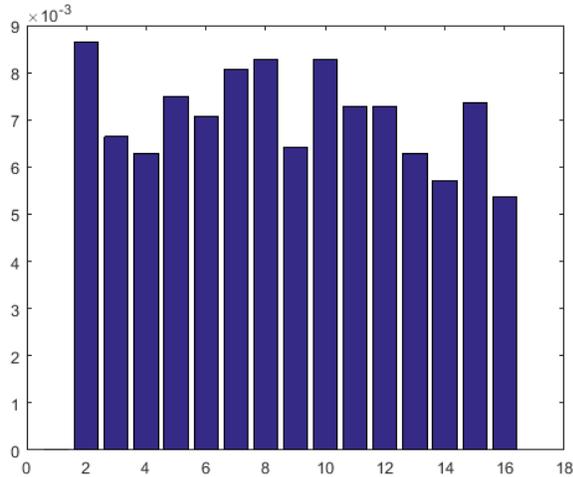Fig. 3. BERs' map of utilized subcarriers with 64 QAM and SNR=23.

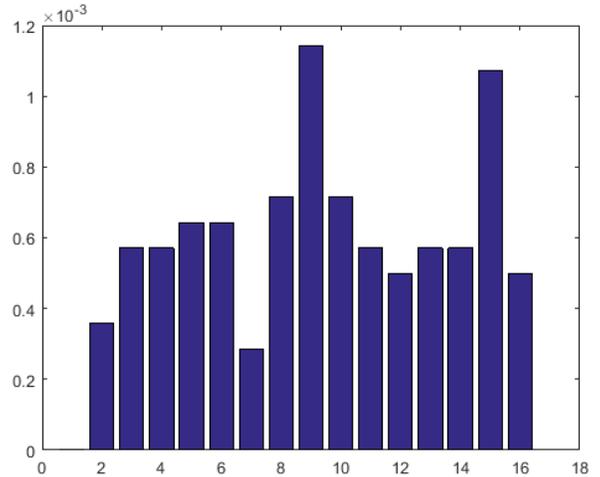Fig. 4.    BERs' map of utilized subcarriers with 128 QAM and SNR=23.



Fig. 6.    BERs' map of utilized subcarriers with 128 QAM, SNR=26.

This level of the SNR, as is shown in Fig. 7 is, however, still not able to support the accepted limit of the BER for the 256 QAM, since the improved power yet less than the induced interference among the samples at the receiver side.

Consequently, to obtain an acceptable amount of errors rate, as is seen in Fig. 8, the SNR is further expanded to about 29 dB, hence, an extra 3 dB improvement is realized for the employed SNR to sustain the BER limits of applied subcarriers with the higher order modulation (256 QAM).

For more details, the relationship between the overall calculated BER and the set level of the SNR is investigated for variant modulation schemes and over different channel models (AWGN, Rayleigh).

As is clear in Fig. 9, the performance in terms of the SNR and BER is firstly examined for the presently utilized 64 QAM and in presence of both the AWGN and Rayleigh noise channels.

It's, also, noticed that the obtained results of the 64 QAM modulation format are varied depending on the channel model, hence, the electrical signal is received over the AWGN with a lower level of power in comparison with the Rayleigh channel. Thus, the required limit of SNR is increased with the Rayleigh due to the impact of channel response on the transmitted signal.

Regarding the primarily proposed modulation formats (128 QAM), it is observed in Fig. 10, that agreeable limit of the BER calculated for the 128 QAM is achieved with a higher level of the SNR. This basically results from the reception system attempt to recover the signal in good condition. Thus, over the both considered channels, the BER is improved with increasing the SNR due to mitigating the influence of convergent samples at the constellation table.

In addition, the dynamic range of the utilized SNR shows a 3dB difference between the minimum limits of received power for both the currently employed 64 QAM and the proposed 128 QAM.
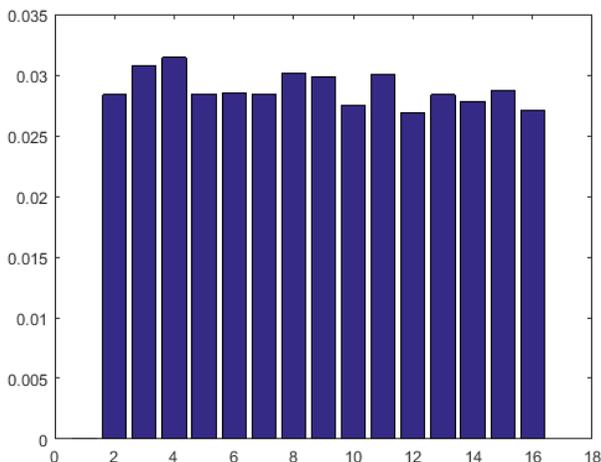


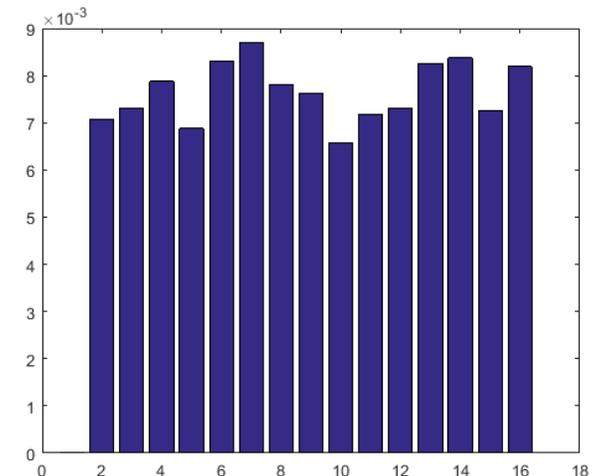Fig. 5.    BERs' map of utilized subcarriers with 256 QAM, SNR=23.



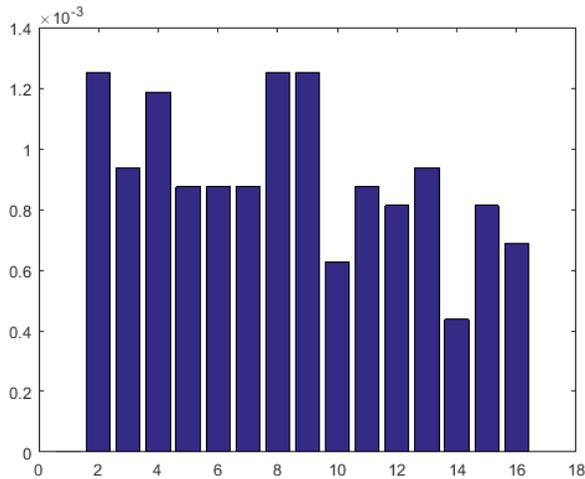Fig. 7.    BERs' map of utilized subcarriers with 256 QAM, SNR=26.

Fig. 8.    BERs' map of utilized subcarriers with 256 QAM, SNR=29.



Fig. 10.    Relation between BER and SNR for 128 QAM over AWGN and Rayleigh.

Regarding the secondly proposed higher order modulation scheme (256 QAM), as is obvious in Fig. 11, the intended limit of receiving power is further increased with promoting the modulation to a better degree of transmission, hence, an extra larger level of the SNR is recorded with the typically acceptable average of error.

Besides, for both adopted channels (AWGN and Rayleigh), whenever, the supplied levels of the SNR are raised, the BER limits are enhanced accordingly. This is fundamentally occurred since the impact of constellations interference is significantly reduced at the receiver side.

It's worth noting that the dynamic range of the applied SNR is even more affected, where about 6 dB variance is registered between the actually used modulation scheme (64 QAM) and the herein proposed (256 QAM).

As is seen in Fig. 12, the trade-off relation between the

gained BER and offered SNR is digitally processed for variant modulation schemes (QAM and PSK). Hence, the overall system performance in terms of the BER and SNR is presented for both the conventional modulation techniques as the BPSK, QPSK and the advanced configurations of modulation like the 128 and 256 QAM over the AWGN channel.

In addition, the simulation results show that, whenever, the received power of the signal is increased the BER is decreased until achieve the accepted criteria of tested modulation scheme.

As is seen in Fig. 13, the channel model can play a big role in degrading the system performance for both the SNR and BER, thus, the necessity for adapting the required power of the receiver at agreed limit of error is appeared. Hence, the required SNR to achieve a good reception with a modulation format over the AWGN is much lower in compare with the signal passes through a noisy mode like the Rayleigh channel model.
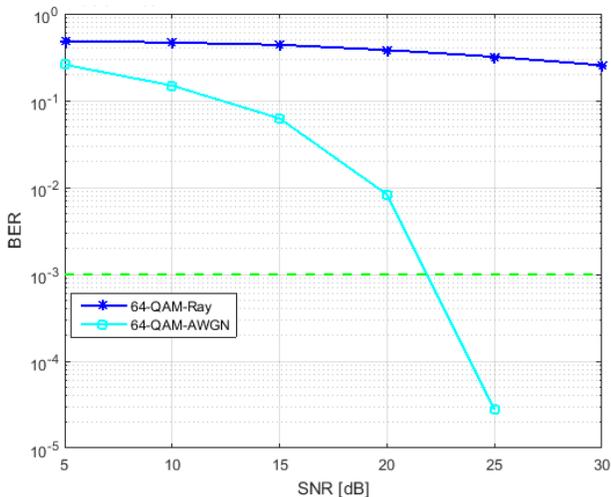


Fig. 9.    Relation between BER and SNR for 64 QAM over AWGN and Rayleigh.



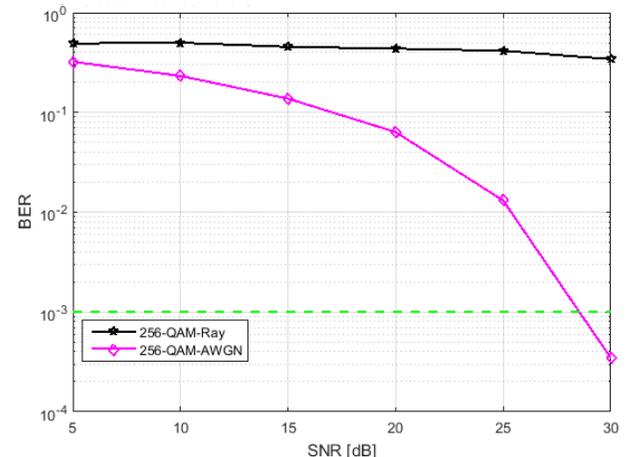Fig. 11.    Relation between BER and SNR for 256 QAM over AWGN and Rayleigh.
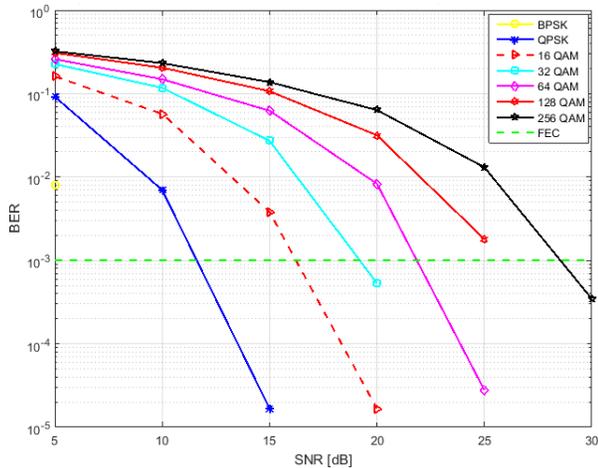
Fig. 12.   System performance in terms of BER and SNR for most common modulation formats including both 128 and 256 QAM over AWGN channel.

This essentially, because of the induced channel response that directly impacts the key parameters of the transmitted signal. Hence, the phase and amplitude of the sample are changed according to the type of channel response.

As a result, to provide a high data rate transmission along with reliable communication system the trade-off relation between modulation schemes and received power is adjusted precisely over the Rayleigh channel.

It's also shown in both Fig. 12, and Fig. 13, that the minimum limits of the SNR that improve the BER to $10^{-3}$ totally differ between transmission over the AWGN than Rayleigh channel. Hence, the modulation formats with the LoS channel can offer a better performance (BER, SNR) compared to the NLoS channel.

Besides, the sharply decreased BER is achieved by increasing the SNR levels in both Rayleigh and AWGN channels
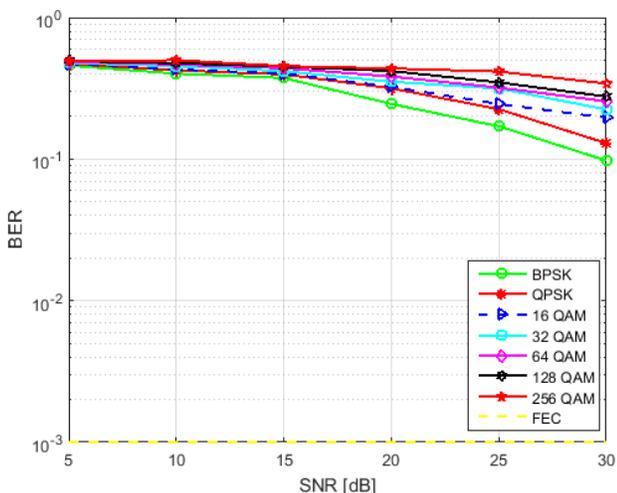
considering that BER is higher with the Rayleigh than the normal AWGN.

From these considered scenarios, it's found that the dynamic range of received power per each transmitted signal depends on modulation format type and utilized channel model.

It's worth noting that, among these 7 investigated modulation formats, the BPSK can be applied with a lower power requirement, while the 256 QAM requires a higher level of power due to the increment of bits number for each transmitted sample. Hence, for those applications which prefer power than the BW efficiency, the low order modulation format can be a more suitable solution than the higher modulation schemes and vice versa. Thus, the modulation techniques are employed depending on the type of applied application.

As such, with the future generation of mobile, the higher modulation schemes are very recommended due to the increased need for a higher bit rate.

This, however, comes up with some technical difficulties such as large-power reception which must be supplied by the conventional big cells to recover the transmitted signal with the higher modulation formats.

To solve such a challenge, nowadays technology, like the small cells can be a good solution for the modern generation of mobile networks which basically aims to provide huge amounts of data efficiently [22].

To test the proposed system, a MATLAB simulation is utilized to check whether the transmission operation with the proposed modulation formats over different channel models is achieved correctly or not.

In this testing process, the constellation table is considered to demonstrate the behaviour of electrical back-to-back transmission system toward the digital signal. Hence, depending on the recognizable receiving of the constellates, the reception can be confirmed, and the level of a suitable receiver power is decided.

The constellation map decision is changed, based on key factors like type of modulation technique and employed channel model.

Accordingly, the required SNR is adapted to treat any probable distortion of the received signal.

As it is mentioned before, the number of probable positions on the constellation map is calculated according to equation (12).

The experimental work shows the examined results for different cases of the constellation table under the most common kinds of modulation formats including the higher order as the 128, and 256 QAM and in presence of variant types of noise like the AWGN and Rayleigh.

As it is clear in Fig. 14, using the phase attribute, the probability of points appearance from transmitting 1 bit over the applied system is only two. These points which basically organized in two options are achieved at the accepted limit of the BER with SNR equivalents to 7 and 58 dB over the AWGN and Rayleigh channels respectively.



Fig. 13.   System performance in terms of BER and SNR for most common modulation formats including both 128 and 256 QAM over Rayleigh channel.

Fig. 14.   BPSK modulation format of receiver.



Fig. 16.   16 QAM modulation format of receiver.

The increased number of bits (2 bits), as is shown in Fig. 15, leads to raising the probability of achieving extra positions at the constellation table to four. Hence, depending on the phase feature, the organized constellations can reflect the accuracy of transmission over channels models. Thus, at a BER equal to $10^{-3}$ the constellation table gives four sorted positions with SNR level variant between 12 dB for the AWGN and 63 dB with the Rayleigh.

In addition, as is seen in Fig. 16, whenever the number of bits is enlarged (4 bits), the appearance chances of receiving new complex numbers are increased resulting in decreasing the distances between the adjacent constellations. Thus, utilizing both phase and amplitude, 16 opportunities are offered in a sorted way at the map, referring that received signal is accepted at the required limit. The recorded SNR for a good reception with 16 QAM and over the AWGN is 16.5 dB while a higher ratio is accounted to the Rayleigh at 67.5 dB.

Besides, as it is noticed in Fig. 17, employing 5 bits for each sample leads to constructing the 32 QAM modulation, which also depends on both phase and amplitude to accommodate the expanded cases of complex numbers at the constellation table. Thus, 32 different cases are received announcing that reached signal is correctly positioned with the gained BER. The agreeable limits of errors are obtained with diverse SNR levels, where 19.5 dB is assigned for the AWGN and 70.5 dB for the Rayleigh.

Moreover, as it is clear in Fig. 18, by raising the number of bits to six, the currently utilized modulation format, 64 QAM is achieved. Thus, employing the aspects of phase and magnitude, 64 arranged shapes are drawn at the constellation map proving that digital signal is received optimally. The registered values of the obtained SNR for this arranged reception over the AWGN and Rayleigh channels are 22.5 dB and 73.5 dB respectively.



Fig. 15.   QPSK modulation format of receiver.



Fig. 17.   32 QAM modulation format of receiver.

Fig. 18.    64 QAM modulation format of receiver.



Fig. 20.    256 QAM modulation format of receiver.



Fig. 19.    128 QAM modulation format of receiver.

Regarding the newly proposed higher order modulation formats, firstly, as it is shown in Fig. 19, seven bits are specified for each transmitted sample to shape a new coordinated constellation table with 128 available options for representing the transmitted data.

Based on the phase and amplitude utilities, the correctly accommodated samples indicate that complex numbers of the transmitted signal are received successfully. The gained levels of the SNR declare that the AWGN needs to about 25.5 dB while the Rayleigh comes up with 76.5 dB to get the acceptable limit of the error's ratio.

Eventually, as it is seen in Fig. 20, eight bits are decided for every transmitted sample introducing an extra number of the selections reach to 256 obtainable codes.

The harmonic reception of the transmitted samples demonstrates that the receiver side is working effectively. Hence, depending on the correct obtaining of both phase and magnitude at the 256 QAM constellation map, the standard limit of errors is accepted.

Higher levels of the computed SNR are assigned with varied models of channels where 29 dB is accounted to the AWGN and around 80 dB is consumed with the Rayleigh.

## V.    CONCLUSION

In this study, new configurations for the PHY of future mobile technology are provided by considering higher approaches of modulation than the currently applied in the telecommunication systems. More specifically, advanced modulation formats like 128 QAM and 256 QAM employing the current waveform (OFDM) is proposed, investigated and evaluated. The performance analysis of the relatively higher schemes of modulation are theoretically and experimentally assessed under different channel conditions (AWGN and Rayleigh). The experimental work shows that higher modulation schemes can supply more data bit rate but le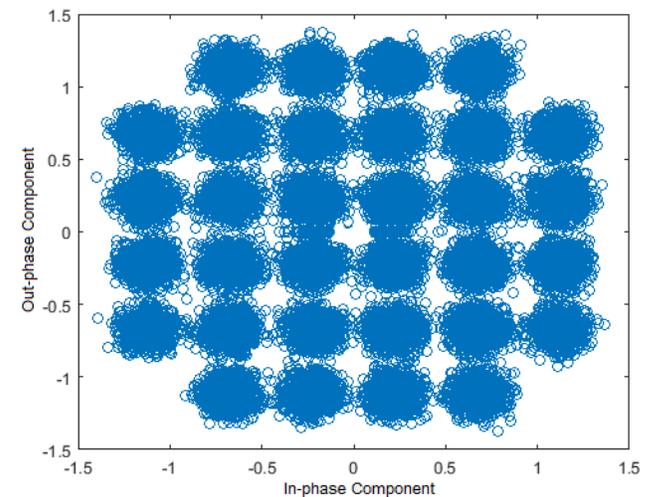ss transmission reliability due to the increased rate of errors compared to lower modulation formats. In addition, it reveals that, for all examined modulation formats, the achieved level of the error rate is minimized with the AWGN while maximized with Rayleigh noise channel. Moreover, the findings clear that the higher order modulations are more suitable for purposes that need great BW efficiency but down power efficiency, while lower order modulations are more recommended for purposes that include less power and error but lower BW efficiency. Thus, the promoted modulation schemes can improve the transmission channel capacity of a future wireless mobile system in comparison with the currently applied format (64 QAM). However, the effect of adding an extra bit for each transmitted sample is directly related to the raised requirement of power for each received complex number due to increased interference among transmitted points in the constellation table. Hence, increased levels of the BER synchronize with the enlarged schemes of modulation. Thus, the gained values of the SNR reflect the close relationship between increased bit rate and required receiving power. As results, the higher modulation formats are not preferred alone for the future generation of mobile networks, since the new

predicated scenarios expected higher channel capacities with lower levels of power. Consequently, sophisticated technology like the small cell is highly recommended, with the high-order modulation formats, due to their efficient ability in increasing the coverage of transmission for future wireless networks of mobile. The achieved results are evaluated mathematically and physically utilizing a MATLAB simulation code.

## REFERENCES

[1] Y. U. E. Xiao, H. Haas, and S. Member, "*Index Modulation Techniques for Next-Generation Wireless Networks,*" vol. 5, pp. 16693–16746, 2017.

[2] M. R. Kadhum, T. Kanakis, A. Al-sherbaz, and R. Crockett, "*Digital Chunk Processing with Orthogonal GFDM Doubles Wireless Channel Capacity,*" Advances in Intelligent Systems and Computing, vol 857, pp. 719-731, Nov. 2018. doi:https://doi.org/10.1007/978-3-030-01177-2-53.

[3] I. Chandran and K. A. Reddy, "*Comparative analysis of various Channel Estimations under different Modulation Schemes,*" vol. 1, no. 3, pp. 832–837, 2017.

[4] K. A. Reddy, "*Evaluation of BER for AWGN , Rayleigh Fading Channels under M-QAM Modulation Scheme,*" pp. 3081–3086, 2016.

[5] M. R. Kadhum, T. Kanakis, and R. Crockett, "*Dynamic Bit Loading with the OGFDM Waveform Maximises Bit-Rate of Future Mobile Communications,*" Advances in Intelligent Systems and Computing, Proceedings of computing conference 2019, London, 2019.

[6] M. Barnela, "*Digital Modulation Schemes Employed in Wireless Communication: A Literature review,*" vol. 2, no. 2, pp. 15-21, 2014.

[7] M. R. Kadhum, T. Kanakis, and R. Crockett, "*Intra-Channel Interference Avoidance with the OGFDM Boosts Channel Capacity of Future Wireless Mobile Communication,*" Advances in Intelligent Systems and Computing, Proceedings of computing conference 2019, London, 2019.

[8] C. U. Ndujiuba, O. Oni, and A. E. Ibhaze, "*Comparative Analysis of Digital Modulation Techniques in LTE 4G Systems,*" vol. 5, no. 2, pp. 60–66, 2015.

[9] W. Jin et al., "*Improved Performance Robustness of DSP-Enabled Flexible ROADMs Free from Optical Filters and O-E-O Conversions,*" IEEE J. Opt. Commun. Netw., vol. 8, no. 8, p. 521, 2016.

[10] T. Hwang, C. Yang, S. Member, G. Wu, S. Li, and G. Y. Li, "*OFDM and Its Wireless Applications: A Survey,*" vol. 58, no. 4, pp. 1673–1694, 2009.

[11] M. R. Kadhum, "*New Multi-Carrier Candidate Waveform for the 5G Physical Layer of Wireless Mobile Networks,*" IEEE, Proceedings of Wireless Days 2019, Manchester, 2019.

[12] N. Sood, "*BER Performance of OFDM-BPSK and -QPSK over Nakagami-m Fading Channels,*" no. 1, pp. 88–90, 2010.

[13] G. H. Im, D. D. Harman, G. Huang, A. V. Mandzik, M. H. Nguyen, and J. J. Werner, "*51.84 Mb/s 16-CAP ATM LAN Standard,*" IEEE J. Sel. Areas Commun., vol. 13, no. 4, pp. 620–632, 1995.

[14] G. Song and L. Cuthbert, "*Physical Layer Simulations of the Broadband Wireless LANs Based on OWSS Signalling Technique,*" pp. 442–445, 2010.

[15] R. Yang, J. Bian, and C. Zhu, "*Fixed point Dual Circular 32-QAM performance for Wireless USB,*" 2011 Int. Conf. Consum. Electron. Commun. Networks, CECNet 2011 - Proc., pp. 1310–1314, 2011.

[16] K. Nagarajan, V. V. Kumar, and S. Sophia, "*Analysis of OFDM Systems for High Bandwidth Application,*" pp. 168–171, 2017.

[17] M. M. Madankar and P. S. Ashtankar, "*Performance analysis of BPSK modulation scheme for different channel conditions,*" no. 2002, 2016.

[18] S. W. Smith, Digital signal processing. 1999.

[19] Q. Yang, S. Chen, Y. Ma, and W. Shieh, "*Real-time reception of multi-gigabit coherent optical OFDM signals,*" vol. 17, no. 10, pp. 873–879, 2009.

[20] M. Ghogho, D. McLernon, E. Alameda-Hernandez, and A. Swami, "*Channel estimation and symbol detection for block transmission using data-dependent superimposed training,*" IEEE Signal Process. Lett., vol. 12, no. 3, pp. 226–229, 2005.

[21] Z. Karaelmas and Ü. Zonguldak, "*Performances of M-PSK and M-QAM Modulated OFDM Signals over AWGN and Rayleigh Fading Channels,*" pp. 523–526, 2010.

[22] H. Haboobi, M. R. Kadhum, and A. Al-sherbaz, "*Utilise Higher Modulation Formats with Heterogeneous Mobile Networks Increases Wireless Channel Transmission,*" Advances in Intelligent Systems and Computing, Proceedings of computing conference 2019, London, 2019.

# Assuring Non-fraudulent Transactions in Cash on Delivery by Introducing Double Smart Contracts

Ngoc Tien Thanh Le[1], Quoc Nghiep Nguyen[2], Nguyen Ngoc Phien[3], Nghia Duong-Trung[4],
Thai Tam Huynh[5], The Phuc Nguyen[6], Ha Xuan Son[7]

CanTho University of Technology, Can Tho city, Viet Nam[1,2,4,7]
Center for Applied Information Technology, Ton Duc Thang University, Ho Chi Minh City, Vietnam[3]
Faculty of Information Technology, Ton Duc Thang University, Ho Chi Minh City, Vietnam[3]
FPT University, Can Tho City, Viet Nam[4,7]
Transaction Technologies PTE. LTD., Singapore[5]
University of Trento, Trento, Italy[6]

*Abstract*—The adoption of decentralized cryptocurrency platforms is growing fast, thanks to the implementation of Blockchain technology and smart contracts. It encourages the novel frameworks in a wide range of applications including finance and payment methods such as cash on delivery. However, a large number of smart contracts developed for cash on delivery suffer from fraudulent transactions which enable malicious participants to break the signed contracts without sufficient penalties. A shipper will involve in the system and place a mortgage to ensure reliability. A buyer also pledges an amount of money when making the order. Our process not only ensures the interests of a seller but also prevents a fraud shipper. The penalties will be made in two scenarios: (i) the buyer refuses to receive the commodities without any reliable reasons; and (ii) the shipper attempts to make any modification on the delivered goods during transportation. To help developers create more secure and reliable cash on delivery system, we introduce double smart contracts, a framework rooted in Blockchain technology and Ethereum, to tackle those mentioned problems. We also contribute our solution as an open source software that developers can easily add to their implementation to enhance functionality.

*Keywords*—*Cash on Delivery (COD); Blockchain; smart contract; Ethereum; e-commerce; online payment*

## I. Introduction

A basic problem for e-commerce is the exchange of digital goods for payment. The earliest solutions to this problem come from at least on the first day of the world wide web [1], e.g. online stores accept credit card payments. Because goods exchange and payment cannot happen simultaneously, there is an inherent tension and consequently, trust in the transaction is required. The seller must trust that the buyer will pay and the buyer must believe that the seller must deliver the goods. Traditionally, this necessity for trust has been solved by introducing a third party, e.g. a credit card.

Cash on Delivery (COD) allows customers to pay in cash when the products are delivered to their home or a location they choose. This is sometimes called a payment system because customers receive goods before making a payment. COD has become increasingly popular in recent years and been considered one of the main payment methods in many countries [2], [3], [4]. However, most published documents about COD have appeared in reports or magazines and/or on

the web, with a few scientific studies to date. Among research articles, most investigated payment methods is in general, rather than focusing on COD in particular. Transfer agents are often used as postal services, but usually, consumer and business shipments will be sent to COD by courier companies, commercial truck forwarders or organizations own delivery services. COD sales usually involve a delivered fee charged by the shipping agents and is usually paid by the buyer. In retail and wholesale transactions, shipments rely on COD-based payment method when the buyer does not have a credit account and the seller does not choose a payment method in advance. COD postal services [5] were first introduced in Switzerland in 1849, India and Australia in 1877, the United States in 1913, Canada in 1922 and the United Kingdom in 1926.

In a contrary direction to previous work, the authors propose an information contract implemented by a seller. A smart contract requires both the shipper and the buyer to place a deposit into an account. In the proposed protocol, the shipper first sends the deposit to get the seller's goods. Then the buyer sends the payment as well as the deposit itself. The seller then sends the key to open the digital commodities. The buyer verifies that the commodities are well received. If the commodities are in good condition, they will send a notification to a smart contract. Deposits made by both parties are only returned after successful transactions have been made, e.g. shippers successful delivered the commodities which are also successfully verified by buyers. To the best of our knowledge, this novel idea is firstly investigated and implemented by the authors.

## II. Related Work

One of the major problems of e-commerce globally is buying and selling between parties on the Internet. Krishnamachari *et. al.* [6] proposed a mechanism to implement a transaction with any asset by using digital keys and these processes did not require a reliable third party. In addition, the authors described a deposit transaction method for fraudulent and delivery transactions between the two parties in which the seller can use digital signatures to verify a transaction. Sellers and buyers use a pair of keys to verify goods. A smart contract is utilized to decide and handle seller-buyer relations

by increasing deposits. But the above article has not analyzed the shipping issue. More specifically, if the delivery does not comply with the commitment, the system cannot resolve.

Other researchers [7], [8] has proposed a mechanism based on the Ethereum Blockchain [9] or Son *et. al.* has introduced a mechanism [10] based on Hyperledger Fabric flatform [11] that relates to product transportation between sellers and buyers. In their approach, the carrier plays an important role. The transportation process consists of 2 steps: (i) a key is shipped with the product and handed over to the buyer, and (ii) the buyer will enter the key to confirm the reception in Smart Contract. Ether [12] will only be placed in the seller's account if the key entered by the buyer matches the key in the Smart Contract. Ether will be forwarded to the seller after successful confirmation. This solution is easy to implement because it is quite simple and depends on the key that the seller gives to the carrier. However, this leads to dependence on the belief that the shipper will not take advantage of that key before handing it to the buyer. Therefore, this solution is not recommended.

Hasan and Salah [13] has introduced a delivery process including buyers, sellers and carriers. If a carrier wants to deliver the commodities of a seller, he/she must place a mortgage payment in advance. This amount is usually double the value of the commodities. If the goods are successfully shipped, the money is paid to the parties. If it fails, the system will resolve the dispute by relying on delivery time, from which the system will make a decision without human intervention. Doubling the mortgage price of the products not only increases the transaction cost but also prevents mortgagees from cheating to avoid loss of money.

In this article, we look at several COD related issues and solve them by Blockchain technology. We introduce double smart contracts to address non-fraudulent transactions. A shipper will involve in the system and place a mortgage to ensure reliability. A buyer also pledges an amount of money when making the order. Our process not only ensures the interests of a seller but also prevents a fraud shipper. The penalties will be made in two scenarios: (i) the buyer refuses to receive the commodities without any reliable reasons and (ii) the shipper attempts to make any modification on the delivered goods during transportation.

## III. MATERIALS AND TECHNICAL BACKGROUND

In this section, the authors summarize the most related technical background that we implement in this work. Readers might also look at some interesting materials in the literature [14], [15], [16].

### A. Blockchain and Blockchain-based Smart Contracts

Blockchain is a list of developing logs, called blocks, linked by encryption. Each block contains the previous block's cryptographic hash function, timestamp, and transaction data. Each block has a block header and a body containing data and hash values of the previous block. The hash value is the result of a hash function. The hash function transforms data of any length into a fixed length string or numeric value, such as 256 bits (32 bytes) with SHA256. Blockchain is a technology that allows secure data transmission based on an extremely complex encryption system, similar to accounting books of a company where cash is closely monitored. In this case, the blockchain is an accounting ledger [17] that works in the digital field. A special feature of blockchain is that transactions are done at a high level of trust without disclosing information.

Blockchain-based smart contracts are proposed contracts that could be partially or fully executed without human interaction [18], [19], [20]. One of the main objectives of a smart contract is an automated escrow. An IMF (International Monetary Fund) staff discussion reported that smart contracts based on Blockchain technology might reduce moral hazards and optimize the use of contracts in general, but "no viable smart contract systems have yet emerged". Due to the lack of widespread use, their legal status is unclear [21]. Smart contract based on blockchain is being considered for many different types of transactions, from ubiquitous devices to real-time operational management structures for industrial products and data transfer in some applications including transaction finance. All types of business and management can participate in the network and use the properties of the Blockchain system to ensure transparency of stakeholders.

### B. Ethereum

Ethereum (ETH) [9] is an open-source, and a blockchain-based distributed computing platform. Ethereum uses functions of smart contracts. Smart contract of Ethereum written by Solidity [22] programming language. A smart contract is a computer protocol intended to digitally facilitate, verify, or enforce the negotiation or performance of a contract. A Smart contract allows the performance of credible transactions without third parties involved. These transactions are trackable and irreversible. A transaction in a smart contract will use Ether [23] unit to pay for the transaction. Like Bitcoin money has Satoshi and Kilobyte, USD money has dollars and cents, Ether is the currency of Ethereum's internal network. Nevertheless, Wei is the smallest unit changed from Ether unit used in a smart contract normally. In an Ethereum network, there are two addresses that we should be noted: an account address and a contract address. Each account address which is an external account has a corresponding personal key (private key). We can treat the private key as a password that we are the only ones who know. We need the address and private key pair to interact with the Blockchain. A contract address is also called a contract account which is controlled by the code stored together with the account. The contract address is determined at the time the contract is created. It is derived from the creator address and the number of transactions sent from that address, the so-called "nonce". Furthermore, every account has a balance in Ether, e.g. in Wei to be exact, 1 ether is $10**18$ wei, which can be modified by sending transactions that include Ether. Furthermore, one unit may be confused with Ether unit which is gas unit. When creating, each transaction is charged by a certain amount of gas, its purpose is to limit the amount of work needed to execute the transaction and pay for this execution at the same time.

### C. Smart Contracts

A cryptocurrency is a decentralized platform that a distributed ledger is used to interact with virtual money. A contract is an instance of a computer program that executes

on the Blockchain. Users transfer money by publishing transactions and interacting with contracts in the cryptocurrency network where information is propagated, data is stored among miners or network's nodes. An underlying cryptocurrency system supports the utilization of smart contracts. A smart contract contains program code, a stored file and an account balance. Any user can submit a transaction to an append-able-only log. When the contracted is created, its program code cannot be changed. An append-able-only log, called a blockchain, which imposes a partial or total arrangement on submitted transactions is the main interface provided by the cryptocurrency. Fig. 1 presents the idea of a decentralized cryptocurrency system and its components.
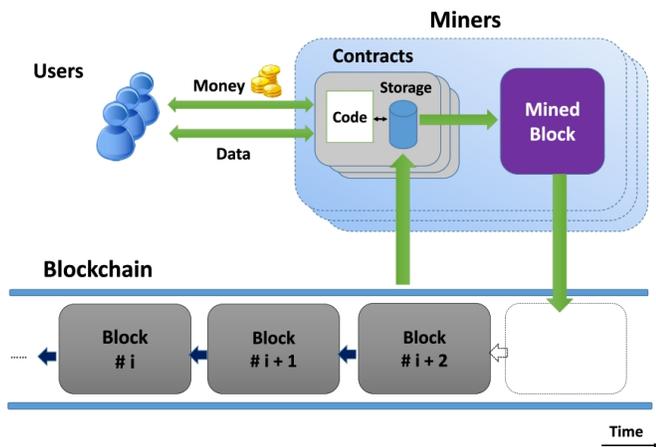


Fig. 1.   An illustration of smart contracts and Blockchain in a decentralized cryptocurrency system [24].

## IV.   PROPOSED METHODOLOGY

In this section, the authors introduce a general overview of the architecture with a highlight of the idea of double smart contracts. Then, we discuss the proposed architecture is more details. Finally, we present several important algorithms that serve as the backbone of our proposed architecture.

### A.  General Architecture

Our proposed COD system consists of several components, e.g. see Fig. 2. A seller gives package information into the system (step 1). In step 2, a buyer sends the request to buy this package. Then, the seller and buyer will sign a smart contract, called smart contract 1, in step 3. In smart contract 1, the buyer must mortgage a price to ensure that if the buyer is attempted to fraudulent or a shady businessman, the money will be sent to the seller. In step 4, the shipper sends a request to deliver the package. After the shipper sends the request, the seller and shipper must sign the smart contract 2 (step 5). In smart contract 2, the shipper deposits an amount of money that is equal to the price of the package to avoid shipper's lost or modification. In step 6, the shipper will send the package to the buyer and confirm a successful transaction. Next, the money deposited in the smart contract 2 will be sent automatically to the seller in step 7. Similar to step 7, smart contract refunds money to the buyer in step 8.



Fig. 2.   General design of our proposed architecture.

### B.  Detailed Design

In this subsection, the authors discuss the proposed COD architecture in more fine-grained details. An illustration of the detailed design is presented in Fig. 3. In step 1, a seller gives a package into the Blockchain system with the information about the package's name and its price. In step 2 and 3, a buyer sent a request to buy the package and deposit the money into the smart contract 1 to ensure personal interests. Next, in step 4 and 5, the seller sends the request to receive the package together with its information and the buyer's address. The seller and shipper will sign the smart contract 2. In this smart contract 2, the shipper must deposit an amount of money equal to the price of the package to ensure that the interest for the seller. Next, the shipper will send this package, step 6, to the buyer and then the buyer must pay money to the shipper in step 7. If the transaction goes well and the shipper receives cash from the buyer, the deposit amount of money in the smart contract 1 will be sent to the buyer in step 8, and the amount of deposit money in the smart contract 2 will be transferred to the seller in step 9. So the delivery has been successful. However, if transportation fails due to the shipper problems such as lost or item damaged during delivery in step 10. The money balance in the smart contract 2 will be transferred to the seller in step 11 and the money balance in the smart contract 1 will be transferred to the buyer in step 12. Furthermore, in step 13, if the buyer does not pay for the shipper, the deposit money in the smart contract 1 will be sent to the seller in step 14. Due to the errors generated by the buyer, deposit money in the smart contract 2 will be sent to the shipper in step 15. The steps that will be discussed in the paper hereafter are referred to these steps in the detailed design.

### C.  Algorithms

The main purpose of the algorithm (1) is to create a package with the request that the seller is the package's owner. In step 2, the package's information will be collected including *id*, *name* and *price*. The information is stored with an auto-id increase in step 3. An exemplification of Algorithm (1) is

Fig. 3.    Detailed Design of our proposed architecture.

---

**Algorithm 1** Create Package

1: **if** the seller is an owner and he/she creates the requirement **then**
2:     Set information: *id*, *package name*, and *price*.
3:     Save the information into an array with id package increasing one unit.
4: **else**
5:     Show an error message to customer.
6: **end if**

---

presented in Table I.

---

**Algorithm 2** Money Deposit

1: Get *package id*.
2: **if** caller address != owner address of the package **then**
3:     set price to deposit.
4:     **if** price == price of the package **then**
5:         set price success.
6:     **else**
7:         show an error message.
8:     **end if**
9:     save *balance* into a smart contract.
10:     **return**   value of *balance* equal to price deposited
11: **end if**

---

TABLE I.       AN EXEMPLIFICATION OF ALGORITHM 1.

| | |
|---|---|
| Transaction hash | 0x55593a8aa7aaee08a6ae3f354bd8eb7ebf8ace5ea 326c48eaf4fe6367bdee36b |
| From | 0xca35b7d915458ef540ade6068dfe2f44e8fa733c |
| To | Sell.AddPackage(string,uint256) 0xbbf289d846208c16edc8474705c748aff07732db |
| Gas | 3000000 |
| Transaction cost | 104565 gas |
| Execution cost | 82397 gas |
| Hash | 0x55593a8aa7aaee08a6ae3f354bd8eb7ebf8ace5ea 326c48eaf4fe6367bdee36b |
| Input | 0xb3a...00000 |
| Decoded input | { "string _name": "car", "uint256 _price": "10" } |
| Call to Sell.idPackage | {"0": "uint256: 1"} |

This main purpose of the algorithm (2) is to get *id* of the package in step 1. Step 3 sets the deposit money with a condition that the price must be equal the price of the package in the algorithm (1). In Table I, the price of the package is 10 Ether. In Table II the amount of deposit money with 10000000000000000000 Wei equal 10 Ether is transferred into a smart contract.

In algorithm (3), the system gets the value of balance in made in algorithm (2). In step 2, the caller of this function confirms the balance transfer to the specified address in step 3. In Table II, readers can see

TABLE II.     AN EXEMPLIFICATION OF ALGORITHM 2.

| | |
|---|---|
| Transaction hash | 0x21c1c983b876728c9607c47ee9c633907714ae4 99e9e077e06a2590f255fbe1e |
| From | 0x14723a09acff6d2a60dcdf7aa4aff308fddc160c |
| To | Sell.ApplyBuy(uint256) 0xbbf289d846208c16edc8474705c748aff07732db |
| Gas | 3000000 |
| Transaction cost | 64531 gas |
| Execution cost | 43076 gas |
| Hash | 0x21c1c983b876728c9607c47ee9c633907714ae4 99e9e077e06a2590f255fbe1e |
| Input | 0x114...00001 |
| Decoded input | { "uint256 id": "1" } |
| Decoded output | { "0": "uint256: 10000000000000000000" } |
| Logs | [] |
| Value | 10000000000000000000 wei |

---

**Algorithm 3** Money Transfer

1: get *balance* of the smart contract.
2: **if** caller address confirmed **then**
3:     transfer *balance* to specified address.
4: **end if**

---

the money deposit equal to 10 Ether unit from address *0x14723a09acff6d2a60dcdf7aa4aff308fddc160c*. So in Table III, the system will transfer money to various address and the balance will be return 0. An instance running of Algorithm (3) is presented in Table III.

## V.   EXPERIMENTAL SETUP

The smart contracts are created and tested using the Remix IDE [25], [26] which provides necessary tools for developing and debugging. As explained in the previous sections, gas is a measure of the amount of expenditure used to calculate the cost need to perform certain activities [24]. Each and every line of code will definitely require a certain amount of gas to calculate. We illustrate three different cases in this article and compare the total gas amount in three cases afterwards. In case 1, we illustrate a normal process and call it Transport Success. Next, in case 2, we discuss transaction errors due to shipper problems. Finally, we present how the system handles transaction errors because the buyer does not accept commodities.

TABLE III.     AN EXEMPLIFICATION OF ALGORITHM 3.

| | |
|---|---|
| Transaction hash | 0x075ea20859f1a0bf554341efc62f48c60010ae14 cf7c536879e94168a86f0d4c |
| From | 0x14723a09acff6d2a60dcdf7aa4aff308fddc160c |
| To | Sell.ConfirmToTransferBuyer() 0xbbf289d846208c16edc8474705c748aff07732db |
| Gas | 3000000 |
| Transaction cost | 30279 gas |
| Execution cost | 9007 gas |
| Hash | 0x075ea20859f1a0bf554341efc62f48c60010ae14 cf7c536879e94168a86f0d4c |
| Input | 0x11f...36795 |
| Decoded input | { } |
| Decoded output | { "0": "uint256: 0" } |
| Logs | [] |
| Value | 0 wei |

### A. Case 1: Transport Success

*1) Seller   0xca35b7d915458ef540ade6068dfe2f44e8fa733c creates a package on the smart contract 1:* see Table IV.

TABLE IV.     CASE 1: STEP 1

| | |
|---|---|
| From | 0xca35b7d915458ef540ade6068dfe2f44e8fa733c |
| To | Sell.AddPackage(string,uint256) 0x692a70d2e424a56d2c6c27aa97d1a86395877b3a |
| Gas | 3000000 |
| Transaction cost | 104693 gas |
| Execution cost | 82397 gas |

*2) Buyer 0xdd870fa1b7c4700f2bd7f44238821c26f7392148 agrees to deposit the amount of money equal to 10 Ether:* see Table V.

TABLE V.     CASE 1: STEP 2

| | |
|---|---|
| From | 0xdd870fa1b7c4700f2bd7f44238821c26f7392148 |
| To | Sell.ApplyBuy(uint256) 0x692a70d2e424a56d2c6c27aa97d1a86395877b3a |
| Gas | 3000000 |
| Transaction cost | 64531 gas |
| Execution cost | 43067 gas |

*3) Seller   0xca35b7d915458ef540ade6068dfe2f44e8fa733c gives the information of the package as well as the buyer's address to the shipper in the smart contract 2:* see Table VI.

TABLE VI.     CASE 1: STEP 3

| | |
|---|---|
| From | 0xca35b7d915458ef540ade6068dfe2f44e8fa733c |
| To | Cod.AddPackage(string,uint256,address,uint256) 0x0dcd2f752394c41875e259e00bb44fd505297caf |
| Gas | 3000000 |
| Transaction cost | 187265 gas |
| Execution cost | 163369 gas |

*4) Shipper 0x14723a09acff6d2a60dcdf7aa4aff308fddc160c agrees to deposit price and deliver the package:* see Table VII.

TABLE VII.     CASE 1: STEP 4

| | |
|---|---|
| From | 0x14723a09acff6d2a60dcdf7aa4aff308fddc160c |
| To | Cod.ApplyDeliver(uint256) 0x0dcd2f752394c41875e259e00bb44fd505297caf |
| Gas | 3000000 |
| Transaction cost | 44397 gas |
| Execution cost | 22933 gas |

*5) Transport   success:*   The   buyer   pays   in   cash directly to   the   shipper.   Then   teh   deposit   money   in smart   contract   1   will   be   transferred   to   the   buyer *0xdd870fa1b7c4700f2bd7f44238821c26f7392148*, see Table VIII.

TABLE VIII.     CASE 1: STEP 5

| | |
|---|---|
| From | 0xdd870fa1b7c4700f2bd7f44238821c26f7392148 |
| To | Sell.ConfirmToTransferBuyer() 0x692a70d2e424a56d2c6c27aa97d1a86395877b3a |
| Gas | 3000000 |
| Transaction cost | 30279 gas |
| Execution cost | 9007 gas |

*6) The deposit money in smart contract 2 is transferred to the seller:* see Table IX.

TABLE IX.      CASE 1: STEP 6

| | |
|---|---|
| From | 0x14723a09acff6d2a60dcdf7aa4aff308fddc160c |
| To | Cod.ConfirmSuccess()<br>0x692a70d2e424a56d2c6c27aa97d1a86395877b3a |
| Gas | 3000000 |
| Transaction cost | 29790 gas |
| Execution cost | 8518 gas |

*B. Case 2: There are transaction errors due to shipper problems*

*1) Seller 0xca35b7d915458ef540ade6068dfe2f44e8fa733c creates a package on smart contract 1:* see Table X.

TABLE X.      CASE 2: STEP 1

| | |
|---|---|
| From | 0xca35b7d915458ef540ade6068dfe2f44e8fa733c |
| To | Sell.AddPackage(string,uint256)<br>0x692a70d2e424a56d2c6c27aa97d1a86395877b3a |
| Gas | 3000000 |
| Transaction cost | 104693 gas |
| Execution cost | 82397 gas |

*2) Buyer 0xdd870fa1b7c4700f2bd7f44238821c26f7392148 agrees to deposit the amount of money equal to 10 Ether:* see Table XI.

TABLE XI.      CASE 2: STEP 2

| | |
|---|---|
| From | 0xdd870fa1b7c4700f2bd7f44238821c26f7392148 |
| To | Sell.ApplyBuy(uint256)<br>0x692a70d2e424a56d2c6c27aa97d1a86395877b3a |
| Gas | 3000000 |
| Transaction cost | 64531 gas |
| Execution cost | 43067 gas |

*3) Seller 0xca35b7d915458ef540ade6068dfe2f44e8fa733c gives the information of the package as well as the buyer's address to the shipper in the smart contract 2:* see Table XII.

TABLE XII.      CASE 2: STEP 3

| | |
|---|---|
| From | 0xca35b7d915458ef540ade6068dfe2f44e8fa733c |
| To | Cod.AddPackage(string,uint256,address,uint256)<br>0x0dcd2f752394c41875e259e00bb44fd505297caf |
| Gas | 3000000 |
| Transaction cost | 187265 gas |
| Execution cost | 163369 gas |

*4) Shipper 0x14723a09acff6d2a60dcdf7aa4aff308fddc160c agrees to deposit price and deliver the package:* see Table XIII.

TABLE XIII.      CASE 2: STEP 4

| | |
|---|---|
| From | 0x14723a09acff6d2a60dcdf7aa4aff308fddc160c |
| To | Cod.ApplyDeliver(uint256)<br>0x0dcd2f752394c41875e259e00bb44fd505297caf |
| Gas | 3000000 |
| Transaction cost | 44397 gas |
| Execution cost | 22933 gas |

*5) The shipper 0x14723a09acff6d2a60dcdf7aa4aff308fddc160c has a problem:* the buyer *0xdd870fa1b7c4700f2bd7f44238821c26f7392148* does not receive the package. The seller *0xca35b7d915458ef540ade6068dfe2f44e8fa733c* will receive money from the balance of smart contract 2, see Table XIV.

TABLE XIV.      CASE 2: STEP 5

| | |
|---|---|
| From | 0xca35b7d915458ef540ade6068dfe2f44e8fa733c |
| To | Cod.ErrorSuccess()<br>0x0dcd2f752394c41875e259e00bb44fd505297caf |
| Gas | 3000000 |
| Transaction cost | 30066 gas |
| Execution cost | 8794 gas |

*6) The money of balance in smart contract 1 will be transferred to the buyer 0xdd870fa1b7c4700f2bd7f44238821c26f7392148:* see Table XV.

TABLE XV.      CASE 2: STEP 6

| | |
|---|---|
| From | 0xdd870fa1b7c4700f2bd7f44238821c26f7392148 |
| To | Sell.ConfirmToTransferBuyer()<br>0x692a70d2e424a56d2c6c27aa97d1a86395877b3a |
| Gas | 3000000 |
| Transaction cost | 30279 gas |
| Execution cost | 9007 gas |

*C. Case 3: There are problems with transactions because the buyer does not accept commodities.*

*1) Seller 0xca35b7d915458ef540ade6068dfe2f44e8fa733c creates a package on smart contract 1:* see Table XVI.

TABLE XVI.      CASE 3: STEP 1

| | |
|---|---|
| From | 0xca35b7d915458ef540ade6068dfe2f44e8fa733c |
| To | Sell.AddPackage(string,uint256)<br>0x692a70d2e424a56d2c6c27aa97d1a86395877b3a |
| Gas | 3000000 |
| Transaction cost | 104693 gas |
| Execution cost | 82397 gas |

*2) Buyer 0xdd870fa1b7c4700f2bd7f44238821c26f7392148 agrees to deposit the amount of money equal to 10 Ether:* see Table XVII.

TABLE XVII.      CASE 1: STEP 2

| | |
|---|---|
| From | 0xdd870fa1b7c4700f2bd7f44238821c26f7392148 |
| To | Sell.ApplyBuy(uint256)<br>0x692a70d2e424a56d2c6c27aa97d1a86395877b3a |
| Gas | 3000000 |
| Transaction cost | 64531 gas |
| Execution cost | 43067 gas |

*3) Seller 0xca35b7d915458ef540ade6068dfe2f44e8fa733c gives the information of the package as well as the buyer's address to the shipper in the smart contract 2:* see Table XVIII.

TABLE XVIII.      CASE 3: STEP 3

| | |
|---|---|
| From | 0xca35b7d915458ef540ade6068dfe2f44e8fa733c |
| To | Cod.AddPackage(string,uint256,address,uint256)<br>0x0dcd2f752394c41875e259e00bb44fd505297caf |
| Gas | 3000000 |
| Transaction cost | 187265 gas |
| Execution cost | 163369 gas |

*4) Shipper 0x14723a09acff6d2a60dcdf7aa4aff308fddc160c agrees to deposit price and deliver the package:* see Table XIX.

TABLE XIX.    CASE 3: STEP 4

| From | 0x14723a09acff6d2a60dcdf7aa4aff308fddc160c |
|---|---|
| To | Cod.ApplyDeliver(uint256) 0x0dcd2f752394c41875e259e00bb44fd505297caf |
| Gas | 3000000 |
| Transaction cost | 44397 gas |
| Execution cost | 22933 gas |

*5) The buyer 0xdd870fa1b7c4700f2bd7f44238821c26f7392148 avoids responsibility and does not receive goods without any acceptable reasons:* The deposit money of the buyer will be transferred to the seller *0xca35b7d915458ef540ade6068dfe2f44e8fa733c* after the seller confirms that the problem is from the buyer, see Table XX.

TABLE XX.    CASE 3: STEP 5

| From | 0xca35b7d915458ef540ade6068dfe2f44e8fa733c |
|---|---|
| To | Sell.ErrorSuccess() 0xbbf289d846208c16edc8474705c748aff07732db |
| Gas | 3000000 |
| Transaction cost | 30044 gas |
| Execution cost | 8772 gas |

*6) The shipper takes back its money in smart contract 2:* see Table XXI.

TABLE XXI.    CASE 3: STEP 6

| From | 0x14723a09acff6d2a60dcdf7aa4aff308fddc160c |
|---|---|
| To | Cod.ConfirmToTransferShipper() 0x0dcd2f752394c41875e259e00bb44fd505297caf |
| Gas | 3000000 |
| Transaction cost | 30135 gas |
| Execution cost | 8863 gas |

## VI.    FINAL REMARKS

In Fig. 4, we observe that the total transaction gas and execution gas are similar in three investigation scenarios. In Ethereum, gas is the concept to discourage over-consumption of resources. The user who creates a transaction must purchase gas by spending currency. Every program instruction consumes some amount of gas during a transaction is executed. Consequently, if a transaction fails in case of the failure caused by other partners, the total gas spent should not be so high. The double smart contracts have successfully penalized undesired failures.

## VII.    CONCLUSION

In this paper, we propose a new framework that uses smart contracts, blockchain and Ethereum to assure non-fraudulent transactions in cash on delivery and enhance the reliability of distributed cryptocurrency platforms. Although there is a lot of thoughtful discussion on the use of smart contracts in distributed cryptocurrency, there have not been many frameworks that would address the non-fraudulent transactions caused by buyers and shippers. The penalties will be made in two circumstances: (i) the buyer refuses to receive the commodities without any reliable reasons and (ii) the shipper attempts to make any modification on the delivered goods during transportation. Our approach is a practical implementation which has been developed and evaluated empirically.
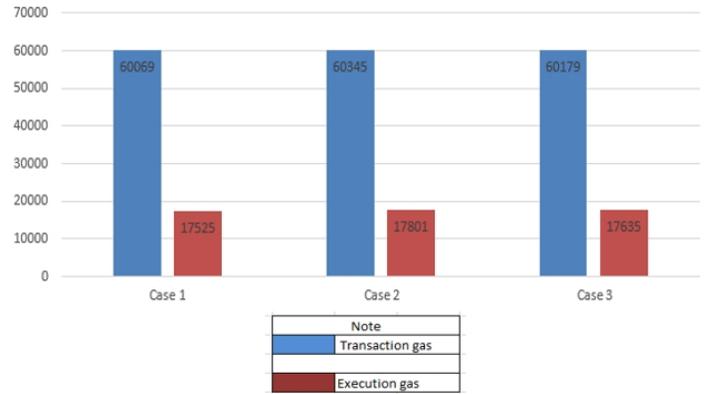


Fig. 4.    Gas consumption in three investigation scenarios.

To facilitate future research endeavors on COD and smart contract programming, we have released the open source codes of our implementation. The materials are available at https://github.com/nghiepnguyen520/Cash-on-delivery.

## REFERENCES

[1] P. Timmers, *Electronic commerce*. John Wiley & Sons, Inc., 1999.

[2] M. Halaweh, "Cash on delivery (cod) as an alternative payment method for e-commerce transactions: Analysis and implications," *International Journal of Sociotechnology and Knowledge Development (IJSKD)*, vol. 10, no. 4, pp. 1–12, 2018.

[3] ——, "Intention to adopt the cash on delivery (cod) payment model for e-commerce transactions: An empirical study," in *IFIP International Conference on Computer Information Systems and Industrial Management*. Springer, 2017, pp. 628–637.

[4] U. Tandon and R. Kiran, "Study on drivers of online shopping and significance of cash-on-delivery mode of payment on behavioural intention," *International Journal of Electronic Business*, vol. 14, no. 3, pp. 212–237, 2018.

[5] J. D. Alie and P. E. Vliek, "International cash-on-delivery system and method," Jul. 24 2007, uS Patent 7,249,069.

[6] A. Asgaonkar and B. Krishnamachari, "Solving the buyer and seller's dilemma: A dual-deposit escrow smart contract for provably cheat-proof delivery and payment for a digital good without a trusted mediator," *arXiv preprint arXiv:1806.08379*, 2018.

[7] "Two party contracts," Feb 2015. [Online]. Available: https://dappsforbeginners.wordpress.com/tutorials/two-party-contracts/

[8] A. M. Antonopoulos and G. Wood, *Mastering ethereum: building smart contracts and dapps*. O'Reilly Media, 2018.

[9] "Ethereum project." [Online]. Available: https://ethereum.org/

[10] H. X. Son, M. H. Nguyen, N. N. Phien, H. T. Le, Q. N. Nguyen, V. D. Dinh, P. T. Tru, and P. Nguyen, "Towards a mechanism for protecting seller's interest of cash on delivery by using smart contract in hyperledger," *International Journal of Advanced Computer Science and Applications*, vol. 10, no. 4, 2019. [Online]. Available: http://dx.doi.org/10.14569/IJACSA.2019.0100405

[11] "Hyperledger fabric." [Online]. Available: https://hyperledger-fabric.readthedocs.io/

[12] G. Hileman and M. Rauchs, "Global cryptocurrency benchmarking study," *Cambridge Centre for Alternative Finance*, vol. 33, 2017.

[13] H. R. Hasan and K. Salah, "Blockchain-based solution for proof of delivery of physical assets," in *International Conference on Blockchain*. Springer, 2018, pp. 139–152.

[14] A. Kosba, A. Miller, E. Shi, Z. Wen, and C. Papamanthou, "Hawk: The blockchain model of cryptography and privacy-preserving smart contracts," in *2016 IEEE symposium on security and privacy (SP)*. IEEE, 2016, pp. 839–858.

[15] S. Nakamoto *et al.*, "Bitcoin: A peer-to-peer electronic cash system," 2008.

[16] Y. Lewenberg, Y. Sompolinsky, and A. Zohar, "Inclusive block chain protocols," in *International Conference on Financial Cryptography and Data Security*. Springer, 2015, pp. 528–547.

[17] G. Wood *et al.*, "Ethereum: A secure decentralised generalised transaction ledger," *Ethereum project yellow paper*, vol. 151, pp. 1–32, 2014.

[18] T. Hamid, "Cash on delivery the biggest obstacle to e-commerce in uae and region," May 2014. [Online]. Available: https://www.thenational.ae/business/technology/cash-on-delivery-the-biggest-obstacle-to-e-commerce-in-uae-and-region-1.604383

[19] F. Idelberger, G. Governatori, R. Riveret, and G. Sartor, "Evaluation of logic-based smart contracts for blockchain systems," in *International Symposium on Rules and Rule Markup Languages for the Semantic Web*. Springer, 2016, pp. 167–183.

[20] M. Alharby and A. van Moorsel, "Blockchain-based smart contracts: A systematic mapping study," *arXiv preprint arXiv:1710.06372*, 2017.

[21] D. He, K. F. Habermeier, R. B. Leckow, V. Haksar, Y. Almeida, M. Kashima, N. Kyriakos-Saad, H. Oura, T. S. Sedik, N. Stetsenko *et al.*, "Virtual currencies and beyond: initial considerations," 2016.

[22] "Solidity." [Online]. Available: https://solidity.readthedocs.io/

[23] E. U. C. Team, "Ethereum unit converter." [Online]. Available: https://etherconverter.online/

[24] K. Delmolino, M. Arnett, A. Kosba, A. Miller, and E. Shi, "Step by step towards creating a safe smart contract: Lessons and insights from a cryptocurrency lab," in *International Conference on Financial Cryptography and Data Security*. Springer, 2016, pp. 79–94.

[25] "Solidity ide." [Online]. Available: https://remix.ethereum.org/

[26] Ethereum, "ethereum/remix-ide," Apr 2019. [Online]. Available: https://github.com/ethereum/remix-ide

# Use of Blockchain in Governance: A Systematic Literature Review

Asad Razzaq[1], Muhammad Murad Khan[*,2], Ramzan Talib[3], Arslan Dawood Butt[4],
Noman Hanif[5], Sultan Afzal[6], Muhammad Razeen Raouf[7]
Department of Computer Science[1,2,3,5,6,7]
Department of Electrical Engineering[4]
Government College University, Faisalabad, Pakistan

*Abstract*—**Blockchain is a distributed network based ledger that is secured by the methods of cryptographic proof. It enables the creation of self-executable digital contracts i.e. smart contracts. This technology is working in collaboration with major areas of research including governance, IoT, health, banking and education. It has anticipated revolutionary ways, which helps us to overcome the problems of governance such as human error, voting, privacy of data, security and food safety. In governance, there is a need to ameliorate the services and facilities with the assistance of blockchain technology. This paper aims to explore the issues of governance which can be resolved with the assistance of Blockchain features. Furthermore this paper also provides the future work directions.**

*Keywords*—*Blockchain; governance; voting; security; privacy*

## I. Introduction

Blockchain (BC) was introduced in 2008 by Satoshi Nakamoto. In order to solve the problem of double spending. Distributed ledger that digitally recorded the transactions however these transactions are encrypted and authenticated by the consensus protocol in the blockchain [1], [2]. In the network of blockchain blocks are created and maintained by using proof-of-work / proof-of-stake. Each block consists of previous block hash address that maintains the historical record in the network [3], [9], [11]. Blockchain was invented as a main technology of Bitcoin, however now its application are increased in a number of fields such as governance, IOT, education, industries and health etc [28]. Moreover blockchain is the technology that remove the centralized third party by distribution of power away from the central / third party in communication, business and even politics or law.

Governance consist of all the actions of governing undertaken by any state. Currently governance is facing many issues like the privacy of data, food safety and voting etc. These issues can be resolved with the assistance of blockchain features such as decentralization, smart contract and immutability [1], [2], [3], [4], [7], [9], [11], [16]. This paper summarize the features of blockchain to resolve the current governance issues.

The challenges of BCT in governance system (such as money laundering and irreversibility) and also determined a deep understanding of governance issues with their solutions [4], [7], [25], [15]. This paper aims to explore how blockchain technology (BCT) is utilized in governance system.

The rest of this paper is ordered as follows. Section 2 covers the literature and features of blockchain. Section 3, research methodology, research questions, and research approach

is examined. In Section 4 solutions of research questions are described. Section 5 discuses about blockchain framework in governance. Section 6 is provide detailed discussion of numerous challenges of validity. Finally in Section 7 conclusion and future directions are discussed.

## II. Review of Literature

This segment grants a number of fundamental thoughts and theories from prevailing research associated with the BCT and numerous functional executions in the governance. It also marks a assessment among the prevailing literature review.

### A. Blockchain

A blockchain is an entire innovative technology for construction and using ledger [11], [9], [16], [21]. Blockchain records data on decentralized node called blocks. Blockchain associations the temper resistant blocks consuming mathematical cryptography and open source software computers networks and incentive mechanism [11], [16], [3], [13]. The blockchain is a trustworthy technology as it does not demand third party verification whereas, BCT verifies the authenticity of record in the network with consensus protocol [11], [21]. Blockchain responds on the distributed system by abolishing centralized authority and execute on the base of decentralized consensus protocol [11], [9]. The blockchain is cryptographic arrival against tempering on transaction stored on blocks. It allows the facility to read the data on nodes [5], [11], [13].

The blockchain is separated into three types in respect to reading the Blockchain data.

*a) Public Blockchain:* Allows all nodes to access the data for reading and new data purpose. Bitcoin and lite coin are the best examples of it [11].

*b) Private Blockchain:* Allows those nodes that are pre-registered by central authorities means only has permission to authorize. Hyper ledger and ripple are the examples of it. Satoshi Nakamoto introduced the idea in 2008 and the data structure "chain of blocks" behind the crypto currency i.e., Bitcoin [11].

*c) Consortium Blockchain:* In this type, blockchain is not control by single authority, whereas controlled by group of approved authorities. Consortium blockchain is additionally titled as semi-decentralized [11].

*1) Features of blockchain:*

*a) Security:* BCT is purely a P2P network technology that eliminates the middle man influence/involvement [1], [2]. The transactions are validated using the consensus mechanism and authentic transaction laid down in a block that contains time-stamp and hashes of previous block [3], [1].Transaction are continuously synchronized with nodes on network and history of transaction remains visible. This P2P and consensus based nature of BCT provides the security to data [1], [2], [11].

*b) Decentralization:* BCT is a digital technology that arises in its distributed peer-to-peer nature that deviates the obligation from central command [1], [7]. It allows to avoid the gathering of power that could let a single party to take control on whole setup. It helps to get rid from intermediaries [9], [16].

*c) Smart contract:* A process of dissemination of digital assets among two or more parties automatically stated by the formula derived on the basis of data that is interrogated at the time of setting up the contract [1]. Smart contract is a electronic program that imposes its accomplishments on blockchain enrolled by the consensus protocol [1], [4], [7], [9]. Consensus mean if all nodes on network uphold the transaction the authentication will be done.
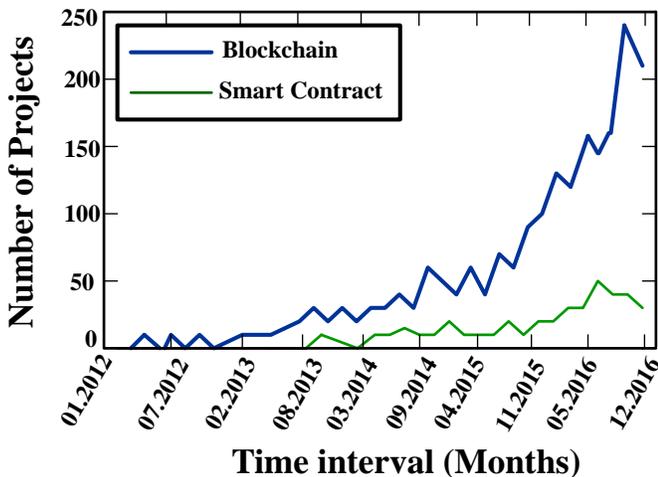


Fig. 1. Blockchain and Smart contracts increasing ratio interval of time [25].

Fig. 1 describes that with the passage of time, the blockchain projects (application) based upon smart contracts are increasing. The greatest benefit for financial engineering is "Smart contracts" computer programmable contracts that are checked and applied electronically, it creates a cheap previously agreed method to ensure that all parts of a contract are fulfilled without interference. The actual benefits of these technologies are their reliability and safety have yet to be assessed [4], [5].

*d) Traceability:* The traceability feature of BC provides the capability to figure out and track the origin of any transaction through digital signature [27].

*e) Immutability:* Blockchain provides a temper-proof environment in which once data has been added in block it can't be altered or removed. Immutability works totally based on 51 percent consensus involves proof of work [1].

Blockchain enables the capability to save records in very strong manners and provide the robust resistant environment to alter the data [4], [5], [6], [7], [9].

B. *Use of Blockchain in Governance*

Governance includes all of the procedures of governing whether undertaken by the government of a state, or by a social system (family, tribe, formal or informal organization) whether through the laws, norms, power or language of a managed society [1] It includes the mechanism needed to balance the member power (related to accountability) and the primary responsibility for enhancing the practicality of organization [20].

UK government released a report that blockchain is potential for government services. Keeping in view on this report government officially recommend the use of blockchain as a service [31].

Estonia government is using decentralized system for providing public notary facilities for their citizens. They can notarize birth certificates and educational documents etc. By using their electronic IDs provided by the state [31].

Blockchain technology can be used in voting system by government. This technology provide the anonymity of voter. The voting record always remain immutable. Danish political party used this in voting system [31].

Government apply blockchain technology for safety of food by using smart farmer market application. This application linked farmers with market, proper record about production and Transportation is grunted. This application ensure the safety and quality of food [31].

There is no hesitation in the significance of blockchain for the development and betterment of governance. Although, there always has been a struggle to find transformation technologies to support in the governance field. The area of governance is a prominent domain [19]. Blockchain governance is about determining who has the authority [12]. Blockchain technology is the amazing innovation which tries to help improve the governance field. BCT provides various benefits and used in various governance fields [13]. New technologies can reduce central administration and government costs to records, transactions, data security and system trust.

There is no cause to discard the value of governance for the progressive evolution of a country. It has been a contest to absorb assist in the governance field innovative technologies like blockchain. There has been number of creativities undertaken by the governance using block-chain technology to store the data. Dubai established blockchain based governance framework [30]. Blockchain provides benefits to Governance as payment using BCT feature smart contracts [5], [7], [17], [22]. The industry extends the concept of transactions to smart contracts. BCT features offers new ways of mutual interest and state management arranged different groups of societies in a more decentralized manner. Once two parties conduct agreement and publish it on Blockchain network that runs without human intervention [5], [7], [22]. The purpose of smart contracts is to make it easy by allowing human intervention to be taken out of the loop, thus allowing full automation [5], [9], [22].

## C. Compare Secondary Studies

Since the last decade blockchain has been the favourite area of the researchers. Governance is one of the main areas of this domain. During the research secondary study could not find any regarding "Governance".

## III. RESEARCH METHODOLOGY

A systematic literature review (SLR) means to identify, evaluate and interpret all available relevant research question [26]. To achieve this goal, The research guidelines gathered from Barbara Kitchenham which are discussed in the following subsections.

### A. The Aim of the Research

Aim of this research is to present an SLR in the field of governance and highlight the issues that are being solved when works with blockchain or its application.Moreover, the issues of governance remains unsolved after BC implementation are highlighted.

### B. Need of Conducting the Systematic Literature Review

The interest of this research is in the field of governance, but the review paper about area with blockchain was not found.

### C. Methodological Approach

The tag-based methodology is to explore interrelated papers. On Google Scholar, these keywords, "Blockchain governance", "Blockchain", "Governance" and "Blockchain governance review" were explored. All the research entitlement are classified from different journals (Springer, IEEE, ACM, WILLEY and some Conference papers) that are well-balanced and downloaded after the evaluation in order to shortlist (see Fig. 2).

### D. Inclusion and Exclusion Criteria

For inclusion, total, 52 papers were searched and shortlisted containing the required keywords "governance and blockchain". For this, the basic concept behind the primary studies was read and analysed. After analysing the abstract of papers certain studies were irrelevant to domain's concentration. Remaining research paper were really helpful and according to the required sectors. For inclusion, 32 primary studies were tagged for data extraction and others were excluded due to irrelevancy to research goal (see Fig. 2).

### E. Motivation and Research Questions

The basic purpose is to present the major features of blockchain which helps to solves the issues of governance. In this modern era of technology, a revolutionary change is necessary in the governance. Therefore, the suggestions that are mandatory for the improvement of current governing system. This paper points out some issues in the governance system and resolve these issues by deep analysis of blockchain. Research questions and motivation are described in Table I.
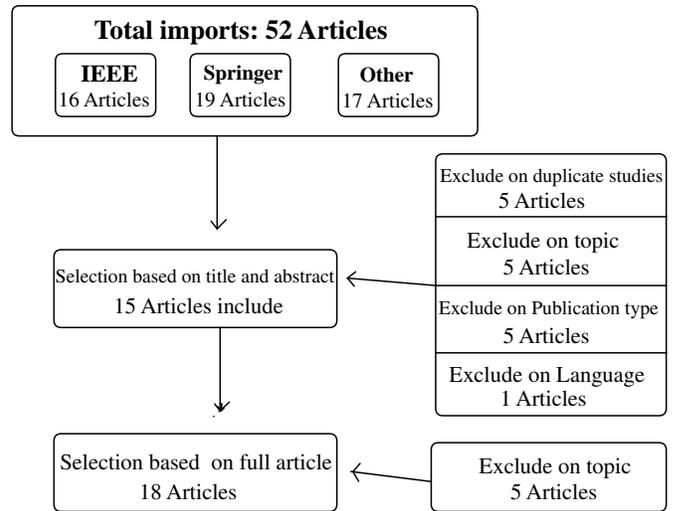


Fig. 2. Search and selection criteria.

TABLE I. RESEARCH QUESTIONS AND MOTIVATION

| # | Research Questions | Motivation |
|---|---|---|
| A | What are the issues regarding governance? | The objective is to grab the issues encountered by governance structure. |
| B | What is the significant feature practiced to decode the problem boundaries? | Ambition to explore the blockchain attribute that precisely resolve the problem faced by present governance structure. |
| C | What are unattended problems? | The purpose is to emphasize on unfold those issues that cloud be secured in future. |

### F. Classification Criteria

In this section, the classification of the shortlisted primary studies is presented. Shortlisted studies were classified according to research questions defined in Table II.

### G. Extraction of Research Data

After delimitation of the studies related to SLR, data extraction process was started. This segment focus on derivation on most relevant data from primary resources that encountered the issues of contemporary governance system and explore features of Blockchain for reformation. The following criteria are used to extract the data from these studies:

- Issues of the Governance
- Features of BC that resolve these issues
- Unaddressed issues

### H. Research Historical Background

Research in blockchain was started with the invention of Bitcoin in 2008. However, the very first paper of BCT was

TABLE II.        CLASSIFICATION AND SUMMARY OF COLLECTED PAPERS

| A. Classification | | |
|---|---|---|
| # | No. of primary studies | Classification criteria |
| A | 10 | Issues in Governance |
| B | 15 | Blockchain Features |
| C | 4 | Blockchain implementation issues |

| B. Summary of Collected Papers | | |
|---|---|---|
| *Issues* | Features | Unaddressed issues |
| [1][2][3][5] [7][10][11] [16][22][23] | [1][2][3][4][5][6] [8][9][11][12][14] [16][18][20][22] | [4][7][15][25] |

published in 2015. As it's not been so long so the publication year was not considered and all the previous work related to the topic was accumulated. Research papers found in year 2016, year 2017 and year 2018 (until March) were 8, 10 and 5, respectively.

## IV.    RESULTS

This section provides an explanation of the research questions that are discussed in three subcategory. The first subcategory classify the common problematic sectors of the governance structure. The second subcategory pinpoint the Blockchain features that settle down the issues of temporary Governance structure and final subcategory highlights the unattended sectors of problem area that could be secured in future.

### A.  What are the Issues Regarding the Governance?

*1) Privacy:* In this digital era, the headache of privacy of governmental records still exists due to central server storage [1], [2], [3], [11]. Data privacy is related to how information is being collected and handled.

*2) Incomplete contracts:* The incompleteness of a contract between authority and contractor due to an unforeseen component which may lead the government towards the problem of writing and enforcing contract cost. [11], [22].

*3) Voting:* It is the process of selecting aspirant in election. The traditional voting mechanism for casting the vote in designated polling station requires the handsome amount of time and money [16]. Existing voting structure enclosed with a serious flaw known as centralized design. In the existing voting system, the bogus vote can be cast that creates disaster or defragmentation in the state [3], [7].

*4) Human Error:* The human involvement in traditional government system plays a key role to manage the physical and digital records. However, the chance of mistake remains where a human being will manage the records manually [14], [29].

*5) Food Safety:* Food Safety is a problem where fresh and healthy agricultural products are not provided to the public [1]. Network becomes complex due to centralized parties like distributors, retailers, transporters and suppliers are stand between producers and end-users [10], [16].

### B.  What is the Significant Feature Practice to Decode the Problem Boundaries?

Table III Highlighting the one by one features of blockchain which are resolving the identified issues of governance.

*1) Security:* BCT is purely a P2P network technology that eliminates the middle man influence/involvement [1], [2]. Transaction are continuously synchronized with nodes on network and history of transaction remains visible. This P2P and consensus based nature of BCT provides the security to data [1], [2], [11]. This feature can overcome the following issues.

*a) Privacy:* Blockchain stores the transactions in network based distributed nodes so central server issue can be eliminated by shifting the data on the blockchain [1], [2]. Blockchain generates the consensus-based tempered proof transactions history that protect the sensitive data exchange on the network [3].

*b) Incomplete Contract:* Instead of trusted third party performing the transactions through their servers with permission, A p2p computer network running the blockchain protocol to checks the transaction by consensus [11], [22].

*2) Decentralized:* BCT is a digital technology that arises in its distributed peer-to-peer nature that deviates the obligation from central command [1], [7].

*a) Privacy:* As nature of blockchain the nodes stored the data in decentralized manner. The threat of data privacy can be minimized by shifting government institutional data on blockchain [16].

*b) Voting:* Due to the influence of third party in existing voting system bogus vote can be cast, decentralization feature can remove the third party so that it can reduce the issue of bogus vote [3], [7].

*c) Food Safety:* As decentralization remove the existence of intermediaries like distributors, retailers, transporters and suppliers. It may help the end-user to get the product directly from producer [10], [16].

*3) Traceability:* The feature of BC provides the capability to figure out and track the origin of any transaction through digital signature [27].

*a) Food Safety:* With the assistance of traceability feature agriculture products can be tracked digitally form ecosystem (How, when and where product was produced) to end-user [8], [9], [18].

TABLE III.    ISSUES VS FEATURES

| Issues | Blockchain features | | | | |
|---|---|---|---|---|---|
| | **Security** | **Decentralization** | **Immutability** | **Smart contracts** | **Traceability** |
| **Privacy** | [1][2][3][11] | [16] | | | |
| **Incomplete contracts** | [11][22] | | | [11] [22] | |
| **Voting** | | [3][7] | | | [16] |
| **Human Error** | | | [1][4][14][29] | | |
| **Food Safety** | | [10][16] | | [1][10][16] | [18] [8][9] |

*b) Voting:* Using traceability only the legal voters can enroll there self through digital signature and their identity can be tracked and it eliminate the threat of bogus vote [16].

*4) Immutability:* Blockchain provides a temper-proof environment in which once data has been added in block it can't be altered or removed [1], [4].

*a) Human Error:* Transaction is checked by all the nodes in consensus mechanism and this mechanism remove the probability of human error [1], [4], [14], [29].

*5) Smart contracts:* Smart contract is a electronic program that imposes its accomplishment on blockchain enrolled by the consensus protocol [1], [4], [7], [9]. Following issues can be resolved through smart contract.

*a) Incomplete Contract:* The Self- executable program helps the contract to not remain incomplete and it bounds the parties to complete the deal otherwise certain payment or transfer of digital asset not be made [11], [22].

*b) Food Safety:* The automation of the smart contract between the product provider and end-user eliminates the intermediary and helps in reducing cost to verify the product quality [1], [10], [16].

*C. What are Unattended Problems?*

*1) Money laundering:* Concealing the origin of illegally obtained money, usually through transfers, involving foreign banks or legal enterprises is known as money laundering as shown in the Fig. 3. Money laundering (ML) is a global problem. Bitcoin virtual currency provides a place for persons who generate, transfer, launder and steal illegal funds with some kind of anonymity [24]. Bitcoin offers many of the same tasks that are associated with other virtual currencies, such as Web Money, and added unique challenges to researchers because of its decentralized nature, even though there is no central Bitcoin server to compromise [25]. Possible Solution of laundering: Fig. 3 is providing two kinds of solution which has betterment or to overcome challenges that crypto coins presently pose for the global anti-money laundering (AML) [25].

*a) Educate Anti money laundering task force:* Firstly, in spite of some shortcomings, the risk assessment approach of the Financial Action Task Force (FATF) provides an effective balance between existing threats and opportunities that currently represent the Crypto-currency. The FATF is an officially intergovernmental organization. FATF Guidance: Three specific critics can be applied within the framework of the FATF approach, based on risk assessment. First, the FATF perpetuates greater dependency on global AML management in the market. Second according to the FATF management tends to rely on solutions based on the magic pool of Technology accordingly, Programming Application interface (APIs) provide client identification information "or third party digital identification systems that need to be adjusted themselves [25], [23].

*b) Totally banned of crypto-currencies:* The virtual nature of crypto-currencies (CCs), users evaluates the currency for many of the same reasons they believe that they can exchange currency-based goods, services or national currency at the date of the lateral [24]. As such, Bitcoins is currently accepted as a form of digital payment method by hundreds of legitimate retailers including vendors selling clothing, games, music, and some hotels and restaurants. If we totally banned crypto currency, money laundering will be stopped completely, but it is not possible because money gram and west union currency are used blockchain technology using virtual currency [25], [23].

*2) Irreversible:* Undoubtedly, the blockchain possesses remarkable properties like a distributed block, such as efficiency, irreversibility, and transparency [7]. However some disadvantages associated with this feature of blockchain [4]. In this mechanism, If data is entered incorrectly, then it can't erase/edit, because blockchain is peer-to-peer network have not reversible or erasable. So, it is a crucial issue of blockchain technology [15]. While the current settlement systems allow the process of voiding and reconciliation, imagine the problems of "rollback" in the world of chains [7].

## V.    CASE STUDY: DUBAI BLOCKCHAIN BASED GOVERNANCE FRAMEWORK

Dubai established Blockchain Global Council to enhance the position of UAE as a leading centre for innovation. Dubai takes the initiative in 2000 for innovation by using blockchain. This initiative involves numerous government departments with many online services being led by Dubai Smart

Fig. 3.    (Money laundering solution structure using blockchain)

Government (DSG). Dubai government had set up a Dubai Vision in 2013 and build a committee which was consist of many departments involved in the vision of Smart Dubai. An executive committee for Dubai Strategy was formed in 2014 consisting governments representatives and private investors.

Finally, in 2015 blockchain based framework and DSG was approved legally for Smart Dubai. *"BitOasis"* and *"Dubai Wills"* were two pilot projects initiated by Dubai Global Counsel to enable the clients inheriting their assets and to register their wills using BCT respectively [30].

## VI.    Threats to Validity

Blockchain is not hot trending area of research yet. A very small number of researches have been extracted for the support of this research till to date but this area is being favourite with the time span. Now a days many contemporary researchers are doing their explanation in concerned subject. Due to novelty of research work in BCT this research could not meet handsome number of reviews to support this research.

Due to lack of literature review it might be possible that any type of work regarding this area will be proposed by any other researcher between the closing time of this research until publication. That is being a mighty threat to this research.

## VII.    Conclusion

The main purpose of this research is to ameliorate the government services with the assistance of blockchain. This paper is intended to highlight the problems of governance such as security, privacy, food safety and incomplete contract. The characteristics of blockchain e.g. decentralization, transparency, and smart contracts are also discussed to overcome the issues of governance. Moreover the issues of governance which are not settled yet with the assistance of blockchain such as money laundering and irreversibility. The solution is also suggested for money laundering. It will be removed by applying the proper rules and regulations as well as blockchain technology implemented on governance.

There is a need of future work on quality standards, available tools for resource allotment and determining legal framework of blockchain implementation on governance.

## References

[1]    SØlnes, J Ubacht, M Janssen, "Blockchain in government: Benefits and implications of distributed ledger technology for information sharing," 2017, Elsevier.

[2]    Zyskind, Guy, and Oz Nathan. "Decentralizing privacy: Using blockchain to protect personal data." In Security and Privacy Workshops (SPW), 2015 IEEE, pp. 180-184. IEEE, 2015.

[3]    Noizat, Pierre. "Blockchain electronic vote." In Handbook of digital currency, pp 453-461, 2015, Elsevier.

[4]    Atzori, Marcella. "Blockchain technology and decentralized governance: Is the state still necessary?" 2015.

[5]    Shermin, Voshmgir. "Disrupting governance with blockchains and smart contracts." Strategic Change 26, no. 5 (2017): 499-509.

[6]    Hou, Heng. "The application of blockchain technology in E-government in China." In Computer Communication and Networks (ICCCN), 2017 26th International Conference on, pp. 1-4. IEEE, 2017.

[7]    Ojo, Adegboyega, and Samuel Adebayo. "Blockchain as a Next Generation Government Information Infrastructure: A Review of Initiatives in D5 Countries." In Government 3.0–Next Generation Government Technology Infrastructure and Services, pp. 283-298. Springer, Cham, 2017.

[8]    Hoggett, R. D. "People, demand, and governance in future energy systems." (2017).

[9]    Qi, Renming, Chen Feng, Zheng Liu, and Nezih Mrad. "BlockchainPowered Internet of Things, E-Governance and E-Democracy." In EDemocracy for Smart Cities, pp. 509-520. Springer, Singapore, 2017.

[10]    Yermack, David. "Corporate governance and blockchains." Review of Finance 21, no. 1 (2017): 7-31.

[11]    Davidson, Sinclair, Primavera De Filippi, and Jason Potts. "Disrupting governance: The new institutional economics of distributed ledger technology," 2016.

[12]    Wang, Sha, Jean-Philippe JP Vergne, and Ying-Ying Hsieh. "The internal and external governance of blockchain-based organizations: Evidence from cryptocurrencies." In Bitcoin and Beyond, pp. 48-68. Routledge, 2017.

[13]    Millard, Jeremy. "European Strategies for e-Governance to 2020 and Beyond." In Government 3.0–Next Generation Government Technology Infrastructure and Services, pp. 1-25. Springer, Cham, 2017.

[14]    Chohan, Usman. "The Decentralized Autonomous Organization and Governance Issues." (2017). Carter, Nic. "A Cross-Sectional Overview of Cryptoasset Governance and Implications for Investors." (2017)

[15]    Wu, Yifan. "An E-voting System based on Blockchain and Ring Signature." Master. University of Birmingham (2017).

[16]    Böhme, Rainer, Nicolas Christin, Benjamin Edelman, and Tyler Moore. "Bitcoin: Economics, technology, and governance." Journal of Economic Perspectives 29, no. 2 (2015): 213-38.

[17]    Tian, Feng. "An agri-food supply chain traceability system for China based on RFID blockchain technology." In Service Systems and Service Management (ICSSSM), 2016 13th International Conference on, pp. 1-6. IEEE, 2016.

[18] Kim, Henry M., Marek Laskowski, and Ning Nan. "A First Step in the Co-Evolution of Blockchain and Ontologies: Towards Engineering an Ontology of Governance at the Blockchain Protocol Level." arXiv preprint arXiv:1801.02027 (2018).

[19] Singh, Munindar P., and Amit K. Chopra. "Violable Contracts and Governance for Blockchain Applications." arXiv preprint arXiv:1801.02672 (2018).

[20] Yli-Huumo, Jesse, Deokyoon Ko, Sujin Choi, Sooyong Park, and Kari Smolander. "Where is current research on blockchain technology?—a a systematic review." PloS one 11, no. 10 (2016): e0163477.

[21] Cong, Lin William, and Zhiguo He. Blockchain disruption and smart contracts. No. w24399. National Bureau of Economic Research, 2018.

[22] Gietzmann, Miles B. "Incomplete Contracts, and the make or buy decision: governance design and attainable flexibility." Accounting, Organizations and Society 21, no. 6 (1996): 611-626.

[23] Moser, Malte, Rainer Bohme, and Dominic Breuker. "An inquiry into money laundering tools in the Bitcoin ecosystem." In eCrime Researchers Summit (eCRS), 2013, pp. 1-14. IEEE, 2013.

[24] Campbell-Verduyn, Malcolm. "Bitcoin, crypto-coins, and global anti-money laundering governance." Crime, Law and Social Change 69, no. 2 (2018): 283-305.

[25] Bartoletti, Massimo, and Livio Pompianu. "An empirical analysis of smart contracts: platforms, applications, and design patterns." In International Conference on Financial Cryptography and Data Security, pp. 494-509. Springer, Cham, 2017.

[26] Khan, Muhammad Murad, Roliana Ibrahim, and Imran Ghani. "Cross domain recommender systems: a systematic literature review." ACM Computing Surveys (CSUR) 50.3 (2017): 36.

[27] Cartier, Laurent E., Saleem H. Ali, and Michael S. Krzemnicki. "Blockchain, Chain of Custody and Trace Elements: An Overview of Tracking and Traceability Opportunities in the Gem Industry." Journal of Gemmology 36, no. 3 (2018).

[28] Sultan A "Internet of Things Security Issues and Their Soluti ons With Blockchain Technology Characteristi cs: A Systematic Literature Review." Am J Compt Sci Inform Technol Vol.6 No.3:27 (2018)

[29] Lomas, Elizabeth. "Information governance: information security and access within a UK context." Records Management Journal 20, no. 2 (2010): 182-198.

[30] Karmakar, Ashmita, and Ummer Sahib. "Smart Dubai: Accelerating innovation and leapfrogging E-democracy." In E-Democracy for Smart Cities, pp. 197-257. Springer, Singapore, 2017.

[31] Alketbi, Ahmed, Qassim Nasir, and Manar Abu Talib. "Blockchain for government services—Use cases, security benefits and challenges." In 2018 15th Learning and Technology Conference (LT), pp. 112-119. IEEE, 2018.

# Swarm Robotics and Rapidly Exploring Random Graph Algorithms Applied to Environment Exploration and Path Planning

Cindy Calderón-Arce[1], Rebeca Solis-Ortega[2]
School of Mathematics,
Costa Rica Institute of Technology,
Cartago, Costa Rica
https://orcid.org/0000-0002-0077-225X[1]
https://orcid.org/0000-0002-3065-8386[2]

*Abstract*—**We propose an efficient scheme based on a swarm robotics approach for exploring unknown environments. The initial goal is to trace a map which is later used to find optimal paths. The algorithm minimizes distance and danger. The proposed scheme consists in three phases: exploration, mapping and path optimization. A cellular automata approach is used for the simulation of the fist two phases. For the exploration phase, a stigmergy approach is applied in order to allow for swarm communication in a implicit way. For the path planning phase a hybrid method is proposed. First an adapted Rapidly-exploring Random Graph algorithm is used and then a scalarized multiobjective technique is applied to find the shortest path.**

*Keywords*—*Swarm robotics; cellular automata; path planning; Rapidly-exploring Random Graph (RRG); scalarized multiobjective optimization*

## I. Introduction

Swarm algorithms have been deeply studied to address problems such as food search and object collection. However, it has been recently used to the exploration and mapping of scenarios, given the importance that this entails in rescue situations. In particular, use swarms robotics for the examination of scenarios in the search for optimal, efficient and safe routes, in a prudential time, reduces the loss and waste of resources in trajectories inspection.

Swarm robotics study the coordination of a large group of relatively simple robots through local rules and implicit communication. This field emerges from the application of swarm intelligence to robots [1]. Swarm intelligence is inspired on insect colonies, bird flocks, fish schools and other types of animal clusters which accomplish complex tasks through simple rules and communication.

It has to be clear that a group of agents must meet specific requirements in order to be considered as a swarm. The agents must be: autonomous, homogeneous, able to sense and actuate in the environment [2]. The aforementioned set of characteristics ensure a distributed and scalable swarm.

There are several swarm applications that have been studied: aggregation, flocking, exploration, foraging, navigation, path formation, object assembly and others [1], [3], [4]. In many of this applications, the use of a swarm is desired given the dangerous nature of the task.

Even though the use of a swarm brings several advantages, the type of communication associated with it can create a possible drawback. There are two possible control schemes: centralized and decentralized. Neither of these control schemes facilitate the supervision of the swarm by a human operator [1].

When a centralized swarm is used its scalability is poor and the swarm becomes sensible to the loss of its central leader. The decentralized approach overcome this main issue, but does not allow to synthesize or access global data unless all individuals are connected to each other. Therefore a human controller cannot access the data thus it cannot predict or alter the behaviour of the swarm.

In addition, path planning methods that minimize not only distances, but also danger, cost, time or energy are of great relevance searching evacuation and access routes in buildings, urban centers or even forests.

In this paper we focus our work on the exploration of unknown, static and dangerous environments. Although these tasks can be executed by humans, the use of swarm robotics will allow to save resources and protect the people in charge of those tasks.

To solve the problem of creating an optimized pathway in an unknown environment, three phases are proposed. In the first stage, a simulated swarm based on a cellular automata scheme will be used for exploring unknown environments. This scheme autonomously propitiate an efficient dispersion of the swarm through the unknown area. In order to perform the task in a more efficient and scalable manner, a decentralized swarm is implemented.

The second part of the solution also uses a cellular automata approach, in which all the visited cells are recorded by an external server in order to trace the map. The non-visited cells are marked as obstacles. Finally, in the third stage, a discrete adapted RRG structures the zone by means of a graph. Then, a Dijkstra algorithm finds the shortest path between two given points.

This paper is structured as follows: in the next section works related with swarm robotics, exploration and path planning are shown. Methods and Materials section presents problem statement and the algorithms used to solve each problem

stage. Section IV shows environments and experimental set up in the simulations were carried out and their final results. Finally, in the last one general conclusions of the work will be discussed.

## II. Related Work

One aim of environments exploration is to know a throughout region, with the purpose of detect some targets distributed randomly in the area [5]. For accomplish this task, the most applied algorithms are bio-inspired. Among the approaches that have been developed, the stigmergy used by many colonies of insects to coordinate their activities, which allows a swarm indirect communication, has been widely researched [5], [6].

Palmieri et al. implemented and tested three biologically inspired coordination strategies: firefly, particle swarm and artificial bee behaviour [5], the better performance was obtained with the firefly-based strategy. This scheme focus on finding targets and recruit robots around them.

In addition, Tan et al. employed a stigmergy method for target search in unknown environments, by means of a swarm [6]. The stigmergy mechanism was employed to guide the robots motion. A pheromone map was used for helping them to reach the target. The most of exploration studies focus only on exploration for finding a target [7], [8]. Therefore, they do not guarantee the recognition of whole explored area.

On the other hand, the standard methods for solving path planning problems are based on using approximated schemes known as sampling based planning methods. These methods employ a random sample of the space (connecting points randomly) and deciding if a route or direction is feasible or if exists a possible collision [9]. Probabilistic RoadMaps (PRM) [10], Rapidly-exploring Random Graph (RRG) [11] and Rapidly-exploring Random Tree (RRT) [12] algorithms are example applications of these methods. RRG generates an undirected graph, possibly containing cycles, and RRT a directed tree. Similarly, PRM algorithms trace a roadmap (graph) which represents a set of collision-free trajectories for computing the shortest path that connects an initial node to a final one [13].

All of these algorithms differ on the process applied to construct a connecting graph [14]. They are probabilistically complete algorithms which have natural support for solving high dimensional complex problems [15]. However, they have the disadvantage of no capacity to stop execution upon failure nor the ability to report when no possible solution exists. Therefore, they are computationally expensive [16]. Another known graph search algorithm is called A*, which uses a discrete space and its success is highly dependent on grid resolution [17].

Dijkstra is one of most famous and simple optimization algorithm [18], with a quadratic time complexity [19]. Also bio-inspired optimization algorithms, as a metaheuristic methods, are often used to approximate an optimization problem solution. They obtain solutions on an efficient way but are not able to meet the real time constraints, neither to reproduce the same solution since they are stochastic [16].

Furthermore, for finding a path, it is possible to optimize not only distance but also other objectives like danger, time or energy needed to cross it. In that sense, it is possible to optimize a multiobjective problem, taking into account several objectives, instead of a problem with a unique objective [20]. Without loss of generality, we consider a biobjective problem, which optimizes distance and dangerousness, but it is also possible to extend the proposed solution to a problem with more than two objectives. A technique implemented to solve multiobjective problems combines the objectives by means of a linear combination of them. Thus, the multiobjetive problem becomes uniobjective and the general optimization methods could be applied to solve the scalarized problem. In which, the objective function incorporates performance indicators of different objectives [21].

## III. Materials and Methods

A simulated swarm of robots and multiobjectives techniques are used as part of the proposed solution, which is organized in three phases: exploration, mapping and path planning, all carried out on simulated environments.

### A. Exploration Phase

A cellular automata approach was used for representing the environment, obstacles and swarm. Von Neumann neighborhood, which is composed by a central cell and its four orthogonally adjacent cells, was used. The states associated to each cell were defined by integer numbers as follows: (0) free cell, (1) cell occupied by an agent and (2) obstacle. The first two states are changeable during the time, but the third one is static.

We implemented two approaches to control the behavior of the swarm: a classic scheme based on a random walk algorithm, and a bio-inspired one based on stigmergy concept. For both, a modification was made which constraints the agent direction.

The stopping criteria of the algorithms was based on the percentage of environment coverage. It was selected just to valid and compare the schemes but for real world applications others criteria can be used like number of iterations or elapsed time.

*1) Random walk algorithm (RW):* This algorithm is used for search strategies for both animal and robots. It is especially useful when the individuals do not know the environmental and do not have cues that can drive the motion, or when their cognitive abilities do not support complex localisation and mapping behaviours [22].

In its simplest form, a random walk can be thought of as a sequence of straight movements and direction changes [22]. Given this an agent can be in one of two states: moving randomly or changing direction for avoiding obstacles [23].

If we take this scheme and adapt it into a cellular automata environment, RW can be modeled assuming that each agent is located in a cell and randomly chooses another one free in its neighborhood.

A modification of RW, called random walk with direction (RWD), adds a priority direction to each agent. This priority will cause that an agent will never be able of choosing a free cell that is opposite to its priority direction. In that sense, an
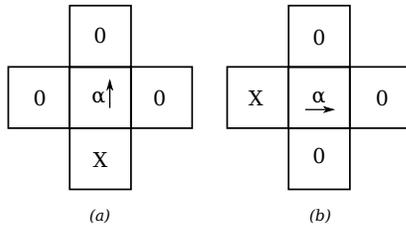
Fig. 1. Representation of the three directional neighbors of the cell $\alpha$ where the agent is located (X represents the non directional neighbor). (a) An agent with direction North can occupy any free cell situated West-North-East from its current position. (b) An agent with direction East can occupy any free cell situated North-East-South from its current position

agent located in $\alpha$ cell, will only have three possible cells to move on (directional neighbors). In Fig. 1 on the example (a) the priority direction is set north, in that case the agent could choose between three cell options: east, north or west. If obstacle collision occurs, the priority will randomly change before to continuing with the algorithm, as shown in Alg. 1.

---

**Algorithm 1:** Exploration algorithm by random walk with direction (RWD)

**Input:** $\rho$
1   $\alpha \leftarrow \rho$ position cell
2   $\tau \leftarrow \rho$ direction
3   **while** *the stop criteria are not satisfied* **do**
4     **foreach** $\rho$ **do**
5       $w^t(\rho_i) \leftarrow$ state of the $i$-directional neighbor at time $t$
6       **if** $\exists\, w^t(\rho_i) = 0$ **then**
7         $j$: randomly choose one directional neighbor
8         $\alpha \leftarrow \rho_j$
9       **else**
10         obstacle detected
11         $\tau \leftarrow$ choose new direction
12       **end**
13     **end**
14     $t++$
15   **end**

---

Even though random walk performs very well on a swarm [23], it does not allow agents to learn about the swarm decisions or to communicate with the environment (beyond the detection of obstacles). In order to take thoughtful decisions for accomplish a faster dispersion through the zone, an alternative algorithm based on a stigmergy approach is proposed.

*2) Virtual pheromone algorithms:* The stigmergy approach, proposed by Grassé and mentioned by Tan et al. in [6] is a mechanism of indirect communication. The agents leave a trace in the environment to propitiate different actions in others agents for leading to a spontaneous and systematic behavior emergence.

Based on this concept, we developed an algorithm called pheromone walk (PW). In this the agents ($\rho$) move through the environment sensing and leaving virtual pheromone on visited cells, with the aims of communicate the space already explored.

The pheromone has two parameters: intensity and evaporation rate. The intensity ($\Upsilon$) determines how strong the initial signal is in each cell $\alpha$ and the evaporation rate indicates how fast this signal will fade away. This pheromone is modeled by an iterative relation, that depends on time $t$ and a decay rate $\kappa$, as follows:

$$\Upsilon^{t+1}(\alpha) = (1 - \kappa) \cdot \Upsilon^t(\alpha) \tag{1}$$

where $\Upsilon^0$ is the initial pheromone intensity and $\kappa$ is a constant selected properly.

Under this approach, the agents choose a free cell with the lowest pheromone intensity, in the case that two or more cells have the same lowest intensity, then one will be chosen randomly.

Similar to the random walks algorithms, a modification of PW called directional pheromone walk (PWD) was implemented. In this case, also a priority direction is assigned to agents, which will select a free cell with the lowest pheromone intensity and located in one of its three directional neighbors. If an obstacle collision is detected, the agent will have to change its direction and move on (see Alg. 2).

---

**Algorithm 2:** Exploration algorithm by directional pheromone walk (PWD)

**Input:** $\rho$
1   $\alpha \leftarrow \rho$ position cell
2   $\tau \leftarrow \rho$ direction
3   **while** *the stop criteria are not satisfied* **do**
4     **foreach** $\rho$ **do**
5       $w^t(\rho_i) \leftarrow$ state of the $i$-directional neighbor of $\rho$ at time $t$
6       **if** $\exists\, w^t(\rho_i) = 0$ **then**
7         $j$: randomly choose one directional neighbor with $min(\Upsilon^t(\rho_i))$
8         $\alpha \leftarrow \rho_j$
9       **else**
10         obstacle detected
11         $\tau \leftarrow$ choose new direction
12       **end**
13     **end**
14     **foreach** $\alpha$ **do**
15       $\Upsilon^{t+1}(\alpha) \leftarrow (1 - \kappa) \cdot \Upsilon^t(\alpha)$
16     **end**
17     $t++$
18   **end**

---

### B. Path Planning Phase

Once the obstacle positions are known and the map is already constructed, an adapted RRG algorithm is used to structure the space. After that, a Dijkstra algorithm is applied to find a shortest path between an initial and goal given points.

*1) Adapted RRG algorithm:* RRG original algorithm operates in a space, in which a configuration is represented by any point on the work space, including obstacles [9], [15]. In this work, each agent is a point into a two dimensional space.

Since the amount of possible configurations is uncountable, the work space is discretized through a rectangular uniform

partition $\mathcal{M}$ of $segh \times segv$ dimension. Thus, adapted RRG algorithm generates a random point on each cell or partition's element.

The adapted RRG algorithm structures the search in a partition of the configuration space, using a discretization and forcing the graph to explore the whole space by including a vertex from each cell of the partition, if it is possible. Let $\mathcal{C}$ be the space configuration, $\mathcal{C}_{free} \subset \mathcal{C}$ the set of collision-free configurations and $\mathcal{G}(V, E)$ the graph defined by the vertices set $V$ and the edges set $E$. The adapted RRG algorithm initializes the graph with $q_{init}$ as a unique vertex without edges. A new configuration is generated creating a random point in $\mathcal{C}$ and looking for the nearest vertex $q_{near} \in \mathcal{G}$, by means of a Breadth-First Search (BFS) and a First In First Out (FIFO) buffer [9], considering possible collisions. The new point and its edge with $q_{near}$ are added to the graph and so on, until each cell of the partition has a vertex into the graph. When a collision between the graph and any obstacle is detected, the algorithm looks for the nearest collision-free configuration in the same direction as $q_{new}$ [24].

A graph $\mathcal{G}$ has a collision if $\mathcal{G}$ goes through a configuration $q \in \mathcal{C} \setminus \mathcal{C}_{free}$. Let $\mathcal{W}$ be the set of obstacle borders in the work space. Now suppose that $\mathcal{G}$ does not have a collision on the $ith$ iteration, but it is obtained in the next one by adding the vertex $q_2$ and the edge $\overline{q_1 q_2}$. Then $q_1 \in \mathcal{C}_{free}$ but there exists at least a point in $\overline{q_1 q_2}$ that is not in $\mathcal{C}_{free}$, i.e, is in $\mathcal{W}$.

Consider the function $f$ defined as follows:

$$f(x,y) = (x - x_1)(y_2 - y_1) - (y - y_1)(x_2 - x_1) \qquad (2)$$

where $(x_1, y_1)$ and $(x_2, y_2)$ are points on the Cartesian plane. The equation $f(x,y) = 0$ defines the locus of the points from the line $l$ passing through $(x_1, y_1)$ and $(x_2, y_2)$, which divides the plane into two regions. Since $f$ is continuous on $\mathbb{R}^2$, if a point $(x,y)$ over one side of $l$ produces $f(x,y) > 0$, then $f(x,y) > 0$ for all points in that region and $f(x,y) < 0$ for all points in the opposite region. By means of that property the intersection between two segments $\overline{AB}$ and $\overline{CD}$ occurs if and only if $f(A) \cdot f(B) < 0$ with respect to $\overline{CD}$ and $f(C) \cdot f(D) < 0$ with respect to $\overline{AB}$.

Identifying when the graph intersects $\mathcal{W}$, it is possible to define when accept or discard $q_2$. In that sense, a COLLISION() function determines if there is an intersection between $\overline{q_1 q_2}$ and $\mathcal{W}$ and helps to determine if a new point $q_2$ is inside or outside an obstacle, based on the previous function (2) and the Jordan curve theorem [24], [25].

Alg. 3 shows how a network is generated where, $D(m)$ is the set of all possible directions to go from the reference point on cell $m$ to an adjacent cell, and $F(m,d)$ indicates the selected cell after applying $d \in D(m)$ from $m$. CELL function provides the index of the cell corresponding to a configuration $q$. All points generated on each cell are contained on $PointList$. On each iteration a center cell is picked from a list called $Queue$ and a point is generated on every adjacent cell following the BFS method. $Queue$ originally contains the index corresponding to $q_{init}$ and adds indices of those adjacent cells from each iteration. Central indices that were already searched are deleted from $Queue$ and added to $Memory$. The GRID_POINT() function gives a point into the cell of index $F(m,d)$, if it is possible. If not, the answer will be $\emptyset$.

The indices are arranged into a sequential order beginning on CELL($q_{init}$) and continues bordering the adjacent cells. The algorithm ends when all the cells have been "visited" by the graph.

---

**Algorithm 3:** Space structuration

**Input:** $\mathcal{W}$, $q_{init}$, $q_{goal}$, $segh$, $segv$
**Output:** $\mathcal{G}(V, E)$

1   $V \leftarrow q_{init}$
2   $E \leftarrow \emptyset$
3   $\mathcal{M} \leftarrow$ Uniform Partition $segh \times segv$
4   $PointList \leftarrow [\ ]$
5   $i \leftarrow$ CELL($q_{init}$)
6   $g \leftarrow$ CELL($q_{goal}$)
7   $PointList(i) \leftarrow q_{init}$
8   $PointList(g) \leftarrow q_{goal}$
9   $Queue \leftarrow \{i\}$
10   $Memory \leftarrow \emptyset$
11 **while** $Queue \neq \emptyset$ **do**
12    $m \leftarrow Queue(1)$
13    $Queue \leftarrow Queue \setminus \{m\}$
14    **if** $m \notin Memory$ **then**
15     $q \leftarrow PointList(m)$
16     **for** $d$ **in** $D(m)$ **do**
17      $q_{new} \leftarrow$ GRID_POINT($\mathcal{W}, q_{init}, ...$
      $...PointList, F(m,d)$)
18      $PointList(F(m,d)) \leftarrow q_{new}$
19      **if** $q_{new} \neq \emptyset$ **and** $q \neq \emptyset$ **then**
20       **if not** COLLISION($\mathcal{W}, q, q_{new}$) **then**
21        $V \leftarrow V \cup \{q_{new}\}$
22        $E \leftarrow E \cup Edge(q, q_{new})$
23       **end**
24      **end**
25      $Queue \leftarrow Queue \cup \{F(m,d)\}$
26     **end**
27     $Memory \leftarrow Memory \cup \{m\}$
28    **end**
29 **end**

---

Space structuration algorithm (Alg. 3) shows an unidirectional search method, which generates just one graph from the initial configuration to structure the complete work space. But also a sequential and parallel bidirectional search was implemented, on which two sub-graphs are generated; one from $q_{init}$ and the other one from $q_{goal}$, and at the end both sub-graphs are joined to structure all the work space. The sequential algorithm alternates each sub-graph generation looking for a balance in the amount of vertices on each one and parallel one generates each sub-graph simultaneously.

*2) Optimization scalarized problem:* Once the graph has been generated, a Dijkstra Shortest Path (DSP) algorithm is used to find the optimal path from $q_{init}$ to $q_{goal}$ [18]. Taking into account not only distance but also dangerousness, by means of the following multiobjective optimization scalarized problem

$$\min_{p \in \mathcal{P}} \left( \omega D_1(p) + (1 - \omega) D_2(p) \right) \qquad (3)$$

where $D_1(p)$ and $D_2(p)$ are the distance and dangerousness values of the path $p$, respectively, $\mathcal{P}$ is the set of all possible paths between $q_{init}$ and $q_{goal}$ through the graph and $\omega \in [0,1]$ is the scalarizing constant.

A filter is applied before using the DSP algorithm, which simplifies the graph removing all the vertices and its corresponding edges with dead ends. Then a cost matrix $CostMat$ is calculated which assigns each element with the objective function value $(\omega D_1(e) + (1-\omega)D_2(e))$, for all edges $e \in E$. If two vertices of $V$ are not directly connected its cost value will be infinity. $D_1(e)$ assigns length and $D_2(e)$ dangerousness of $e$. Since, it is not possible to know a priori a continuous danger function defined on all the environment but it is possible to know information about dangerousness in some places of environment, we simulate danger information as a discretized heat distribution on the explored environment. Then, that information is used to construct a continuous dangerousness function into the work space. A Radial Basis Function (RBF) is selected to construct that function, because the most likely the information has lots of different features and nonlinear behaviour and it is known that RBF has a good performance ajusting nonlinear data. Although there are different kernels to used with RBF, for simplicity without loss of generality we consider a polyharmonic spline (PHS) kernel. Let $\psi : \mathcal{C} \to \mathbb{R}$ be the resulting heat RBF, which is defined as follows

$$\psi(x,y) \;=\; \sum_i \left( \eta_i \cdot \phi_i(r) \right) \qquad (4)$$

with $r = \|(x,y) - (x_c,y_c)\|_2$, $(x_c,y_c)$ represents where the RBF is centered, $\phi_i$ each PHS function, $\eta_i$ its corresponding weights and $\phi(r)$ is defined as

$$\phi(r) \;=\; \begin{cases} r^k \ln(r) & \text{; for an even } k \\ r^k & \text{; otherwise} \end{cases} \qquad (5)$$

where $k$ represents the PHS order [26], [27]. Thus, $\psi$ is used to assign dangerousness rate on each edge $e$ of graph $\mathcal{G}$, through a line integral over $\psi$ as shown in (6), which is approximated by means of numerical methods.

$$D_2(e) \;=\; \left| \int_\psi e \, ds \right| \qquad (6)$$

Finally, the DSP algorithm is executed to obtain the optimal path, using *CostMat* defined above.

## IV. Experiments

The work space was portrayed as a two-dimensional map composed by different kind of obstacles, randomly distributed.

### A. Experimental Setup

Three environments have been created in order to analyze the execution of the system. We take under consideration different scenarios like closed and open spaces, convex and concave obstacles, dead ends and others.
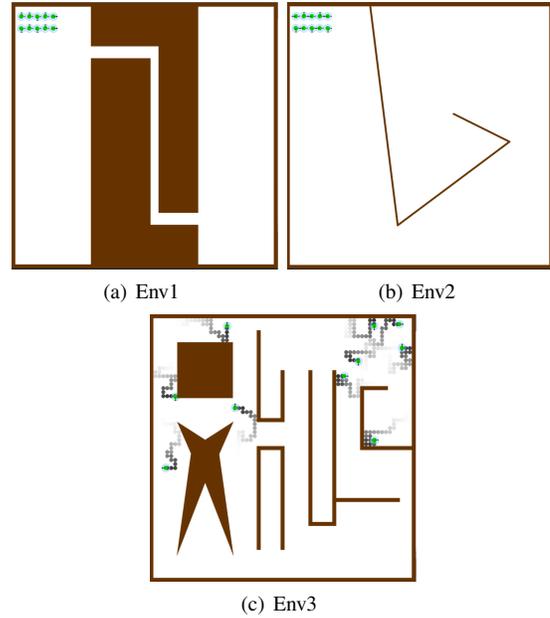


(a) Env1  (b) Env2

(c) Env3

Fig. 2. Representation of environments used for analysing the execution of the proposed solution.

The simulations were carried out using the software Processing 3.3.7, with a $1005 \times 1005$ pixel map for all the environments. Brown shapes represent the obstacles, while the agents of the swarm (which are a 1:67 scale of the map) are shown as two concentric circles that portray the body and the range of sensing. Also, the pheromones was represented by a circle in the corresponding cell (Fig. 2).

For the exploration phase we employ three different sized swarm of 10, 15 and 20 agents. The swarm, in each case, was deployed from the same area. For the PWD and PW algorithms the parameters selected for the evaporation rate was $\kappa = 1/\Upsilon^t(\alpha)$. Also five different initial pheromone intensities ($\Upsilon^0$ = 100, 300, 500, 700, 900) where defined.

All the agents are autonomous, homogeneous and have the ability to communicate with the environment through pheromones. Each robot not only walks around the environment avoiding obstacles, but also leaves a trace behind it of pheromone and senses the virtual substance of others and determines the intensity of it.

For the optimization phase, work space partitions of $5 \times 5$, $10 \times 10$, $15 \times 15$ and $20 \times 20$ were used in order to generate the possibles optimal paths. Three order PHS functions $\phi_i$ are using to compute $\psi$, by means of a RBF with the training data as centres and the heat distribution shown in Fig. 3.

The heat distribution is the same for all environments and *CostMat* is calculated with five different values for the scalarizing constant $\omega$: 0, 0.125, 0.25, 0.375 and 0.5.

Problem shown in (3) is solved 100 times on each environment for each $\omega$ values defined above and for the three work space partitions, generating the graph with three different searches: unidirectional, sequential bidirectional and parallel bidirectional, denoted Graph 1, Graph 2 and Graph 3, respectively.

Fig. 3. Heat distribution used to compute $D_2$, higher values represent a higher dangerousness on the environment.
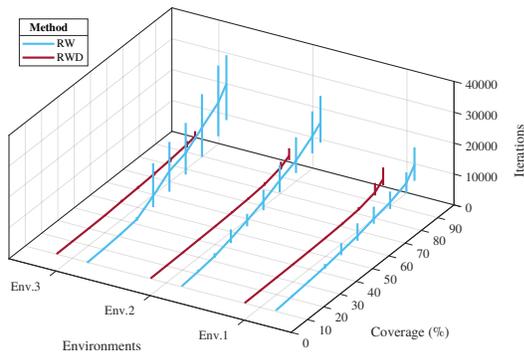


Fig. 4. Comparison between the performance between RW and RWD in all three environments, the graph shows the number of iterations need it to cover certain percentages of the terrain.
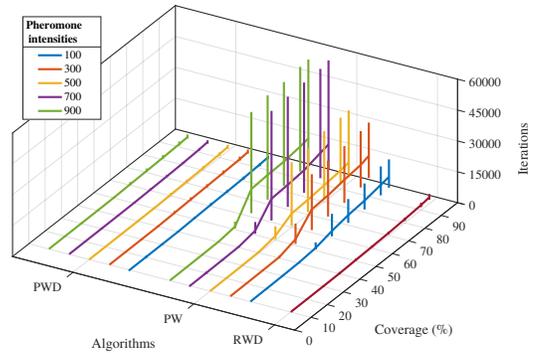
### B. Results and Discussion

Comparing the results of the RW and RWD algorithms, it is clear that adding a direction to the agent helps to improve the performance of the swarm. In Fig. 4 is shown that in the three environments RWD reduces the iterations needed for cover each of them by more than 50%.
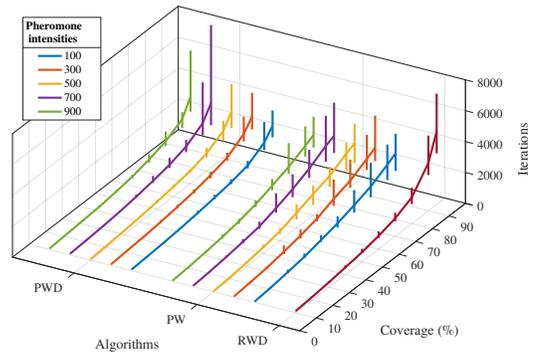
Analyzing RWD against the two pheromone methods, we can see that under certain configuration of $\Upsilon^0$, the performance of PW can be similar to RWD, however is PWD the one that stand out (Fig. 5). Given this results we can assure that if we add direction to the movement of an agent and give him the ability to communicate its path through the environment with pheromones, the dispersion and coverage of the swarm in the area can be accomplish in a more faster and efficient way.

Fig. 5 also shows that, as expected, while the coverage percentage of the environments increases, the amount of iterations necessary to achieve it also increase. This situation occurs because unexplored points are often very sparse in the environment, which makes difficult for the swarm to find them. This fact also causes that the standard deviation (SD) grows along with the percentage of coverage. A special result is obtained with the Env1 (Fig. 5(a)) in which the data presented a big SD. This occurs because the structuring of the environment, in which, there is only one connection path between two large unexplored areas .
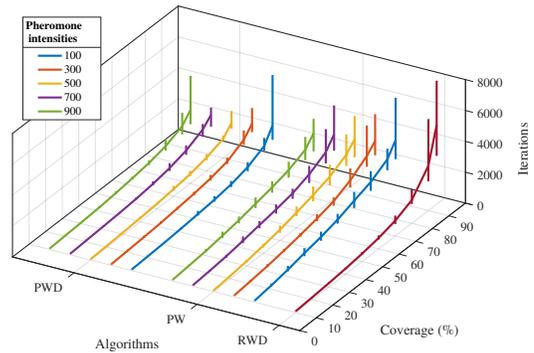
Fig. 6 presents a summary of the PWD results in every



(a) Env1



(b) Env2



(c) Env3

Fig. 5. Results of iterations needed it to cover certain percentages of the environments by the algorithms RWD, PW, PWD
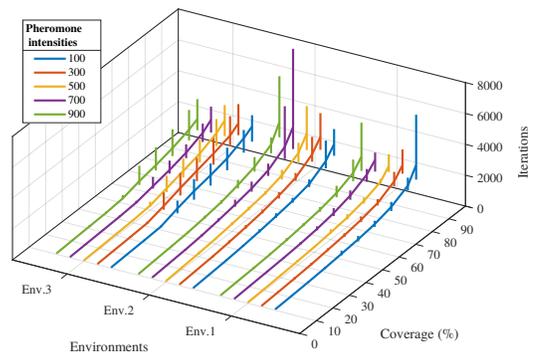


Fig. 6. Summary of the results of PWD in all three environments

(a) Env1 (10 agents and 1800 iterations)

(b) Env1 (20 agents and 912 iterations)

(c) Env2 (10 agents and 2806 iterations)

(d) Env2 (20 agents and 1388 iterations)

(e) Env3 (10 agents and 2298 iterations)

(f) Env3 (20 agents and 845 iterations)

Fig. 8. Cellstep heatmaps that shows the number of visits for each cell of the environment. The size of the swarm and the number of iterations need it to cover the 95% of the area is described in each case.
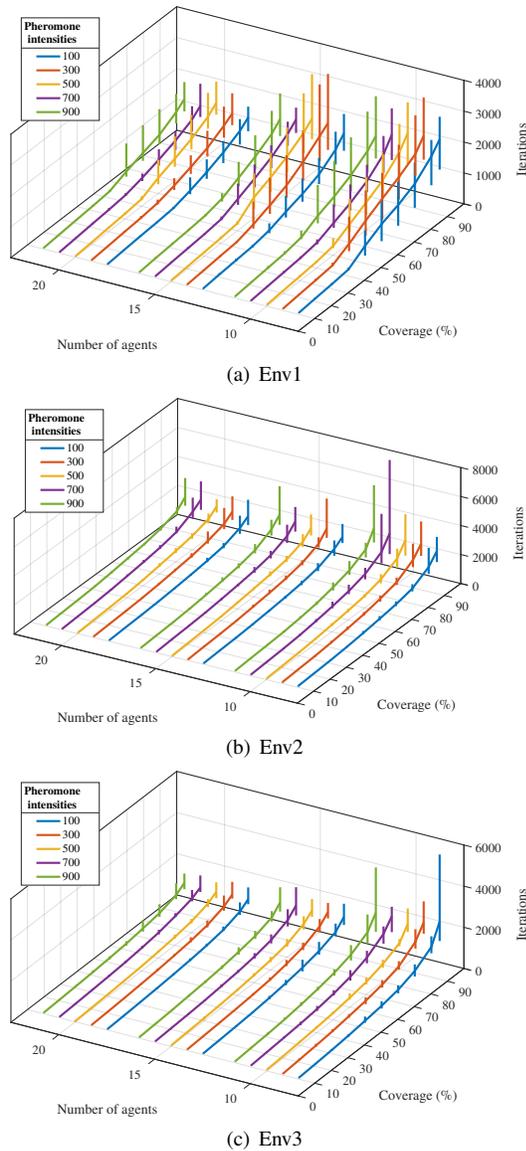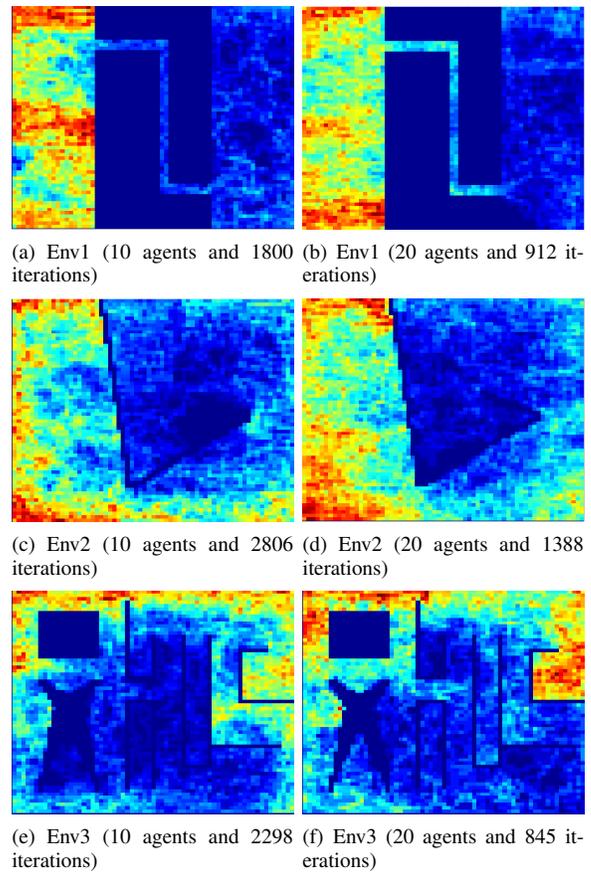


(a) Env1



(b) Env2



(c) Env3

Fig. 7. Results of the iterations need it to cover certain percentages of the environments by the algorithm PWD with three swarm of 10, 15 and 20 agents

environment and in certain percentages of coverage. In this we can see that there is not an ideal initial pheromone intensity that gives best results in all the configurations, because this will vary depending the environment an the percentage coverage. However we can hypothesize that a low level of $\Upsilon^0$ will produce similar results to the RWD while highest one will create a very saturated space and will not allow a faster dispersion.

All exploration results were performed by a swarm of ten agents. For comparing the changes that can produce the introduction of a bigger group of simulated robots, the same experiments were preformed by a swarm of fifteen and twenty agents. As expected, while the number of agents increases the faster is the coverage and dispersion of the area (Fig. 7). Also it can be seen that the SD of the data and the differences between the pheromone intensities decreases.

A cellstep heatmap, that shows the number of visits for

each cell of the environment, was also created in all the experiments in order to analyze the behavior of the agents. The results show that there are areas that are constantly being explored, this creates what we call over-exploration of the environments. This situation causes that the agents invest time on exploring areas that are already known instead of visiting unexplored ones, what can result on a increases in the iteration need it to accomplish the overall task. In the examples showed in Fig. 8, it can be seen that the area where the agents were deployed are usually highly explored. Also even though the size of the swarm augmented and the number of iterations reduced, the over-exploration is not necessarily decreased.

At the same time that exploration is occurring, a external server stores all visited cells for construct an environment map. Fig. 9 shows the evolutionary process of creating the map on Env3, guided by the coverage percentage of area. Once the exploration is finished, the Adapted RRG algorithm proposed takes this map to structure the space and find a path, minimizing the distance but taking into account the dangerousness.

Fig. 10 shows a small variations on graph generation time by means of unidirectional and sequential bidirectional search, but variations resulting of parallel bidirectional search are slightly larger. Regardless of search type used, execution time of graph generation increases as mesh becomes more refined. Execution time with parallel bidirectional search is always
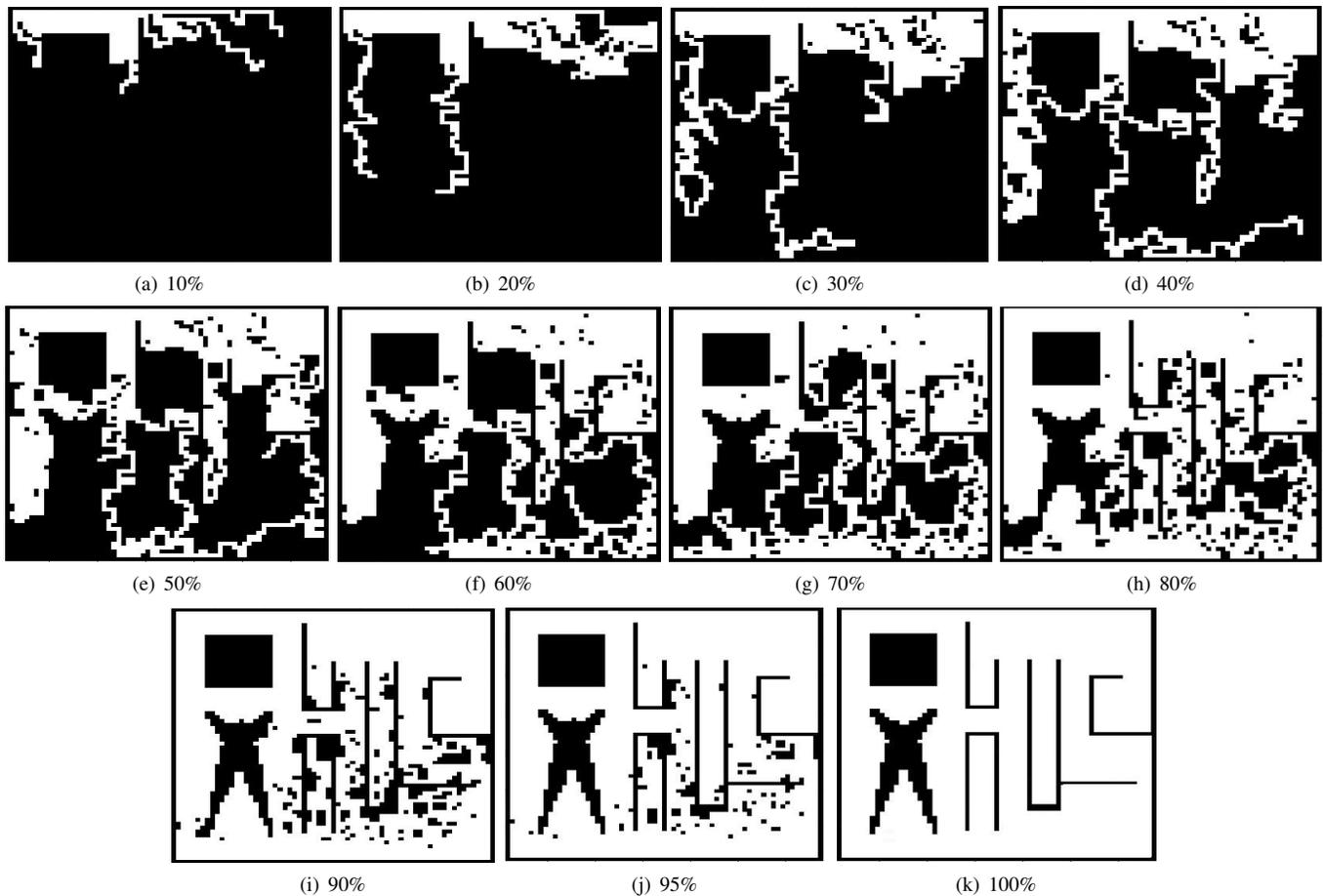
(a) 10%          (b) 20%          (c) 30%          (d) 40%

(e) 50%          (f) 60%          (g) 70%          (h) 80%

(i) 90%          (j) 95%          (k) 100%

Fig. 9.    Example of the map creation process of Env3 carried out during the exploration phase, according the coverage percentage of the area

TABLE I.    AVERAGE ELAPSED TIME IN 100 EXPERIMENTS BY FINDING THE OPTIMAL PATH, WITH DIFFERENT PARTITION SIZES AND $\omega$ VALUES ON ENV3

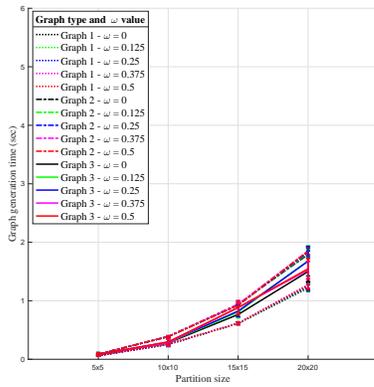| Partition size | Graph type | $\omega = 0$ | | $\omega = 0.125$ | | $\omega = 0.25$ | | $\omega = 0.375$ | | $\omega = 0.5$ | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | mean | std | mean | std | mean | std | mean | std | mean | std |
| | Graph 1 | 0.4268 | 0.1236 | 0.4127 | 0.1135 | 0.4403 | 0.1431 | 0.4406 | 0.1004 | 0.4393 | 0.2025 |
| $5 \times 5$ | Graph 2 | 0.5344 | 0.1158 | 0.4380 | 0.0517 | 0.6308 | 0.1609 | 0.5717 | 0.1844 | 0.4810 | 0.0845 |
| | Graph 3 | 0.6740 | 0.1957 | 0.7324 | 0.2828 | 0.7331 | 0.2666 | 0.7029 | 0.2407 | 0.7215 | 0.2549 |
| | Graph 1 | 5.2220 | 0.2101 | 5.2524 | 0.2006 | 5.2481 | 0.2114 | 5.3179 | 0.1993 | 5.3145 | 0.2004 |
| $10 \times 10$ | Graph 2 | 5.4593 | 0.1860 | 5.5821 | 0.2236 | 5.6681 | 0.1879 | 5.6588 | 0.2200 | 5.6591 | 0.2328 |
| | Graph 3 | 6.8302 | 0.4058 | 7.2021 | 0.7134 | 6.8418 | 0.3448 | 6.9627 | 0.4594 | 7.1016 | 0.6783 |
| | Graph 1 | 19.7240 | 0.4468 | 20.5150 | 0.5672 | 20.6680 | 0.4568 | 20.5630 | 0.5486 | 20.2910 | 0.5838 |
| $15 \times 15$ | Graph 2 | 19.0830 | 0.5272 | 20.0130 | 0.5316 | 20.0340 | 0.5244 | 20.4890 | 0.5410 | 20.2180 | 0.5674 |
| | Graph 3 | 22.1300 | 1.4991 | 24.4360 | 2.2988 | 25.1070 | 2.0405 | 25.0750 | 2.0488 | 25.4720 | 2.8845 |
| | Graph 1 | 41.7990 | 1.1626 | 44.9960 | 1.0414 | 45.3290 | 0.8837 | 45.9970 | 0.9466 | 46.6880 | 0.9544 |
| $20 \times 20$ | Graph 2 | 47.0210 | 1.1203 | 51.5370 | 1.0564 | 50.4570 | 1.2023 | 51.2210 | 1.1727 | 51.7330 | 0.9533 |
| | Graph 3 | 45.3380 | 4.0262 | 43.4460 | 4.0006 | 42.3570 | 2.3770 | 42.8490 | 2.6175 | 45.2010 | 4.1474 |

lower than execution time with unidirectional search. Elapsed graph generation time on Env3 is always greater than the elapsed on Env1 and Env2.

Once the graphs have been generated, DSP algorithm is applied on the same way for three cases: unidirecional, sequential and parallel bidirectional graph search. Table I summarizes the elapsed time at path planning process on Env3, we can see that the variations on Graph 3 (parallel bidirectional search) are also greater than Graph 1 and 2. There are no significant differences in path planning time with different $\omega$ values into the same partition, but there are significant differences between path planning time into different partition sizes. Path planning time increases as partition size increases. Results for Env1 and
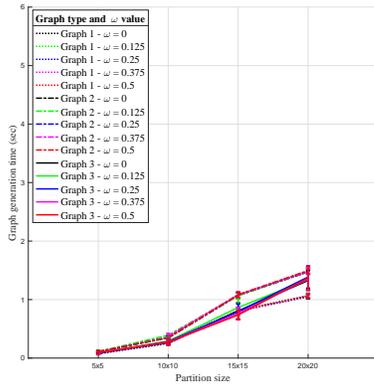
Env2 are quite similar.

The path distance depends on the environment as shown in Fig. 11. Although there are also slight differences between path distance obtained with different $\omega$ values, the greatest differences appear with different partition size. Path distance presents larger variations with a coarse mesh, i.e., with a small partition sizes. Better results are obtained as the number of cells increases.
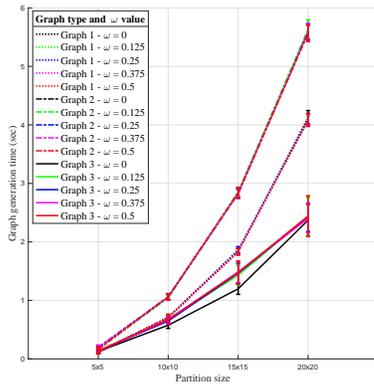
Other general results were made, equivalent to those showed in [24]. The shortest path found improve according to refinement of the mesh, but it also depends on the type of environment. For example $5 \times 5$ partition is not always

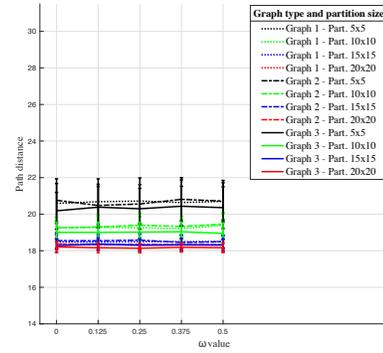(a) Elapsed time by adapted RRG algorithm on Env1

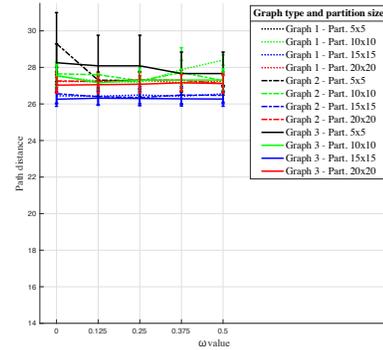

(b) Elapsed time by adapted RRG algorithm on Env2



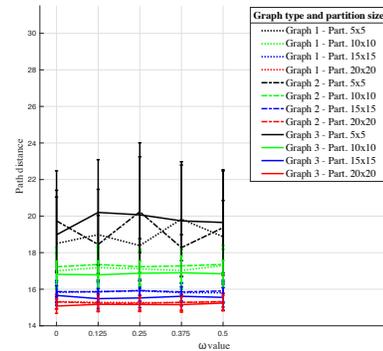(c) Elapsed time by adapted RRG algorithm on Env3

Fig. 10. Average elapsed time of 100 experiments generating graph on each environment, with uni and bidirectional search methods and based on different partition sizes



(a) Path distance obtained on Env1



(b) Path distance obtained on Env2



(c) Path distance obtained on Env3

Fig. 11. Average path distance of 100 experiments in graph generation on each environment, with uni and bidirectional search methods and based on different $\omega$ values

successful, not only because in some cases it is not possible to structure all space with the graph and obtain a path, but also because sometimes it is possible to obtain a path but not necessarily the optimal one. Moreover, it is possible to obtain different paths based on different partitions. Even though a finer mesh allows being closer to the optimal solution, it produces a denser and more complex graph.

Fig. 12 shows the performance of graph generation and path planning process in Env3. It presents some differences when using a scalarized objetive function, which combines distance and dangerousness, and when using an objective function
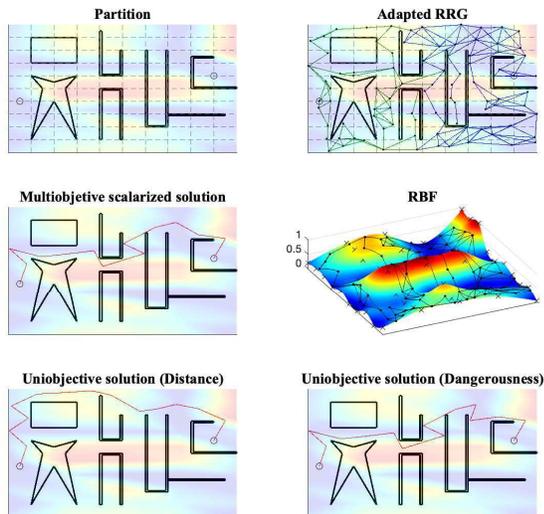
Fig. 12. Experimental results on Env3 with $10 \times 10$ as a partition size, $\omega =0.375$ and an parallel bidirectional search

that optimize just one of these elements. It is possible to see the RBF resulting with which dangerousness is calculated and the correponding line integrals.

## V. Conclusions

A solution that involves swarm robotics and multiobjectives techniques for finding paths, minimizing distance and danger in an unknown environment, was presented. The solution was composed by three phases: exploration, mapping and path planning, and was made under computer simulations with swarms of ten, fifteen and twenty agents.

Exploration and mapping phases were implemented by means of a cellular automata, in which the environments were mapped into a two dimensional map. For the first phase four algorithms, organized into two approaches, were presented and applied in three different environments.

The best performance scheme was PWD, which provides a direction to the agent movements and give him the ability to communicate its path through the environment. PWD achieves the complete tasks faster and more efficiently than the others. Even though all the experiments were carried out with five different initial amount of pheromone intensity, the ideal value varies depending on the structure of the space and the coverage percentage.

As was expected, the coverage time decreases with increase the number of agents. Besides, as the number of agents increases, differences between results with different pheromone intensities also decrease.

On the other hand, $\omega$ values does not interfere with the graph generation and graph type (unidirectional, sequential and parallel bidirectional search) does not affect the optimization process. But both, graph generation and optimization process, always depend on partition sizes.

A fine partition allows a complete work space structuring and to find the optimal path. Shapes and quantity of obstacles scattered, through the environment, will affect the effectiveness of the partition dimension. Results show that even though a more refined partition produces a shorter path, the graph generation will take longer execution time and produce a denser and more complex graph.

Parallel bidirectional search has better performance than unidirectional search, but sequential bidirectional search sometimes has a shorter run-time than parallel one, depending on the environment.

As future work a more robust algorithm has to be created in order to reduce the over-exploration of the environments and a genetic multiobjective optimization algorithm will be used, without scalarizing.

## References

[1] J. C. Barca and Y. A. Sekercioglu, "Swarm robotics reviewed," *Robotica*, vol. 31, no. 3, pp. 345–359, 2013.

[2] I. Navarro and F. Matía, "An introduction to swarm robotics," *Isrn robotics*, vol. 2013, 2012.

[3] L. Bayındır, "A review of swarm robotics tasks," *Neurocomputing*, vol. 172, pp. 292–321, 2016.

[4] L. Garattoni and M. Birattari, "Autonomous task sequencing in a robot swarm," *Science Robotics*, vol. 3, no. 20, p. eaat0430, 2018.

[5] N. Palmieri, X.-S. Yang, F. De Rango, and S. Marano, "Comparison of bio-inspired algorithms applied to the coordination of mobile robots considering the energy consumption," *Neural Computing and Applications*, pp. 1–24, 2017.

[6] Q. Tang, F. Yu, Y. Zhang, L. Ding, and P. Eberhard, "A stigmergy based search method for swarm robots," in *International Conference on Swarm Intelligence*. Springer, 2017, pp. 199–209.

[7] D. A. Lima and G. M. Oliveira, "A cellular automata ant memory model of foraging in a swarm of robots," *Applied Mathematical Modelling*, vol. 47, pp. 551–572, 2017.

[8] C. R. Tinoco, D. A. Lima, and G. M. Oliveira, "An improved model for swarm robotics in surveillance based on cellular automata and repulsive pheromone with discrete diffusion," *International Journal of Parallel, Emergent and Distributed Systems*, pp. 1–25, 2017.

[9] S. M. LaValle, *Planning algorithms*. Cambridge university press, 2006.

[10] R. Valencia and J. Andrade-Cetto, "Path planning in belief space with pose slam," in *Mapping, Planning and Exploration with Pose SLAM*. Springer, 2018, pp. 53–87.

[11] J. Faigl, "On self-organizing map and rapidly-exploring random graph in multi-goal planning," in *Advances in self-organizing maps and learning vector quantization*. Springer, 2016, pp. 143–153.

[12] A. Bircher, K. Alexis, U. Schwesinger, S. Omari, M. Burri, and R. Siegwart, "An incremental sampling-based approach to inspection planning: the rapidly exploring random tree of trees," *Robotica*, vol. 35, no. 6, pp. 1327–1340, 2017.

[13] A. D. Mali, "Probabilistic roadmaps with higher expressive power," *International Journal on Artificial Intelligence Tools*, vol. 25, no. 04, p. 1650027, 2016.

[14] S. Karaman and E. Frazzoli, "Sampling-based algorithms for optimal motion planning," *The International Journal of Robotics Research*, vol. 30, no. 7, pp. 846–894, 2011.

[15] I. Noreen, A. Khan, and Z. Habib, "A comparison of rrt, rrt* and rrt*-smart path planning algorithms," *International Journal of Computer Science and Network Security (IJCSNS)*, vol. 16, no. 10, pp. 20–27, 2016.

[16] B. Saicharan, R. Tiwari, and N. Roberts, "Multi objective optimization based path planning in robotics using nature inspired algorithms: A survey," in *International Conference on Power Electronics, Intelligent Control and Energy Systems (ICPEICES)*. IEEE, 2016, pp. 1–6.

[17] G. Sharon, R. Stern, A. Felner, and N. R. Sturtevant, "Conflict-based search for optimal multi-agent pathfinding," *Artificial Intelligence*, vol. 219, pp. 40–66, 2015.

[18] H. Wang, Y. Yu, and Q. Yuan, "Application of dijkstra algorithm in robot path-planning," in *Second International Conference on Mechanic Automation and Control Engineering (MACE)*. IEEE, 2011, pp. 1067–1069.

[19] R. Arjun and P. Reddy, "Shama, and m. yamuna,"research on the optimization of dijkstra's algorithm and its applications"," *International Journal of Science Technology and Management*, vol. 4, no. 1, pp. 304–309, 2015.

[20] C. Calderón-Arce and P. Alvarado-Moya, "Optimización multiobjetivo con funciones de alto costo computacional. revisión del estado del arte," *Revista Tecnología en Marcha*, pp. 16–24, 2016.

[21] N. Gunantara, "A review of multi-objective optimization: Methods and its applications," *Cogent Engineering*, vol. 5, no. 1, pp. 1–16, 2018.

[22] C. Dimidov, G. Oriolo, and V. Trianni, "Random walks in swarm robotics: an experiment with kilobots," in *International Conference on Swarm Intelligence*. Springer, 2016, pp. 185–196.

[23] R. Morlok and M. Gini, "Dispersing robots in an unknown environment," in *Distributed Autonomous Robotic Systems 6*. Springer, 2007, pp. 253–262.

[24] C. Calderón-Arce, R. Solís-Ortega, and T. Bustillos-Lewis, "Path planning on static environments based on exploration with a swarm robotics and rrg algorithms," in *38th Central America and Panama Convention (CONCAPAN XXXVIII)*. IEEE, 2018, pp. 1–6.

[25] T. C. Hales, "The jordan curve theorem, formally and informally," *American Mathematical Monthly*, vol. 114, no. 10, pp. 882–894, 2007.

[26] H. Yu, T. Xie, S. Paszczyñski, and B. M. Wilamowski, "Advantages of radial basis function networks for dynamic system design," *Transactions on Industrial Electronics*, vol. 58, no. 12, pp. 5438–5450, 2011.

[27] G. A. Barnett, N. Flyer, and L. J. Wicker, "An rbf-fd polynomial method based on polyharmonic splines for the navier-stokes equations: Comparisons on different node layouts," *arXiv preprint arXiv:1509.02615*, 2015.

# Social Media Cyberbullying Detection using Machine Learning

John Hani[1], Mohamed Nashaat[2], Mostafa Ahmed[3],
Zeyad Emad[4], Eslam Amer[5]
Department of Computer Science,
Misr International University, Cairo, Egypt

Ammar Mohammed[6]
Department of Computer Science, ISSR
Cairo University, Cairo, Egypt

*Abstract*—With the exponential increase of social media users, cyberbullying has been emerged as a form of bullying through electronic messages. Social networks provides a rich environment for bullies to uses these networks as vulnerable to attacks against victims. Given the consequences of cyberbullying on victims, it is necessary to find suitable actions to detect and prevent it. Machine learning can be helpful to detect language patterns of the bullies and hence can generate a model to automatically detect cyberbullying actions. This paper proposes a supervised machine learning approach for detecting and preventing cyberbullying. Several classifiers are used to train and recognize bullying actions. The evaluation of the proposed approach on cyberbullying dataset shows that Neural Network performs better and achieves accuracy of 92.8% and SVM achieves 90.3. Also, NN outperforms other classifiers of similar work on the same dataset.

*Keywords*—*Cyberbullying; machine learning; neural network*

## I. Introduction

With the increasing number of users on social media leads to a new way of bullying. The later term, is defined as an intentional or an aggressive acts which are carried out by person or groups of individuals using repeatedly communication messages over time against a victim who cannot easily defend him or herself [1]. Bullying has always been a part of society. With the inception of the internet, it was only a matter of time until bullies found their way on to this new and opportunistic medium. Using services like email and instant messenger, bullies became able to do their nasty deeds with anonymity and great distance between them and their targets. According to Cambridge dictionary the term cyberbullying is defined as the activity of using the internet to harm or frighten another person, especially by sending them unpleasant messages. The main factor that separates cyberbullying from traditional bullying is the effect that it has on the victim. Traditional bullying may end in physical damage as well as emotional and psychological damage, as opposed to cyberbullying, where it is all emotional and psychological.

Given the consequences of cyberbullying on victims, it is urgently needed to find a proper actions to detect and hence to prevent it. One of the successful approaches that learns from data and generates a model that automatically classifies proper actions is machine learning. Machine learning can be helpful to detect language patterns of the bullies and hence can generate a model to detect cyberbullying actions. Thus, the main contribution of this paper is to propose a supervised machine learning approach for detecting and preventing cyberbullying. The proposed approach is evaluated on

a cyberbullying dataset from kaggle which was collected and labeled by the authors Kelly Reynolds et al. in their paper [2]. The performance of SVM and Neural Network classifiers are compared on both TFIDF and sentiment analysis feature extraction methods. Furthermore, experiments were made on different n-gram language model. 2-gram, 3-gram and 4-gram has been taken into consideration during the evaluation of the model produced by the classifiers. Finally, we evaluate our proposed approach with previous related work who used the same data.

The rest of the paper is organized as follows. Section II shows several related work. Section III describes the proposed approach. Section IV shows the experimental results and the evaluation of the proposed approach. Finally, Section V concludes the paper.

## II. Related Work

There are many approaches that proposes systems which can detect cyberbullying automatically with high accuracy. First one is author Nandhini et al. [3] have proposed a model that uses Naïve Bayes machine learning approach and by their work they achieved 91% accuracy and got their dataset from MySpace.com, and then they proposed another model [4] Naïve Bayes classifier and genetic operations (FuzGen) and they achieved 87% accuracy. Another approach by Romsaiyud et al. [5] they enhanced the Naïve Bayes classifier for extracting the words and examining loaded pattern clustering and by this approach they achieved 95.79% accuracy on datasets from Slashdot, Kongregate, and MySpace. However, they have a problem that the cluster processes doesn't work in parallel. Moreover, in the approach proposed by Bunchanan et al. [6] they used War of Tanks game chat to get their dataset and manually classified them and then compared them to simple Naïve classification that uses sentiment analysis as a feature, their results were poor when compared to the manually classified results. Furthermore, Isa et al. [7] proposed an approach after getting their dataset from kaggle they used two classifier Naïve Bayes and SVM. The Naïve Bayes classifier yielded average accuracy of 92.81% while SVM with poly kernel yielded accuracy of 97.11%, but they did not mention their training or testing size of the dataset, so the results may not be credible. Another Approach by Dinakar et al. [8] that aimed to detect explicit bullying language pertaining to (1) Sexuality, (2) Race & Culture and (3) intelligence, they acquired their dataset from YouTube comment section. After applying SVM and Naïve Bayes classifiers, SVM yielded accuracy of 66%

and Naïve Bayes 63%. Moving on to Di Capua et al. [9], they proposed a new way for cyberbullying detection by adopting an unsupervised approach, they used the classifiers inconsistently over their dataset, applying SVM on FormSpring and achieving 67% on recall, applying GHSOM on YouTube and achieving 60% precision, 69% accuracy and 94% recall, applying Naïve Bayes on Twitter and achieving 67% accuracy. Additionally, Haidar et al. [10] proposed a model to detect cyberbullying but using Arabic language they used Naïve Bayes and achieved 90.85% precision and SVM achieved 94.1% as precision but they have high rate of false positive also the are work on Arabic language.

Another type of approaches using Deep Learning and Neural Networks. One of the proposed methods is Zhang et al. [11] in their paper uses novel pronunciation based convolution neural network (PCNN), thereby alleviating the problem of noise and bullying data sparsity to overcome class imbalance. 1313 messages from twitter, 13,000 messages from formspring.me. Accuracy of the twitter dataset wasn't calculated due to it being imbalanced. While Achieving 56% on precision, 78% recall and 96% accuracy, while achieving high accuracy their dataset was unbalanced, so that gives false results and that reflects in precision score which is 56%. The authors Nobata et al. [12] showed that using abusive language has increased recently, They used a framework called Vowpal wabbit for classification, and they also developed a supervised classification methodology with NLP features that outperform deep learning approach, The F-Score reached 0.817 using dataset collected from comments posted on Yahoo News and Finance.

Zhao et al. [13] proposed framework specific for cyberbullying detection, they used word embedding that makes a list of pre-defined insulting words and assign weights to obtain bullying features, they used SVM as their main classifier and got recall of 79.4%. Then another approach was proposed by Parime et al. [14] they got their dataset from MySpace and manually marked them and they used the SVM Classifier for the classification. Moreover, Chen et al. [15] proposed a new feature extraction method called Lexical Syntactic Feature and SVM as their classifier and they achieved 77.9% precision and 77.8% recall. Furthermore, Ting et al. [16] proposed a technique based on SNM, they collected their data from social media and then used SNA measurements and sentiments as features. Seven experiments were made and they achieved around 97% precision and 71% as recall. Furthermore, Harsh Dani et al. [17] introduced a new framework called SICD, they used KNN for classification. Finally, they achieved 0.6105 F1 score and 0.7539 AUC score.

SVM classifier was one of the approaches used in the research papers. Dadvar et al. [18][19][20][21] have constructed in the first and second paper a Support Vector Machine classifier using WEKA, their dataset was collected from Myspace. They achieved 43% on precision, 16% in recall and they didn't mention the accuracy, the only difference between the two papers is that they used gender information in classification in the second paper. Moreover, in their second paper 4626 comments from 3858 distinct users were collected. The comments were manually labelled as bullying (9.7%) and non-bullying (inter-annotator agreement 93%). SVM classifier was applied by them and were able to reach results of up to 78% on precision

and 55% on recall. Finally, in their third paper they applied 3 models for their dataset gathered from YouTube comment section: Multi-Criteria Evaluation Systems (MCES), machine learning: (Naïve Bayes classifier, decision tree, SVM), Hybrid approach. The MCES score 72% on accuracy, while Naïve Bayes scored the highest out of the three with 66%. Moving on to another author, Potha et al. [22] have also used the SVM approach and achieved 49.8% result on accuracy. While Chavan et al. [23] used two classifiers: logistic regression and support vector machine. The logistic regression achieved 73.76 accuracy and 60% recall and 64.4% Precision. While for the support vector machine they achieved 77.65% accuracy and 58% recall and 70% precision's and they got their dataset from Kaggle.

### III. PROPOSED APPROACH

The proposed approach, as seen in Fig. 1, contains three main steps: Preprocessing, features extraction and classification step. In the preprocessing step we clean the data by removing the noise and unnecessary text. The preprocessing step is done in the following:

- Tokenization: In this part we take the text as sentences or whole paragraphs and then output the entered text as separated words in a list.

- Lowering text: This takes the list of words that got out of the tokenization and then lower all the letters Like: 'THIS IS AWESOME' is going to be 'this is awesome'.

- Stop words and encoding cleaning: This is an essential part of the preprocessing where we clean the text from those stop words and encoding characters like \n or \t which do not provide a meaningful information to the classifiers.

- Word Correction: In this part we used Microsoft Bing word correction API [24] that takes a word and then return a JSON object with the most similar words and the distance between these words and the original word.
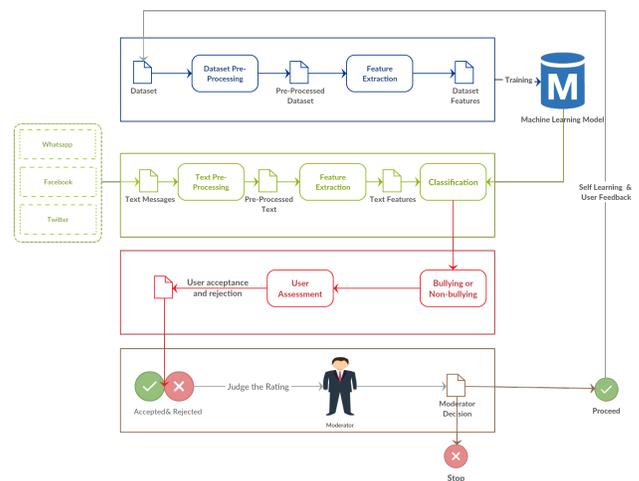


Fig. 1. Proposed Approach

The second step of the proposed Model is the features extraction step. In this step the textual data is transformed into a suitable format applicable to feed into machine learning algorithms. First we extract the features of the input data using TFIDF[25] as and put them in a features list. The key idea of TFIDF is that it works on the text and get the weights of the words with respect to the document or sentence. In Addition to TFIDF, we use sentiment analysis technique[26] to extract the polarity of the sentences and add them as a feature into the features list containing the TFIDF features. The polarity of the sentences means that if the sentence is classified as positive or negative. For that purpose we extract the polarity using Text Blob library[27] which is a pre-trained model on movie reviews. In addition to the feature extraction using TFIDF and sentiment polarity extraction, the propose approach uses N-Gram[28] to consider the different combinations of the words during evaluation of the model. Particularly, we use used 2-Gram, 3-Gram and 4-Gram.

The last step in the proposed approach is the classification step where the extracted features are fed into a classification algorithm to train, and test the classifier and hence use it in the prediction phase. We used two classifiers, namely, SVM (Support Vector Machine) and Neural Network. The neural network contains three layers: Input, hidden, output layer. In the input layer, it consists of 128 nodes. In the hidden layer, it contains 64 neurons. The output layer is a Boolean output.

Generally, the evaluation of classifiers is done using several evaluation matrices depends on the confusion matrix. Among of those criteria are Accuracy, precision, recall and f-score. They are calculated according to the following equations:

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} \qquad (1)$$

$$Precision = \frac{TP}{TP+FP} \qquad (2)$$

$$Recall = \frac{TP}{TP+FN} \qquad (3)$$

$$F - Score = \frac{2*precision*recall}{precision+recall} \qquad (4)$$

Where TP represents the number of true positive, TN represents the number of true negatives, FP represents the number of false positives, and FN represents the number of false negatives classes.

## IV. EXPERIMENTAL RESULTS

This section describes the experimental results on the proposed approach. We evaluate the proposed approach on the cyberbullying dataset from kaggle. In the following we describes the Data and the results.

### A. Data Description

We have used cyberbullying dataset from Kaggle which was collected and labeled by the authors Kelly Reynolds et al. in their paper [2]. This dataset contains in general 12773 conversations messages collected from Formspring. The dataset contains questions and their answers annotated with either cyberbullying or not. The annotation classes were unbalanced distributed such that 1038 question-answering instances out of 12773 belongs to the class cyberbullying, while 11735 belongs to the other class. First, to remedy the data unbalancing, we take the same number instances of both classes to measure the accuracy. We also removed from the data big size conversations and remove the noisy data. We ended up with total 1608 instance conversations where 804 instances belongs to each class. Table I summarizes the statistics of dataset.

TABLE I.    STATISTICS OF THE DATASET

| Total number of Conversations | 1608 |
|---|---|
| Number of cyberbullying | 804 |
| Number of non-Cyberbullying | 804 |
| Number of distinct words | 5628 |
| Number of token | 48843 |
| Maximum Conversation size | 773 Characters |
| Minimum Conversation size | 59 Characters |

### B. Results

After preprocessing the dataset, we follow the same step presented in Section III to extract the features. We then split the dataset into ratios (0.8,0.2) for train and test. Accuracy, recall and precision, and f-score are taken as a performance measure to evaluate the classifiers. We apply SVM as well as Neural Network (NN) as they are among the best performance classifiers in the literature. We run several experiments on different n-gram language model. In Particular, we take into consideration 2-gram, 3-gram, and 4-gram during the evaluation of the model produced by the classifiers. Table II summarizes the accuracy of both SVM and NN. The SVM classifier achieved the highest percentage using 4-Gram with accuracy 90.3% while the NN achieved highest accuracy using 3-Gram with accuracy 92.8%. It is found that the average accuracy of all n-gram models of NN achieves 91.76%, while the average accuracy of all n-gram models of SVM achieves 89.87%. Fig. 2 depicts the accuracy results of both classifiers.

TABLE II.    THE ACCURACY OF SVM AND NN IN DIFFERENT LANGUAGE MODEL

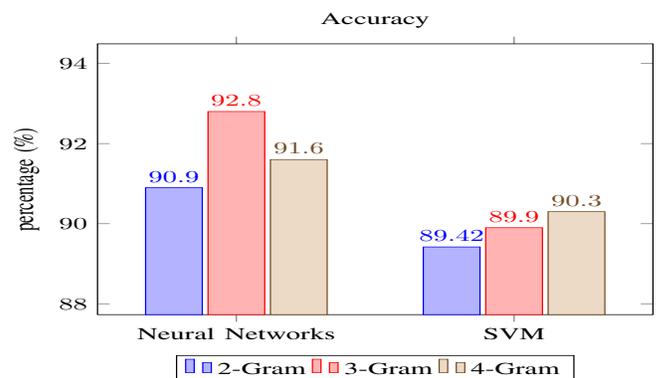| Classifier | 2-Gram | 3-Gram | 4-Gram | Average |
|---|---|---|---|---|
| SVM | 89.42% | 89.9% | **90.3%** | 89.87% |
| Neural Network | 90.9% | **92.8%** | 91.6% | 91.76% |



Fig. 2.   Comparison between SVM and Neural Network in Terms of Accuracy

In addition to accuracy, Table III and Table IV show the evaluations of both classifiers in terms of precision and recall respectively for each language model. The trade-off between recall and precision is shown in Table V which represents the f-score of both classifiers in the different language model. Table V summarizes the f-score of both SVM and NN. The SVM classifier achieved the highest f-measure using 4-Gram with f-score 90.3% while the NN achieved highest f-measure using 2-Gram with f-score 92.2%. It is found that the average f-score of all n-gram models of NN achieves 91.9%, while the average f-score of all n-gram models of SVM achieves 89.8%. Fig. 3 summarizes the f-score of the classification of the SVM and Neural Network. The results of average accuracy as well as the average f-score indicate that NN performs better than SVM.

In addition to the previous experiments, we evaluate and compare our classifiers on the proposed approach with the work of [23]. In this work, they used logistic regression and SVM for classification and used the same data. Moreover, we have calculated the average accuracy, recall, precision and F-score of our two classifiers. The summary of results is shown in Table VI. To compare the work, it is found that our proposed NN model outperforms all other classifiers and is ranked as the best results in terms of average accuracy and F-Score achieving accuracy 91.76% and f-score 91.9%. In Fig. 4 we are comparing between our best classifier with their best classifier in case of accuracy. Finally, here in Fig. 5 we are comparing between our best classifier with their best classifier in case of F-Measure.

TABLE III. RECALL OF SVM AND NN

| Classifier | 2-Gram | 3-Gram | 4-Gram | Average |
|---|---|---|---|---|
| SVM | 89.42% | 90.3% | **90.8%** | 90.1% |
| Neural Network | 91.6% | 91.5% | **92%** | 91.7% |

TABLE IV. PRECISION OF SVM AND NN

| Classifier | 2-Gram | 3-Gram | 4-Gram | Average |
|---|---|---|---|---|
| SVM | 89.42% | 89.5% | **90%** | 89.6% |
| Neural Network | **93%** | 92.5% | 91.7% | 92.4% |

TABLE V. F-SCORE OF SVM AND NN

| Classifier | 2-Gram | 3-Gram | 4-Gram | Average |
|---|---|---|---|---|
| SVM | 89.42% | 89.8% | **90.3%** | 89.8% |
| Neural Network | **92.2%** | 91.9% | 91.8% | 91.9% |

TABLE VI. COMPARISON WITH RELATED WORK

| | Classifier | Avg. Accuracy | Avg. Recall | Avg. Precision | Avg. F-Score |
|---|---|---|---|---|---|
| Vikas S Chavan | Logistic regression | 73.76 | 61.47% | 64.4% | 62.9% |
| | SVM | 77.65% | 58.29% | 70.29% | 63.7% |
| Current Results | Neural Network | **91.76%** | **91.7%** | **92.4%** | **91.9%** |
| | SVM | 89.87% | 90.1% | 89.6% | 89.8% |



Fig. 4. Comparison between the Best Classifiers in Terms of Accuracy



Fig. 3. Comparison between SVM and Neural Network in Terms of F-Measure
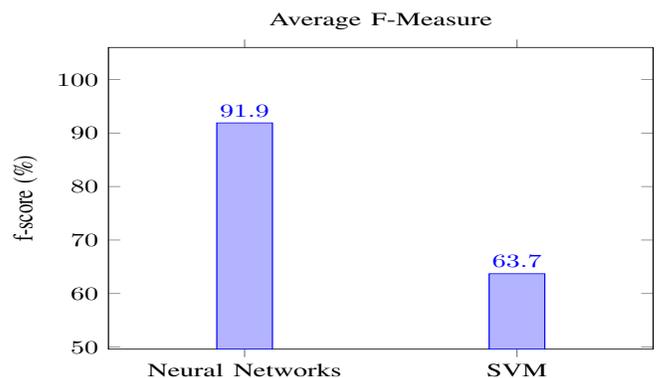


Fig. 5. Comparison between the Best Classifiers in Terms of F-Measure

## V. Conclusion

In this paper, we proposed an approach to detect cyberbullying using machine learning techniques. We evaluated our model on two classifiers SVM and Neural Network and we used TFIDF and sentiment analysis algorithms for features extraction. The classifications were evaluated on different n-gram language models. We achieved 92.8% accuracy using Neural Network with 3-grams and 90.3% accuracy using SVM with 4-grams while using both TFIDF and sentiment analysis together. We found that our Neural Network performed better than the SVM classifier as it also achieves average f-score 91.9% while the SVM achieves average f-score 89.8%. Furthermore, we compared our work with another related work that used the same dataset, finding that our Neural Network outperformed their classifiers in terms of accuracy and f-score. By achieving this accuracy, our work is definitely going to improve cyberbullying detection to help people to use social media safely. However, detecting cyberbullying pattern is limited by the size of training data. Thus, a larger cyberbullying data is needed to improve the performance. Hence, deep learning techniques will be suitable in the larger data as they are proven to outperform machine learning approaches over larger size data.

## References

[1] Peter K Smith, Jess Mahdavi, Manuel Carvalho, Sonja Fisher, Shanette Russell, and Neil Tippett. Cyberbullying: Its nature and impact in secondary school pupils. *Journal of child psychology and psychiatry*, 49(4):376–385, 2008.

[2] Kelly Reynolds, April Kontostathis, and Lynne Edwards. Using machine learning to detect cyberbullying. In *2011 10th International Conference on Machine learning and applications and workshops*, volume 2, pages 241–244. IEEE, 2011.

[3] B Nandhini and JI Sheeba. Cyberbullying detection and classification using information retrieval algorithm. In *Proceedings of the 2015 International Conference on Advanced Research in Computer Science Engineering & Technology (ICARCSET 2015)*, page 20. ACM, 2015.

[4] B Sri Nandhini and JI Sheeba. Online social network bullying detection using intelligence techniques. *Procedia Computer Science*, 45:485–492, 2015.

[5] Walisa Romsaiyud, Kodchakorn na Nakornphanom, Pimpaka Prasertsilp, Piyaporn Nurarak, and Pirom Konglerd. Automated cyberbullying detection using clustering appearance patterns. In *Knowledge and Smart Technology (KST), 2017 9th International Conference on*, pages 242–247. IEEE, 2017.

[6] Shane Murnion, William J Buchanan, Adrian Smales, and Gordon Russell. Machine learning and semantic analysis of in-game chat for cyberbullying. *Computers & Security*, 76:197–213, 2018.

[7] Sani Muhamad Isa, Livia Ashianti, et al. Cyberbullying classification using text mining. In *Informatics and Computational Sciences (ICICoS), 2017 1st International Conference on*, pages 241–246. IEEE, 2017.

[8] Karthik Dinakar, Birago Jones, Catherine Havasi, Henry Lieberman, and Rosalind Picard. Common sense reasoning for detection, prevention, and mitigation of cyberbullying. *ACM Transactions on Interactive Intelligent Systems (TiiS)*, 2(3):18, 2012.

[9] Michele Di Capua, Emanuel Di Nardo, and Alfredo Petrosino. Unsupervised cyber bullying detection in social networks. In *Pattern Recognition (ICPR), 2016 23rd International Conference on*, pages 432–437. IEEE, 2016.

[10] Batoul Haidar, Maroun Chamoun, and Ahmed Serhrouchni. A multilingual system for cyberbullying detection: Arabic content detection using machine learning. *Advances in Science, Technology and Engineering Systems Journal*, 2(6):275–284, 2017.

[11] Xiang Zhang, Jonathan Tong, Nishant Vishwamitra, Elizabeth Whittaker, Joseph P Mazer, Robin Kowalski, Hongxin Hu, Feng Luo, Jamie Macbeth, and Edward Dillon. Cyberbullying detection with a pronunciation based convolutional neural network. In *2016 15th IEEE International Conference on Machine Learning and Applications (ICMLA)*, pages 740–745. IEEE, 2016.

[12] Chikashi Nobata, Joel Tetreault, Achint Thomas, Yashar Mehdad, and Yi Chang. Abusive language detection in online user content. In *Proceedings of the 25th international conference on world wide web*, pages 145–153. International World Wide Web Conferences Steering Committee, 2016.

[13] Rui Zhao, Anna Zhou, and Kezhi Mao. Automatic detection of cyberbullying on social networks based on bullying features. In *Proceedings of the 17th international conference on distributed computing and networking*, page 43. ACM, 2016.

[14] Sourabh Parime and Vaibhav Suri. Cyberbullying detection and prevention: Data mining and psychological perspective. In *Circuit, Power and Computing Technologies (ICCPCT), 2014 International Conference on*, pages 1541–1547. IEEE, 2014.

[15] Ying Chen, Yilu Zhou, Sencun Zhu, and Heng Xu. Detecting offensive language in social media to protect adolescent online safety. In *Privacy, Security, Risk and Trust (PASSAT), 2012 International Conference on and 2012 International Confernece on Social Computing (SocialCom)*, pages 71–80. IEEE, 2012.

[16] I-Hsien Ting, Wun Sheng Liou, Dario Liberona, Shyue-Liang Wang, and Giovanny Mauricio Tarazona Bermudez. Towards the detection of cyberbullying based on social network mining techniques. In *Behavioral, Economic, Socio-cultural Computing (BESC), 2017 International Conference on*, pages 1–2. IEEE, 2017.

[17] Harsh Dani, Jundong Li, and Huan Liu. Sentiment informed cyberbullying detection in social media. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 52–67. Springer, 2017.

[18] Maral Dadvar and Franciska De Jong. Cyberbullying detection: a step toward a safer internet yard. In *Proceedings of the 21st International Conference on World Wide Web*, pages 121–126. ACM, 2012.

[19] Maral Dadvar, de FMG Jong, Roeland Ordelman, and Dolf Trieschnigg. Improved cyberbullying detection using gender information. In *Proceedings of the Twelfth Dutch-Belgian Information Retrieval Workshop (DIR 2012)*. University of Ghent, 2012.

[20] Maral Dadvar, Dolf Trieschnigg, Roeland Ordelman, and Franciska de Jong. Improving cyberbullying detection with user context. In *European Conference on Information Retrieval*, pages 693–696. Springer, 2013.

[21] Maral Dadvar, Dolf Trieschnigg, and Franciska de Jong. Experts and machines against bullies: A hybrid approach to detect cyberbullies. In *Canadian Conference on Artificial Intelligence*, pages 275–281. Springer, 2014.

[22] Nektaria Potha and Manolis Maragoudakis. Cyberbullying detection using time series modeling. In *Data Mining Workshop (ICDMW), 2014 IEEE International Conference on*, pages 373–382. IEEE, 2014.

[23] Vikas S Chavan and SS Shylaja. Machine learning approach for detection of cyber-aggressive comments by peers on social media network. In *Advances in computing, communications and informatics (ICACCI), 2015 International Conference on*, pages 2354–2358. IEEE, 2015.

[24] Youssef Bassil and Mohammad Alwani. Post-editing error correction algorithm for speech recognition using bing spelling suggestion. *arXiv preprint arXiv:1203.5255*, 2012.

[25] Akiko Aizawa. An information-theoretic perspective of tf–idf measures. *Information Processing & Management*, 39(1):45–65, 2003.

[26] Theresa Wilson, Janyce Wiebe, and Paul Hoffmann. Recognizing contextual polarity in phrase-level sentiment analysis. In *Proceedings of Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing*, 2005.

[27] Steven Loria, P Keen, M Honnibal, R Yankovsky, D Karesh, E Dempsey, et al. Textblob: simplified text processing. *Secondary TextBlob: Simplified Text Processing*, 2014.

[28] William B Cavnar, John M Trenkle, et al. N-gram-based text categorization. In *Proceedings of SDAIR-94, 3rd annual symposium on document analysis and information retrieval*, volume 161175. Citeseer, 1994.

# A Blockchain-based Value Added Tax (VAT) System: Saudi Arabia as a Use-Case

Ahmad Alkhodre[1]
Islamic University of Madinah

Salman Jan[3]
Malaysian Institute of Information Technology
Universiti Kuala Lumpur & University of Peshawar

Shah Khusro[5]
University of Peshawar

Toqeer Ali[2]
Islamic University of Madinah

Yazed Alsaawy[4]
Islamic University of Madinah

Muhammad Yasar[6]
Malaysian Institute of Information Technology
Universiti Kuala Lumpur

*Abstract*—**Businesses need trust to confidently perform trade among each other. Centralized business models are the only mature solutions available to perform trades over the Internet. However, they have many problems which includes but are not limited to the fact that these create bottleneck on the server as well as requires trusted third parties. Recently, decentralized solutions have gained significant popularity and acceptance for future businesses. The wide acceptance of such systems is indeed due to the trust management among various untrusted business stakeholders. Many solutions have been proposed in this regard to provide de-centralized infrastructure for various business models. A standard solution that is acceptable to the industry is still in demand. Hyperledger umbrella Blockchain projects, that are supported by IBM and many other industry big players are gaining popularity due to its efficient and pluggable design. In this study, the author present the idea of utilizing Blockchain to design a Value-Added Tax (VAT) system for Saudi Arabia's newly introduced tax system. The reason to select this business model for VAT is twofold. First, it provides an untampered distributed ledger, which cannot be deceived by any party. Each transaction in the system cannot go unnoticed by the smart contract. Secondly, it provides a transparent record, and updates all involved parties regarding each activity performed by stakeholders. The newly proposed system will provide a transparent database of VAT transactions according to our smart contract design and at each stage of supply chain, tax will be deducted and stored on peer-to-peer network via consensus process. The author believes that the proposed solution will have significant impact on VAT collection in the Kingdom of Saudi Arabia.**

*Keywords*—*VAT; hyperledger; blockchain; consensus; decentralized network*

## I. Introduction

Value Added Tax (VAT), is a self credit mechanism of indirect accumulative tax related to consumption of goods or services which is imposed on the end user as chain of business [1], [2]. VAT is collected from individual customers and is paid back to the Government. VAT collection is an important process of a state administrative authority that allows a government to generate revenue from active honest tax payers and if the tax payment is not made properly then the dishonest payers may reduce the tax liabilities to their own level of choices. Consequent to non compliance in tax payment and tax fraud, countries and governments receive losses in tax collection [3], [4]. It has also been observed that significant

time and resources are consumed during the audit when the tax gap is identified.

VAT is implemented at the stage of production and distribution to increase the inland revenue. The noncompliance of VAT implementation results in the VAT gap [5] that represents the difference between estimated tax and actual tax collected. There are a number of reasons for the VAT gap which may include envision, fraud, bankruptcies, and insolvencies relating finance, etc [6].

In Saudi Arabia, VAT is introduced as a standard rate of 5% with effect from January 2018 (Rabi Al-Thani 14, 1439) with exception of some goods and services. The sectors for which VAT is avoided includes health, education and transportation [7]. However, VAT is implemented on services and things like healthcare and treatment [8], medicine and medical equipments, import and exports [9] clothes, fuel i.e. food, petrol, & diesel, utility bills and hotel rooms etc. Failing to pay VAT, General Authority of Zakat & Tax of Saudia Arabia may like to impose appropriate penalties on taxpayers for violating VAT rules as set forth by the Law or implementing regulations. However, there are difficulties due to the manual implementation of the tax collection procedure which cannot be neglected [10].

Generally tax invoices are proof of tax collection. A tax invoice is a proof that the seller has collected the VAT from the buyer. A VAT previously paid by a buyer of good is partly recovered from the next buyer of the same good and services. During the VAT reporting period, the difference is calculated among tax invoice paid and tax invoice received. Only the taxable person for VAT purposes can generate tax invoices [11].

### A. VAT Mechanism as Adopted by Saudi Arabia

There are 135,906 VAT registered businesses in the Saudi Arabia. Upon purchase of taxable service or goods, buyer pays VAT to the seller along with price of the service or good. This information is reflected in the invoice issued by the seller. The buyer can then use this invoice to show he/she is an active tax payer and further to reduce her overall tax where required.

At the point when a VAT-enlisted business offers a decent or benefit, it charges – accepting a standard case – an additional

5% of VAT over the business cost. The business will apply that 5% to every single qualified sales independently from its income with the end goal to later dispatch a bit of it to the administrative authority or government. The VAT that a business gathers on its deals is called Output VAT. The VAT that a business pays back to its providers is called Input VAT. With the end goal to ascertain the amount they owe to GAZT, every business will take note of the amount of VAT it has gathered from clients and subtract from it the aggregate VAT it paid in a similar period. A detailed process of VAT collection in KSA is described in Fig. 1 while an example is also provided in Table I.
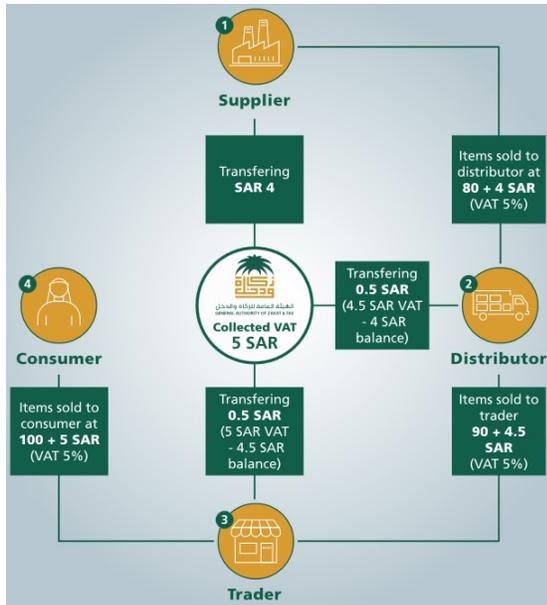


Fig. 1. VAT collection mechanism in KSA

### B. Centralized vs. Decentralized Businesses

Famous examples of centralized business-to-business and customer-to-business e-commerce systems include ebay, Amazon [12], [13], etc. However, to bring trust into these systems, third party authorities are required. Due to these trusted entities customers make product selection, provide credit card information, make checkouts and payments so all tasks are indeed centralized and all the actions are initiated by users keeping in view trust components. Similar is the case with banks that are in fact centralized and trusted entities which are in turn authorized by some other entities till it reaches the Government bodies. In short, all types of businesses that includes business-to-business, business-to-customers, customer-to-business and customer-to-customer are indeed centralized systems and involve trust parties. So, a chain-of-trust is developed which originates in some place and ends in other that leads to no direct trust involved with clients [14]. Further, it includes the bottleneck on the central entity.

We have on the other hand decentralized systems like the Blockchain, that is in fact an evolution business oriented decentralized and heterogeneous solution involving a chain-of-trust among various peer entities. It has decentralized trust in a manner that all communication is allowed in this paradigm while there is no need of a third party, i.e. it ensures that communication among the parties or nodes of the Blockchain is trusted [15], [16], [17]. In fact, it stores all the communication among the two or more untrusted parties on an untampered ledger. This ledger is verified by many peer nodes before storing it. This simple mechanism gives trust to all parties involved in the business transaction that nobody can forge with the information.

Cryptocurrencies [18], for example, bitcoin [19], etherium [20] infrastructure that provides peer-to-peer business model. All are decentralized environments and involves trust among nodes without the need for a third party. It is also important to mention here that decentralized entities, such as, Blockchain has been tested with many models [21]. Various models have been presented to implement blockchain technologies. However, ought most priority need to be paid in opting a model for the businesses that is adopted and backed by various industry giants.

In [22] provides an implementation of Blockchain that has features of solving consensus and can manage and resolve disputes among various entities. However, it could not be adopted by large industrial players as those implementations were too specific. The author selected a model implementation of Blockchain that is adopted by the industry, the Hyperledger Fabric. Majority of banks, Intel and Google, etc. have support in these umbrella projects.

Value Added Tax (VAT) is recently imposed by the government of Saudi Arabia for the first time. However, there is a need for an automated system that can keep track of millions of transactions generated from a supply-chain in the country for VAT. A traditional solution is to design and implement a centralized system, however, due to aforementioned reasons it is necessary to adopt new technologies which have more flexibility. In this regard, Blockchain is quite feasible for VAT implementation because most of the future application are going towards it. All the stakeholders in the supply-chain system will become part of Blockchain network and all transaction will go via the Blockchain system, setup by the government authorities. Smart-contract will be introduced to keep track of valid transaction and propagation on the Blockchain network.

The aim of conducting the research is to critically examine the existing VAT requirements and to provide implementation of model that can effectively collect VAT without the need of a third party and aligned with vision 2030 of Saudi Arabia while the objectives of studies are provided hereunder:

1) To Design Hyperledger Fabric and Composer infrastructure for supply chain management.
2) To design an end-to-end VAT implementation on a decentralized network.
3) To implement newly designed VAT system on Hyperledger composer and hyperledger fabric for implementation of VAT and setting up various peers and order nodes on Blockchain network.

Decentralized network and VAT are newly introduced in Saudi Arabia. For the implementation of the same authors adopted decentralized system. We opt Blockchain and within the Blockchain the author adopted a generalized system that

TABLE I.    EXPLANATION OF VAT PROCESS THROUGH PURCHASE OF COMPUTERS

| VAT (transactions) | Price of Item | Tax per item | Items bought or sold | Total Tax | Name of Tax |
|---|---|---|---|---|---|
| Computers bought from Manufacturer | SAR 1,000 | SAR 1,000 x 5% = SAR 50 | 100 | SAR 50 x 100 PCs = SAR 5,000 | Input-VAT |
| Computers sold to customer | SAR 1,200 | SAR 1,200 x 5% = SAR 60 | 100 | SAR 60 x 100 PCs = SAR 6,000 | Output-VAT |
| VAT owed to Govt | VAT due = Output VAT – Input VAT (SAR 6,000 – SAR 5,000 = SAR 1,000) | | | | |

is accepted by industry. Moreover, it is more execute efficient, more secure, solves the consensus and is a standardized operating system that have consensus pluggable architecture. Through consensus pluggable architecture, one can use any consensus algorithm.

### C. Reason for Opting Blockchain

Blockchain conveys solid ongoing data from numerous layers to a large stakeholders, just like the case with tax collection, particularly on a worldwide level. Blockchain is utilized because its distributed ledger which means that Blockchain offers transparency and all supply chain are registered on the system to record all transactions involved in exchange of a product or bulk of products from seller to purchaser. As a result, Blockchain accommodates and maintains all history of transactions performed and nothing can be excluded during performing any transaction. Because of using Blockchain, any organization can improve its renew as the digitized transactions are recorded on untampered ledgers which can never be altered. In case of any alteration, the chain of trust is corrupted that causes the whole network to be compromised. In case of corruption of the network, the peer nodes are not affected as they have their own network which includes ordering nodes. Upon generation of transactions, the ordered nodes computes hashes of the same. The hashes are stored and appended with previous transactions. The final hash generated is disseminated among the other nodes which upon receiving, calculates their own final hash based on their own records. The transactions are accepted or rejected based on the final hash values generated at the nodes. In case of matching of the final hash values, transaction is generally accepted otherwise, the new transaction is rejected.

Focus of any organization functioning in today's competitive marketplace is to gain and sustain competitive advantage. With the huge volumes of data stored in databases, data marts and data warehouses coupled with advanced data analysis tools, managers are now in a better position to make smart and effective decisions which result in competitive advantage for their organizations. Business Analytics (BA) is a new and upcoming area of advanced data analysis that has emerged as a significant area of study for both researchers as well as practitioners over the last two decades. BA is the process of transforming huge volumes of data into new knowledge through analysis and using that knowledge to for effective decision making and problem solving which ultimately results in value-creating competitive actions.

The proposed research work corresponds to the priority areas of Saudi Arabia. According to vision 2030 of Saudi Arabia, the government is investing on Blockchain. Keeping in view the importance and the features it offers, Blockchain is implemented in rest of the world. This research studies contributes in terms of utilizing Blockchain for VAT in Saudi Arabia so that all possible solutions are shifted to Blockchain infrastructure.

The author propose a solution for VAT and supply-chain interface that is built on standardized Hyperledger interface. The proposed VAT model is on execute-order-architecture. Specifically, Hyperledger composer is used to design VAT system and implemented the interface for all the involved parties as shall be explained in the upcoming sections. The

Section II provides background studies conducted regarding VAT collection and various employed techniques for its collection. Section III elaborates methodology of our own proposed solution and provides comprehensive implementation details. Moreover, Blockchain employment in various domains is explained in the section which is followed by details of various configuration and implementation related to employing Blockchain for VATA collection in the Kingdom of Saudi Arabia. Section V concludes the paper.

## II.  BACKGROUND

Along with its partners across the GCC, the Kingdom of Saudi Arabia has chosen to implement a standard VAT tax rate of 5%. This is one of the lowest rates in in the world. See the examples below of standard VAT rates in other countries. Most of the countries chooses to employ some sort tax rate within countries as depicted in Fig. 2.

As can be ascertained from vision 2030 of Saudi Arabia, Government is investing on Blockchain. This research studies contributes in terms of utilizing Blockchain for VAT in Saudi Arabia.

Modern technologies such as artificial intelligence, Internet of Things, Blockchain are among the key pillars of the development and progress of countries and serve as major contributor to the Industry 4.0 and the development of GDP. Encompassed in this vision is Saudi Arabia's ambition to
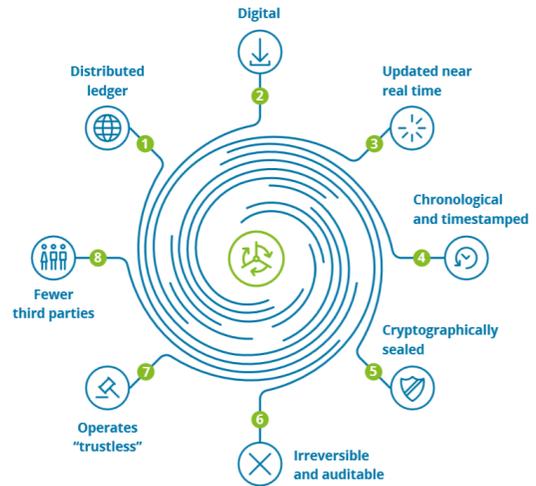
Fig. 2.    VAT rates in various countries



Fig. 3.    Blockchain Key Features [24]

deliver an improved e-government, drive digital transformation and attract foreign investment. With a focus to also improve the Kingdom's ease of doing business ranking, achieving these ambitions will require challenges to conventional practices. Implementing modern technologies such as artificial intelligence, Internet of Things, and Blockchain is core to enable this transformation. Furthermore, this path to digital transformation will play a fundamental role in shaping the way challenges such as competition, cost and budget pressures, changing citizen and resident demands, are undertaken. Blockchain, in particular, has the potential to transform government services while disrupting industries such as healthcare, education, banking, and real estate. The technology could reduce VAT Fraud and can reduce VAT Loss [23].

Blockchain offers transparency and all supply chain are registered on the system to record all transactions involved in exchange of a product from seller to purchaser. As a result, Blockchain accommodates and maintains all history of transactions performed and properly records any action/transaction that cannot be excluded or omitted. Because of using Blockchain, any organization can improve its renew as the digitized transactions are recorded on untampered ledgers which can never be altered. In case of any alteration, the chain of trust is corrupted making the whole network compromised. However, in case of corruption of the network, the peer nodes are not affected as they have their own network which includes ordering nodes. The ordering nodes upon generation of transactions compute hashes of the same. The hashes are stored and appended with previously generated transactions. The final hash generated is disseminated among the other nodes which upon receiving calculates their own final hash based on their on records. The transactions are accepted or rejected based on the final hash values generated at the nodes. In case of matching of the final hash values, transaction is generally accepted otherwise, the new transaction is rejected. Peer nodes perform evaluation, validation, they provide a fault tolerant environment in case of a node which is down. Because of the feature under discussion and as provided in the Fig. 3, we feel that the proposed infrastructure can improve the newly adopted VAT manual collection system in terms of generating revenue in Saudi Arabia.

### A. Analysis and Comparison Results of Employing Various Consensus Algorithms

The existing solutions for Value Added Tax (VAT) collection are manual and electronic methods. One of these is depicted in Fig. 4.
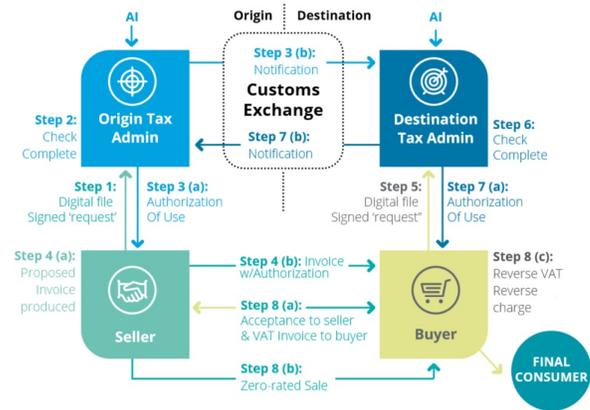


Fig. 4.    Electronic Tax Collections [24]

Researchers have provided solutions that are based on Blockchain [11], [25], [26], [22], [27], [28], [29], however, they are using mining as consensus algorithm and are not proven to be very efficient neither they are standardized. Mining based architecture are also known as order-execute-architecture. Wijaya et al. [11] implemented Blockchain through Proof of Work (PoW) as consensus mechanism for implementation and collection of VAT. However, the PoW is not very efficient i.e. transaction throughput is too slow (21 transactions per second).

Wijaya et al. [11] have proposed a convention permissioned based Blockchain that turns around the way toward dealing with fake tax invoices. Their strategy requires Tax Payable for VAT Purposes $(TPVP)$ shall not be allowed to make tax invoices solicitations to different TPVPs without first procuring legitimate tax credit i.e. A TPVP needs to pay some cash to get the tax credits, at that point they can make a duty

receipt as an approach to exchange the tax credits to another TPVP. By turning around the procedure, it is ensured that the tax invoices in the framework are portrayal of genuine assessment cash flowing through VAT framework. i.e. invoice is issued after checking credit and debit record of VAT through Blockchain. Thus, the paper based mechanism is converted to Blockchain based electronic VAT collection process.

The Blockchain provides an on-line environment for maintaining distributed database of records, ledgers and events that are executed and these belongs to a number of participating bodies. The public ledger and records gets updated and modified with the permission of majority of participants and is not subject to subsequent crash, modification or manipulation. The famous controversial bitcoin is based on Blockchain flawless technology. Blockchain provides a democratic place where it can further facilitate in provision of applications in the domains of finance and others which are briefly explained in the following subsections.

### B. The Trust Aspect

The current business activities are performed through "Trust" on third parties for ensuring transactions are performed flawlessly. However, it is possible that these third parties are hacked, misused or mismanaged. Such issues are better coped through Blockchain wherein at any point in time, the transactions carried out at as the assets are verifiable in future.

### C. Smart Contracts

Smart contract is another use case of Blockchain which are digital programs that keep conditional records of participants and upon occurring specific event, the contract elements are executed which may include payment to a party.

### D. Smart Property

Blockchains comes with tremendous opportunities. Use cases of Blockchain also includes "smart property" where the property can be physically or nonphysical. Examples of physical entities may include car, house, plot and things like that. Blockchain has been remained the most secure test environment [30].

### E. Health and the Banking Sector

Most of the banking sectors are exploring ways regarding how to migrate towards Blockchain technology in order to secure their transactional records. Health records, legal documentation, private securities and marriage licenses are maintained in the block-chain effectively. These assets are best protected through digital fingerprints of the assets instead the actual asset [31], [32], [33].

### F. Understanding Security Mechanism of Blockchain as Case Study of Bitcoin:

In Blockchain, security of transactions is ensured through public and private cryptographic keys i.e. transactions are protected through digital signatures. So in order to spend money, the owner of the cryptocurrency needs to prove his ownership of the private key Every transaction is broadcasted to the rest of nodes and then after verification the transaction is

recorded in the ledger that is public in nature. In any case, there is question of keeping up the request of these exchanges that are communicated to each other hub in the Bitcoin network. The Bitcoin tackled this issue by an instrument that is currently famously known as the Blockchain framwork. The Bitcoin framework orders exchanges by setting them in blocks and then connecting these in a chain. The exchanges in a single block are considered to have occurred in the meantime. These blocks are connected to one another (like a chain) in an appropriate straight, sequential manner with each block comprising the hash of previous one.

Any node can however contain unconfirmed transactions from which blocks can be created and broadcasted among other nodes. Bitcoin tackles this issue by presenting a scientific riddle: each block will be added in the Blockchain given that it contains a response to an extremely extraordinary numerical query. It takes ten minutes for a node in the Bitcoin system to make a correct speculation and produce a block. A little likelihood exists that more than one blocks are created. To cope this problem, it is checked which node solves the problem that node will be allowed to broadcast the block to all other nodes in the network.

The next section provides necessary details of methodology for the proposed solution of VAT implementation in Saudi Arabia through Blockchain.

### III. METHODOLOGY

The methodology for our proposed solution alongwith selected implementation details are provided in the following subsections:

### A. Selection of Blockchain Solution for VAT Collection

Researchers have provided solutions for businesses on peer-to-peer networks, however, those cannot be adopted as a standard for all. Similarly, they do not have support from industry. Generally, there is always support from industry for the acceptance of a product otherwise failure comes to such products in the long run. The fact can be ascertained by looking to Android mobile operating system which is supported by Google and it has achieved an acceptance and received most of the mobile market share. There were many other operating systems that could not get much popularity due to unavailability of support from industry big players. In this research, the author have opted Blockchain, a solution that can easily be adopted by industry and has quite big support for business applications. After selection of the Blockchain infrastructure, the author designed and implemented VAT system over the Blockchain for Saudi Arabia as a use-case .

### B. Proposed Solution and Selected Implementation Details

The proposed model is primarily designed to implement VAT system for the Kingdom of Saudi Arabia through state-of-the-Art Blockchain. Ideally, any supply chain management that relates to public and private organizations, service providers, consumers, large and small businesses will essentially be employed through Blockchain. For collection of VAT, the architectural framework is built upon a Hyperledger Fabric operating system, connecting a number of supply chain management organizations. The network comprises of a cluster of

servers that possess *orderer* nodes and certificate authorities (CAs), *peer* nodes. Saudi VAT collection mechanism is hosted on the cluster that manages and control the entire logical components including implementation of mechanism for consensus and certification authority. Likewise a group of interconnected servers implementing supply chain management are connected to the cluster through peer nodes as depicted in Fig. 5.
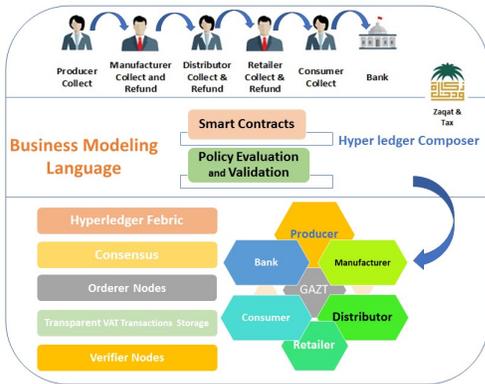


Fig. 5. Complete Proposed VAT System on Blockchain Infrastructure

The Hyperledger Fabric is the main structure or backbone that runs the composer over itself. In the Hyperledger Fabric, a network is composed consisting of cluster of servers for hosting Saudi VAT collection departments and in similar fashion, collection of servers hosting tax payable organizations which are interconnected with each others. The upcoming organization that intends to join the network or become part of the Blockchain are issued certificates by Certificate Authority (CA) of Saudi Arabia for secure communication among peers. These certificates are issued to organizations for a particular role and based on those roles, the transactions are controlled in the Blockchain.

After successful deployment of the Hyperledger Fabric, Hyperledger composer Business Network is deployed on each of the node and likewise issued connection profile for each of the relevant participant in individual organizations. Through these connection profiles, the composer-rest-server API creates URLs for different services including the creation of records or deletion of some records. The organizational personnel can trigger smart-contract transactions to deduct and send VAT amount from their system invoices to the distributed ledger. In-response, new automatic transactions are triggered that update VAT amount on-behalf of organization and the payee. In the next section, implementation related details regarding the proposed solution is provided.

## IV. IMPLEMENTATION

In this section, technical details are provided related to VAT implementation initiating with configuration of various services.

### A. Configuration of Various Services

Inside the fabric sample, the directory contents are downloaded and modified to increase number of joining organizations in the Count variable of the file namely, crypto-config.yaml. To configure the proposed model application's

services on each peer which includes CA, orderer, Peer, couchdb and CLI, the settings are re-configured in the docker-compose.yaml file. In this file, the certification key are added for the organizations that intends to join the network. The author also added extra-hosts section after the volume section in all of our services except Couchdb and CLI. After creation of separate docker-compose.yaml file for each of the other organization. These file will have the configuration of various services for the other organizations. Here for supply chain organization two services will be added, i.e. peern.example.com and couchDB. For deployment on other organization peers the crypto-config folder is copied, docker-composer-peern.yml, env folder and start-peern.sh. To start network on each peer command ./start-peern.sh is utilized.

With a running Hyperledger fabric network and multiple peer nodes, the business network is deployed using Hyperledger composer. To do so first Peer Admin Cards are created for managing the network. In this connection profile, the author added all pool of organizational IPs which are taking part in the network and updated service certificates for them. Further, added these settings in the connection profile in the section of JSON CONNECTION PROFILE SECTION.

To determine whether a particular transaction is valid, Validation System Chaicode (VSCC) is invoked that is accompanied by the transaction's chain code. When a peer receives a transaction, it invokes the VSCC associated with the transaction's Chaincode as part of the transaction validation flow to determine the validity of the transaction.

### B. Creation of Participation and Assigning Identities

In Composer a Participant is just a data item, specifically an object in a Participant Registry. A participant cannot access the Business Network on the Fabric until an Identity has been Issued to and bound to that Participant. Identities are generated by the CA which belongs to an Organization. A user (administrator) with an Identity can create Participants if they have the ACL access to do so, but only an Identity with specific rights in the CA can issue Identities.

### C. Design of Blockchain based VAT System

Fig. 6 outlines the Structure of Business Network Definition that runs on composer base architecture. It depicts a complete data model of the proposed architecture that includes Assets, Participants, and Transactions. All stackholders can trigger their transactions, and can explain how data will be managed.

For implementation of the proposed infrastructure, *Blockchain* environment is built. The state of art tool, the hyperledger fabric is opted to define business logic and configure orderer nodes. Peers have their own stakeholders which includes various supply chains like manufacturer, mediocre.Further, as per our feasibility study, the deployment of such a system is highly effective.

For the purpose, a business-oriented modeling language is opted as the hyperledger fabric interface is so complex for writing applications. The author opted hyperledger composer that is a step further and application oriented. Through this, applications can be modeled more appropriately. The VAT
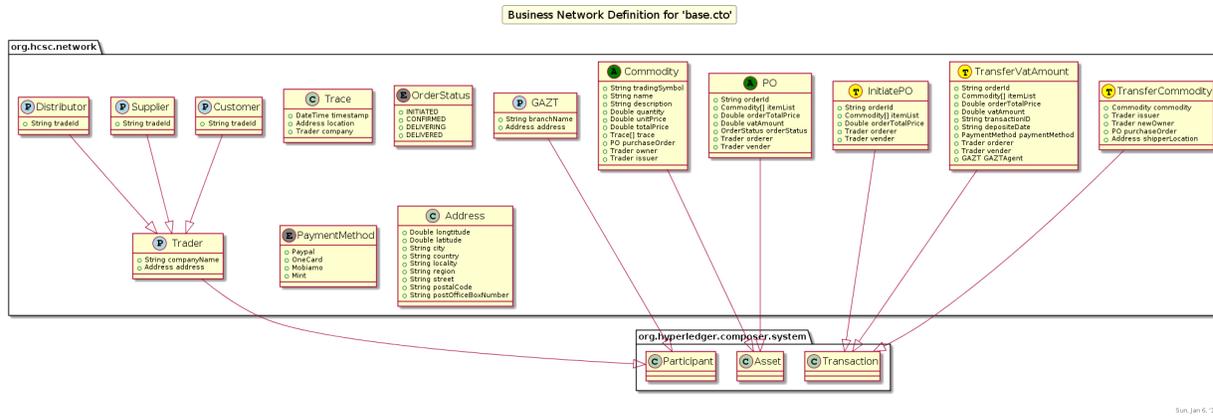
Fig. 6.   Class diagram of Business Network Archive

system is designed through hyperledger composer which shall define all components and how they will operate on the modeling language as depicted in Fig. 7.

### D. Deployment of Hyperledger Composer and Peers

The proposed model is primarily designed to implement Value Added Tax system of the Kingdom of Saudi Arabia through state-of-the-Art Blockchain. Ideally, any supply chain management that relates to public and private organizations, service providers, consumers, large and small businesses and industries will essentially be employed through Blockchain. For collection of VAT, our framework is built upon a Hyperledger Fabric network, which connects all the supply change management organizations in Saudi Arabia. A cluster of servers is framed that possess orderer nodes and certificate authorities (CAs). Saudi VAT collection mechanism is hosted on the cluster that shall essentially be a managing platform to adopt and implement the entire logical and other necessary details of the consensus mechanism and certification authority. Likewise a group of servers as depicted in Fig. 7, which are interconnected with each others implementing supply chain management which are connected to the cluster through peer nodes.

In the Hyperledger Fabric, a network is composed consisting of cluster of servers for hosting Saudi VAT collection departments and collection of servers hosting tax payable organizations which are interconnected with each others. The upcoming organization intended to join the network or become part of the Blockchain are issued certificates by CA of Saudi Arabia. These certificates are issued to organizations for a particular role and based on those roles, the transactions are controlled in the Blockchain.

### E. Hyperledger Business Composer Network

After successful deployment of the Hyperledger Fabric, Hyperledger composer Business Network is deployed on each of the node and likewise connection profiles are issued for each of the relevant participant in individual organizations. Through these connection profiles, the composer-rest-server API creates URLs for different services like creation and deletion of records. Organizational personnel can trigger smart-contract transactions to deduct and send VAT amount from

their system invoices to the distributed ledger. In-response new automatic transactions are triggered that update VAT amount on-behalf of organization and the payee.

Inside the fabric sample downloaded from the directory contents are modified to increase number of joining organizations in the Count variable of the file namely, crypto-config.yaml. To configure our proposed model application's services on each peer which includes CA, orderer, Peer, couchdb and cli, the settings are re-configured in the docker-compose.yaml file. In this file the certification key are added for the organizations that intends to join the network. Extra-hosts section are also added after the volume section in all of our services except Couchdb and CLI. Further, docker-compose.yaml creates separate file for each of the other organization. These file will have the configuration of our services for the other organizations. Here for supply chain organization, only two services are added which are peern.example.com and couchDB. For deployment on other organization peers the folder of crypto-config is copied, docker-composer-peern.yml, env folder and start-peern.sh. To initiate the network on each peer, the command ./start-peern.sh executes the network.

### F. Deploying Business Network through Hyperledger Composer

With a running Hyperledger fabric network and multiple peer nodes, the business network is deployed using Hyperledger composer. To do so. first Peer Admin Card is created that manages the network. In this connection profile, authors have added all pool of organization IPs which are taking part in the network and update our services certificates. Further, added these settings in the connection profile in JSON connection profile section.

## V. Conclusion

Tax-returns and settlements are ascertained over a fixed period, for instance over monthly or quarterly basis and the counts are not based on exchanges, yet rather on self-assertive dates for instance receipt dates. The process is troublesome for governments, if certainly feasible, to track VAT installments. The computerized age is additionally forming tax collection systems into a totally extraordinary shape, by not just changing the connection among citizens and tax authorities, yet in
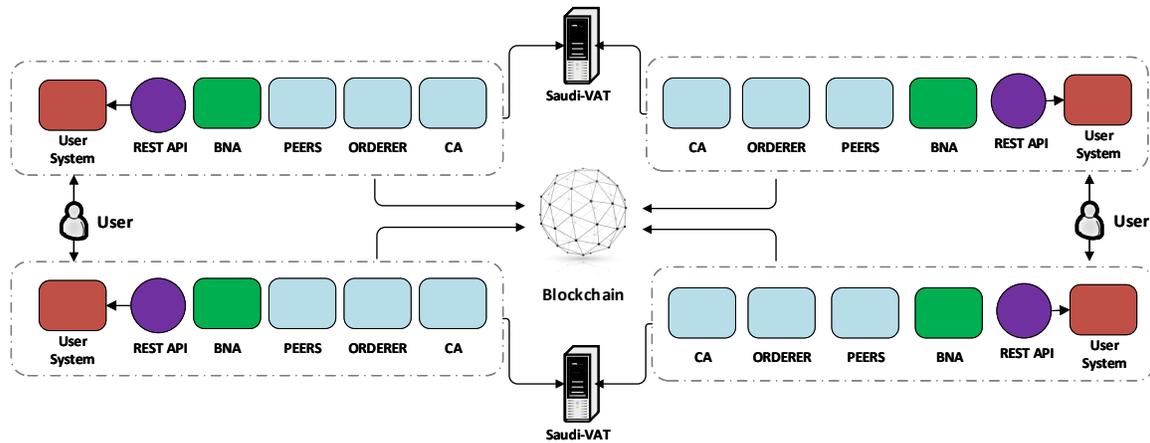
Fig. 7. Proposed Implemented Architecture for Value Added Tax Collection in KSA

addition modifying the manner in which government taxes are covered, submitted and stored. The capability of digitizing VAT has been seen by numerous nations, and new arrangements emerge including SAF-T in Europe or real time invoicing mechanism in South America and Brazil. This search studies provides details of various methods applied for VAT collection through Blockchain and analyses the presented solutions. The studies contributes by presenting a more effective and efficient VAT collection solution through Blockchain. The author, provided a proof-of-concept implementation on Hyperledger composer which is an execute-order architecture to satiate efficient transaction processor over Blockchain.

## VI. Future Work

Blockchain-based VAT system is a step forward in transparent and secure Tax collection systems. The solution so far presented in the literature are based on Proof-of-Work consensus model that only works with Bitcoin architecture. However, in the proposed architecture VAT is designed and implemented on PBFT (Practical Byzantine Fault Tolerance) consensus model. It is recommended that other consensus model should be enabled on Hyperledger Fabric operating system and run the proposed architecture. That will give some more interesting results on performance, security and storage model of VAT based transactions.

## Acknowledgment

## References

[1] A. Schenk, V. Thuronyi, and W. Cui, *Value added tax.* Cambridge University Press, 2015.

[2] J. Olsen, C. Kogler, J. Stark, and E. Kirchler, "Income tax versus value added tax: A mixed-methods comparison of social representations," *Journal of Tax Administration*, vol. 3, no. 2, pp. 87–107, 2017.

[3] M. Cooper and M. Knittel, "Partial loss refundability: How are corporate tax losses used?" *National Tax Journal*, pp. 651–663, 2006.

[4] D. Butler *et al.*, "Superannuation: Managing tax losses in an smsf," *Taxation in Australia*, vol. 49, no. 11, p. 699, 2015.

[5] L. Barbone, M. Bonch-Osmolovskiy, and G. Poniatowski, "Study to quantify and analyse the vat gap in the eu member states," 2015.

[6] K. Pavlov, "Can the general reverse charge mechanism combat missing trader fraud and provide for secure vat collection?" 2017.

[7] F. Times, "Saudi arabia and uae introduce 5[Online]. Available: https://www.ft.com/content/b1742920-efd0-11e7-b220-857e26d1aca4

[8] GAZT, "HealthCare Guideline, version 1, February 2018 by Kingdom of Saudi Arabia," 2018, available at: https://www.vat.gov.sa/sites/default/files/2018-03/VAT%20%20Healthcare%20Guide%20English.pdf.

[9] GAZTAX, "Imports and Exports Guideline, September 18, 1st Edition by Kingdom of Saudi Arabia," 2018, available at: https://www.vat.gov.sa/sites/default/files/2018-09/VAT_import_Export_Guideline_English.pdf.

[10] I. Bostan, C. Popescu, C. Istrate, I.-B. Robu, and I. Hurjui, "The impact of taxation of the domestic economic transactions on the vat collection through electronic fiscal devices," *Amfiteatru Economic*, vol. 19, no. 45, p. 581, 2017.

[11] D. A. Wijaya, J. K. Liu, D. A. Suwarsono, and P. Zhang, "A new blockchain-based value added tax system," in *International Conference on Provable Security.* Springer, 2017, pp. 471–486.

[12] P. Resnick, R. Zeckhauser, J. Swanson, and K. Lockwood, "The value of reputation on ebay: A controlled experiment," *Experimental economics*, vol. 9, no. 2, pp. 79–101, 2006.

[13] G. Linden, B. Smith, and J. York, "Amazon. com recommendations: Item-to-item collaborative filtering," *IEEE Internet computing*, no. 1, pp. 76–80, 2003.

[14] M. A. Morid and M. Shajari, "An enhanced e-commerce trust model for community based centralized systems," *Electronic Commerce Research*, vol. 12, no. 4, pp. 409–427, 2012.

[15] J. Yli-Huumo, D. Ko, S. Choi, S. Park, and K. Smolander, "Where is current research on blockchain technology?—a systematic review," *PloS one*, vol. 11, no. 10, p. e0163477, 2016.

[16] G. W. Peters and E. Panayi, "Understanding modern banking ledgers through blockchain technologies: Future of transaction processing and smart contracts on the internet of money," in *Banking Beyond Banks and Money.* Springer, 2016, pp. 239–278.

[17] X. Xu, C. Pautasso, L. Zhu, V. Gramoli, A. Ponomarev, A. B. Tran, and S. Chen, "The blockchain as a software connector," in *2016 13th Working IEEE/IFIP Conference on Software Architecture (WICSA).* IEEE, 2016, pp. 182–191.

[18] P. Walker and P. J. Venables, "Cryptographic currency for securities settlement," Jul. 11 2017, uS Patent 9,704,143.

[19] S. Nakamoto, "Bitcoin: A peer-to-peer electronic cash system," 2008.

[20] C. H. Lee and K.-H. Kim, "Implementation of iot system using block chain with authentication and data protection," in *2018 International Conference on Information Networking (ICOIN).* IEEE, 2018, pp. 936–940.

[21] A. Baliga, "Understanding blockchain consensus models," *Persistent*, 2017.

[22] R. T. Ainsworth and M. Alwohaibi, "Blockchain, bitcoin, and vat in the gcc: The missing trader example," 2017.

[23] V. 20-30, "Saudi Arabia's ambition to deliver an improved e-government," available at: https://ameinfo.com/money/economy/saudi-3/.

[24] S. Underwood, "Blockchain beyond bitcoin," *Communications of the ACM*, vol. 59, no. 11, pp. 15–17, 2016.

[25] R. T. Ainsworth and A. Shact, "Blockchain (distributed ledger technology) solves vat fraud," 2016.

[26] M. Swan, *Blockchain: Blueprint for a new economy*. " O'Reilly Media, Inc.", 2015.

[27] R. T. Ainsworth and M. Alwohaibi, "The first real-time blockchain vat-gcc solves mtic fraud," 2017.

[28] P. van der Zwan, "Cryptocurrencies & vat," *Tax Professional*, vol. 2018, no. 32, pp. 12–14, 2018.

[29] M. R. Hoffman, "Can blockchains and linked data advance taxation," in *Companion of the The Web Conference 2018 on The Web Conference 2018*. International World Wide Web Conferences Steering Committee, 2018, pp. 1179–1182.

[30] J. MICHAEL, A. COHN, and J. R. BUTCHER, "Blockchain technology," *The Journal*, 2018.

[31] X. Yue, H. Wang, D. Jin, M. Li, and W. Jiang, "Healthcare data gateways: found healthcare intelligence on blockchain with novel privacy risk control," *Journal of medical systems*, vol. 40, no. 10, p. 218, 2016.

[32] M. Mettler, "Blockchain technology in healthcare: The revolution starts here," in *e-Health Networking, Applications and Services (Healthcom), 2016 IEEE 18th International Conference on*. IEEE, 2016, pp. 1–3.

[33] N. J. Witchey, "Healthcare transaction validation via blockchain proof-of-work, systems and methods," Nov. 19 2015, uS Patent App. 14/711,740.

# Navigation Application with Safety Features

Andrew Usama[1], Moustafa Waly[2], Habiba Elwany[3],
Monica Medhat[5], Youssef Mobarak[6]
Misr International University
Computer Sciences
Cairo, Egypt

Mohamed H. ElGazzar[4]
IoT Senior Solutions Architect
Vodafone International Services
Cairo, Egypt

*Abstract*—In 2017, the number of car accidents that occurred was astronomically high, even though, infrastructural road systems are being continuously built and renewed to make it more efficient. But a significant problem which still remains is that a staggering number of accidents is exactly what should be avoided. In order to address this issue, this paper will serve to survey and discuss some of the solutions proposed, both software and hardware, for this problem. This will include some of the implemented safety features while also exploring how to make the system more interactive and smooth to meet user needs.

*Keywords—Traffic safety; navigation systems; neural networks*

## I. Introduction

In this survey paper, we gathered several papers with the specific purpose and aim of building a navigation application with safety features and how to implement them while also keeping the application interactive and user-friendly. Safety features include how to detect road junctions and how to identify straight road network sections using software such as famous machine learning algorithms or hardware such as the Velodyne rangefinder. This paper further aims to discuss several algorithms with their advantages and disadvantages, in order to reach a consensus on which algorithms are best suited for the tasks presented. This paper is organized in the following order Abstract, Introduction, Related Work, Conclusion, and lastly the references.

## II. Related Works

### A. Efficient Filtering and Clustering Mechanism for Google Maps

The work presented in [1] aims to solve the problem of displaying a large number of markers on the map without greatly affecting the interactivity and the reaction time for the application.

The algorithms used to solve this problem are Grid-based Filtering algorithm and DBSCAN, density-based spatial clustering of application with noise, algorithms [2], [3], and [4] which are used to be able to display large-scale geospatial data on the map. The first algorithm used is Grid-Based filtering which is used to prevent markers from overlapping when displaying hundreds or even thousands of them. The Algorithm divides the viewport into several grid cells for effective filtering. After that, it will get the top K places from each grid cell and performs the algorithm. Then, we give each cell a unique reference key given using Cartesian coordinate systems which save computation time when defining the reference from

each point's longitude and latitude. And then select the top k ranked places in each of these grid squares and display them to avoid overlapping markers. But even then, sometimes we can still get several markers that overlap which makes it difficult for the user to pick a marker if it's being overlapped and in that's where DBSCAN is used. DBSCAN is known for good performance. It is suited for tasks relating to large datasets, with noise, and can identify clusters with different sizes and shapes. It works by clustering regions with a high density of points. This algorithm needs three input parameters: For an experiment, a dataset of 74000 restaurant locations was gathered to test the algorithms proposed and then divided into hotspots for clearer results. One of the hotspots tested, was around the Busan area and the result are shown in Table I.

TABLE I. Results of the Busan area experiment [1]

|  | All | Naive Top-k | Proposed Algorithm |
|---|---|---|---|
| #markers | 370 | 125 | 77 |
| #clusters | - | - | 6 |
| Distance | 5 km | 5 km | 5 km |
| Zoom level | 13 | 13 | 13 |

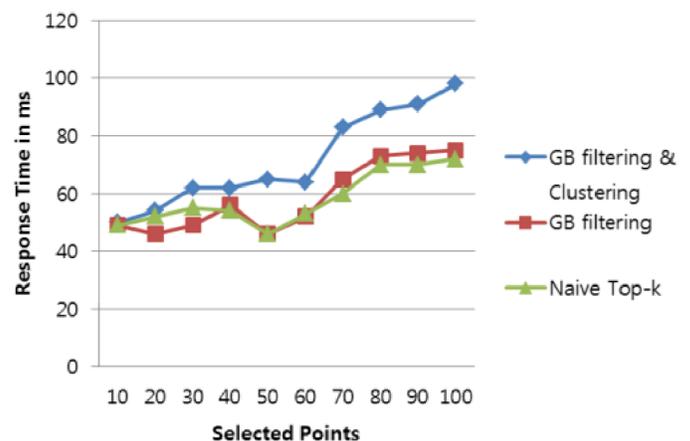Using these algorithms significantly improved the response time by 20ms as shown in Fig. 1.



Fig. 1. Response time for selected points [1]

The main drawback of this paper is that it does not try any algorithms other than the ones proposed. Instead, it compares a hybrid of Grid-based filtering and DBSCAN with naive plotting and while the results are significantly better; it is not a fair comparison since naive plotting is too basic to be compared to 2 algorithms.

### B. Cloud Aided Safety-based Route Planning

How to use the work described in [5] in order to create a neural network and use clustering techniques to produce a reliable risk assessment for safe route planning. While achieving the goal of guiding the driver toward the route with a low accident risk.

Developing a safety-based multi-objective safety route planning is done through an assessment a road risk module, which is proposed through features like road risk index, crash prediction, neural network, hybrid neural network model, and Crash rate history in order to model the risk level. Furthermore, another approach for road risk index emerged in DE Leur's and Sayed's work [6] which studied the driver assessment of existing road risks. A function for exposure of crash rates is often modeled as road risks, in which the amount of crashes occurrence are the exposure representation. Finally, the crash rate is the number of crashes per unit exposure and the most common measurement of exposure is the miles a vehicle has traveled.

Using these factors which were gathered from the road segment and crash data, a crash rate prediction was modeled. This model has a multi-state database that includes accident data, roadway inventory files, state supplemental inventory containing curves, and grade data. Furthermore, 70 percent of the samples in the dataset was used for training while 15 percent were reserved for the testing and validation of the neural network. The mean square error is the measurement of the quality of the module and for the [7] R-value regression is the relation between the expected output and the actual output. The mean square error used was 9.21 and the R-value regression was 0.80. Consequently, a hybrid neural network model was developed that aims to improve neural network models. This model was developed by partitioning the raw input-output data into three clusters using the fuzzy C-means clustering algorithm. For each cluster, we get a neural network model as shown in Fig. 2. Finally, the accident rates can be
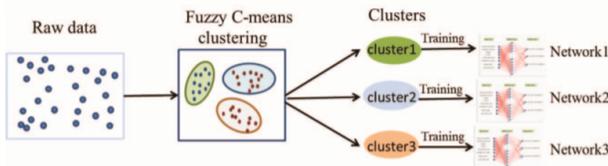


Fig. 2. A hybrid neural network model [5]

predicted, and the risk index is generated using the following formula:

$$RRI(i) = F(\sum_{i=1,2,3} n_{ij} S_j / AADT_i) \qquad (1)$$

,where:

- i is the road segment number
- RRI(i) is the risk index of road segment i
- $n_{ij}$ is the predicted number of accidents
- $S_j$ is the cost of an accident of type
- The $AADT_i$ represents the annual average daily traffic on the road segment

A real-world route experiment, shown in Fig. 3, was performed which is a route planning from Scioto Downs Inc to Delaware Ohio by first abstracting into graph the road network. The main road intersection was represented as nodes. The aim is to test the models and see the result from the first node to the final node with the minimum cost in the time expected to travel and the risk index. The output was computed by Cplex for $\alpha$ = 0 and $\alpha$ = 0.2, where $\alpha$ is the weight on cumulative RRI reflecting the driver's safety preferences. for $\alpha$ = 0, for the travel time expected was 42 minutes and the total risk index was 161.86. While when $\alpha$ = 0.2, the travel time expected was 44 minutes and total risk index was 103.57. The second
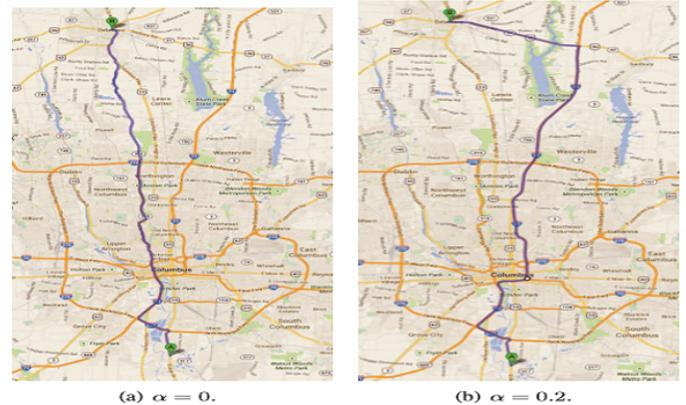


(a) $\alpha = 0$.  (b) $\alpha = 0.2$.

Fig. 3. Optimal routes with different $\alpha$ s [1]

route had a 36 percent higher risk when compared to the first route but required 2 more minutes which shows that the safe route often differs from the fast route. Finally, the deficiency of modeling a road risk index using crash rate cannot be used for roads with no historical crash data. And it cannot predict the dependence on changing factors like traffic density or weather condition.

### C. Identification of Curves and Straight Sections on Road Network from Digital Vector Data

According to [8], a method was proposed for the automatic identification of curves and their characteristics, particularly the horizontal size of the curves and their length on roads. According to the World Road Association (WRA), accidents tend to occur 1.5 to 4 times on curves than on straight roads. The method proposed was built on the principle of determination of the radius of an osculating circle [9].

The problem is that there is no ideal data available, due to the fact that breakpoints are somewhat equally apart from one another. The data is neither "dense" nor "sparse". And there are only very small measurement errors, but that type of data is not available. Therefore, the circumscribed circle calculation

TABLE II. The values in percentage indicate how many break points were correctly classified [8]

| Threshold value of radius (m) | II/446 | | II/444 | |
| --- | --- | --- | --- | --- |
| | Circumscribed circle | Osculating circle | Circumscribed circle | Osculating circle |
| 100% | 65% | 65% | 54% | 56% |
| 300% | 73% | 72% | 74% | 76% |
| 500% | 74% | 73% | 80% | 80% |
| 1000% | 70% | 78% | 80% | 81% |
| 1500% | 70% | 74% | 78% | 80% |
| 2000% | 72% | 75% | 77% | 78% |

of the radius may lead to the incorrect determination of curves due to the relatively small radii assignment. That's why a method called the osculating circle was used, it has none of the above drawbacks because it is merged with the Douglas-Peucker algorithm.

They compared their method with 3 other methods which are based on the identification of the circle circumscribed radius of the road segment, or on the calculation of the break-points angle formed by a road line. In this paper, data from the Road and Motorway Directorate of the Czech Republic (RMD CR) was used because they are directly acquired by the roads administrator. The authors collect the maps on which the curves exist, and then adds a geometric line generalization which reduces and simplifies the data without losing important information.

Then, the Douglas-Peucker algorithm was used which connects the curve's start and end points. Then, finds a point formed by all the connected points with the biggest distance from the line segments. While also considering the Euclidean distance formula and then comparing the distance retrieved with the selected and acceptable tolerance. And, if the accuracy tolerance is more than the distance, terminate the calculation if not, then it finds a simplified curve with fewer points. And then pass the result to an osculating circle calculation that detects the curves to approximate the line curvature at a specific point. The change of the tangent to the curve slope directly changes the radius.

Therefore, by drawing a circle on each point and running the equation using the three consecutive points of a circumscribing circle radius and the cumulative angle of three points; they are able to identify the curve.

Using the osculating circle, the authors had a positive predictive value of 95.9 percent which allows them to identify the curves correctly. The method is based on calculating the radius of the osculating circle and comparing it with three methods alternatively. But, the methods of calculating the radius of the curve significantly appeared better than those based on the lines of the break angle calculation. Table II shows the success rate of the given classified values for the boundary of the radii.

On the other hand, the determination of margins for identifying curves was left for further research hence it was not properly addressed in this work.

### D. Detecting Road Junctions by Artificial Neural Networks

Nowadays, Road junction detection has become an important and essential factor of navigation systems while being very difficult to detect. And, most of the early approaches were not successful. This paper presents a method for road junction detection using medium to a high-level resolution aerial images. This method should be able to deal equally well with three and four arm junctions based on raster and vector information. In this paper, image processing algorithms such as edge and line detection were used.

In [10], artificial neural networks (ANN) were used for object detection. Because neural networks have proven their universality in several technical fields such as image processing. They are used to recognize signalized control points in photogrammetric images and to detect airplanes in images as mentioned in many papers that used ANN. But the results were unsatisfactory so the Lanser algorithm [11] was added which is a rotation invariant version of the Deriche edge finder for edge detection. Then, the Ramer algorithm was used for edge smoothing and a circle centered on the junctions was used. Moreover, the Levenberg-Marquard algorithm was used because it's known for its high efficiency. Then, ambiguous image samples were removed due to wrong lighting condition issues, trees, and buildings causing the shadowing of the image which produced clean data. Then, the neural network was trained using that data from the images. Afterward, the ANN was left to detect the junctions. The detection of the junctions was accomplished by passing a predefined size window over the image. The result of the detection was either it's a junction or a non-junction. The window has a center pixel that the image stores the calculated radius of the circle drawn in it.



Fig. 4. Detected junctions in the other part of the study area [10]

An ANN classifier was used and a feed-forward neural network was chosen because it is frequently implemented in many programs. The testing was done using 0.4 M. resolution, black-and-white orthoimages which covered a region near Germany. The result was that the program was able to detect several potential junctions. All recognized junctions were marked with circles as shown in Fig. 4 and 5 and had their radius computed and any additional potential crossings

were marked with squares. It also detected bridges because they resemble crossing road segments. Finally, the advantages of this paper are that they made it easy to detect road junctions using ANN along with multiple other algorithms.
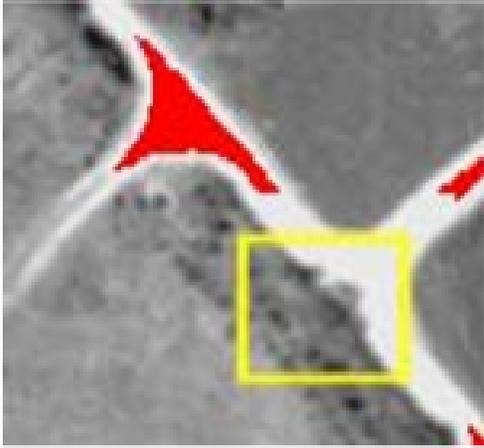


Fig. 5. Enlargement of figure 4 with tree shadows [10]

The disadvantages of this paper are that the ANN was trained on clean data to detect the junctions, which is inapplicable in real life. Additionally, the number of training samples was small which means that the result might not have been accurate. Although a larger training sample with more variety could have been used since one of the main advantages of ANN is its ability to handle large datasets. Consequently, it detected bridges as road junctions.

### E. Road Junction Detection from 3D Point Clouds

In this paper [12], the main focus falls on how to detect the changing traffic conditions because they're of vital importance when it comes to the safety of autonomous cars that navigate urban environments. One of the most critical traffic problems that require a lot of attention is road junctions. This work presents several machine learning techniques in order to detect junctions including Support vector machines, Adaptive boosting, and Artificial neural networks. During operation, 3D range finders create frames, which are point clouds created at regular intervals and each of those frames contains a set of environmental readings from a specific and singular location. The Objective algorithm's main objective is to classify any single frame into either a road junction or a road. Three phases are required to accomplish this task which are: feature extraction, initial classification, and finally Structured Classification. A Laser rangefinder collects three-dimensional point cloud environmental data. Features are, then, extracted from the point cloud using the method described in [13]. This method was chosen because it is rotation invariant since cars might change directions at junctions. Then, the features are passed on to the second phase: classification which as mentioned before includes the ANN, SVM, and AdaBoost

TABLE III. TEST RESULTS [12]

| Dataset | Classifier | Accuracy | Precision | Recall |
|---------|-----------|----------|-----------|--------|
| SC2 | ANN | 0.7088 | **0.870** | 0.8100 |
| | AdaBoost | **0.8700** | 0.6478 | **0.8466** |
| | SVM | 0.7220 | 0.7854 | 0.4987 |
| K2 | ANN | 0.8049 | 0.6606 | 0.6867 |
| | AdaBoost | **0.9027** | **0.8736** | **0.8139** |
| | SVM | 0.3626 | 0.8267 | 0.3009 |

to see which one will yield the best results. And then the third phase comes in which is Structured Classification

Artificial neural networks [14] are a nonlinear statistical machine learning method that is mostly used in pattern recognition. An ANN's objective is to adjust all the edges' weights and all the weights within each neuron, the weights of the previous layers are added and then they are put through an activation function like Sigmoid functions which restricts the output to either 0 or 1.

Adaptive Boosting [15] belongs to a family of methods called boosting that also trains boosted classifiers. Adaptive boosting uses dataset D for training, each iteration of t adds an extra layer to a t-1 layer classified $f_{t-1}$. These new layers that are constructed by Adaboost are each assigned a weight At to a possible weak classifier's hypothesis h, so that the error Et is minimized:

$$E_t = \sum_{x \in D} E[F_{t-1}(\mathbf{X}) + \alpha_t h(\mathbf{X})] \qquad [17] \qquad (2)$$

The aforementioned process neglects any classifiers that don't improve the predictive percentage of the model which also helps avoid overfitting.

Support Vector Machine [16] is a supervised learning model that can be used for classification. Any data given to an SVM is represented in the form of points in a space. SVM then attempts to fit a hyperplane through the input space. And the hyperplane that leaves the largest distance between different classes is the best fit.

The result of using the aforementioned three classifiers which are ANN, AdaBoost, and SVM yielded the results shown in Table III.

The next step is the structured classification in which the sensor produces 10 frames per second so that a number of consecutive frames belonging to the same class can be taken. This is done so that the classifier is given noisy predictions. In the experiments and results phase, 2 separate tests were run, the first using Velodyne 32 laser collected by members of the Carina 2 project and the second using a Velodyne 64, which is a more advanced version of Velodyne 32, which generates denser point clouds with twice as many per frame.

The tests were done to compare the performance of the three base classifiers' used, each algorithm was trained and tested separately and the accuracy, precision, and recall were calculated.

TABLE IV. TEST RESULTS WHEN USING CROSSROAD DETECTION THROUGH CURB DATA AND CROSSROAD DETECTION THROUGH ROAD SURFACE DATA [17]

(A) 404-2

| | | Predicted | |
|---|---|---|---|
| | | CR | NCR |
| Actual | CR | 340 | 22 |
| | NCR | 27 | 461 |
| Accuracy | | 94.24 % | |
| MSE | | 0.11517 | |

(B) 404-5-2

| | | Predicted | |
|---|---|---|---|
| | | CR | NCR |
| Actual | CR | 327 | 35 |
| | NCR | 25 | 463 |
| Accuracy | | 92.94 % | |
| MSE | | 0.13592 | |

(C) 404-10-2

| | | Predicted | |
|---|---|---|---|
| | | CR | NCR |
| Actual | CR | 345 | 17 |
| | NCR | 33 | 455 |
| Accuracy | | 94.12 % | |
| MSE | | 0.11853 | |

(D) 210-2

| | | Predicted | |
|---|---|---|---|
| | | CR | NCR |
| Actual | CR | 334 | 4 |
| | NCR | 1 | 112 |
| Accuracy | | 98.67 % | |
| MSE | | 0.01941 | |

(E) 210-5-2

| | | Predicted | |
|---|---|---|---|
| | | CR | NCR |
| Actual | CR | 328 | 10 |
| | NCR | 2 | 112 |
| Accuracy | | 97.35 % | |
| MSE | | 0.05046 | |

(F) 210-10-2

| | | Predicted | |
|---|---|---|---|
| | | CR | NCR |
| Actual | CR | 332 | 6 |
| | NCR | 1 | 113 |
| Accuracy | | 98.45 % | |
| MSE | | 0.03331 | |

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \qquad [12] \quad (3)$$

$$Precision = \frac{TP}{TP + FP} \qquad (4)$$

$$Recall = \frac{TP}{TP + FN} \qquad (5)$$

Where TP are the true positives, TN are the true negatives, FP are the false positives, and FP are the false negatives.

The Drawbacks of this paper is that only 3 specific classifiers were tested and not more which gave no room for the possibility of other algorithms that might have performed better.

### F. Crossroad Detection using Artificial Neural Networks

According to [17], Autonomous ground mobiles may be the key to increasing traffic flow and highway capacity because they're capable of optimizing cars' navigation in the streets. Another one of their benefits is their ability to avoid collisions and therefore reducing road accidents. Autonomous vehicles must have a steadfast environment perception system to rely on, in order to safely navigate the streets and avoid collisions. This paper proposed an approach that takes advantage of our autonomous car's perception components to identify crossroads using a neural network. In order to detect curbs, obstacles are detected by analyzing the scan ring compression. Different obstacles have and create different compression patterns. For example, curbs make the sensor return distances between rings that are smaller than the ones flat terrain generates.

$$\begin{aligned}
\triangle r_i^{plane} &- r_{i-1}^{plane} \\
&= \frac{h}{\tan \theta_i} - \frac{h}{\tan \theta_{i-1}} \\
&= h((\tan \theta_i)^{-1} - (\tan \theta_{i-1})^{-1})
\end{aligned} \qquad (6)$$

$$I_i = [\gamma \triangle r_i^{plane}, \delta \triangle r_i^{plane}] \qquad [18] \quad (7)$$

The output of both the curb and the road surface detector were adapted to feed the neural network. And, in order to fix the number of points that will be inputted in the network, the points obtained from both detectors were placed in a grid comprising of cells that had the same dimensions. The first step of building the aforementioned grid is to store the number of points that fall in, in each cell. Then, this number is divided by the maximum number of points observed in the grid's cells Which results in a proportionate grid and limits any cell's value into a [0,1] range. So that cells with fewer points have lower values while cells with a larger amount of points have higher values. Training the aforementioned neural network for the classification of crossroads was performed in two ways: The first was by using road curb examples and the second was by using road surface examples. Each data set was separately trained in three different network topologies which are x-2, x-5-2, and x-10-2, where the input layer size x was 210 for road surface data and 404 for curb data. Then the dataset was split into 2 parts: training and testing. Training was assigned 66.6% of the dataset while testing was assigned the remaining 33.3% . And after testing the result were as follows:

Through the "Receiver operating characteristic", the crossroad classification performance comparison was conducted between curb and road surface data. 0.940589 was the AUC obtained for curb while 0.981209 was the AUC obtained for road surface data. In this method, it was deduced that using

the data obtained from the road surface was more appropriate in the verification of a crossroad's presence in the street as shown by the results in Table IV. And since the AUC was near 1.0, It was demonstrated that this classifier provided adequate true and false positives proportions. It was confirmed, by Analyzing distinct neural network topologies, that the single-layer perceptron brings higher accuracy and lower MSE to classify unseen patterns.

The drawbacks of this paper are that it only tested its data in a single-layer perceptron classifier. It, also, didn't test other classifiers; it only tested different data.

### G. Parallel Hyper-Heuristic Algorithm for Multi-Objective Route Planning in a Smart City

According to [18], the commercial navigation application for route planning only focuses on parameters such as distance or time while completely neglecting safety which is one of the most important parameters. The safety parameter consists of many things like crime rate which are one of the risk issues found around the world. Most people, tourists in particular, often don't have information about safe or dangerous areas which is why people need an application that gives them the safest and fastest route. So, first, we must know the crime risk index which is extracted through several features, some temporally related and others spatially related. Spatially related features include the number of police stations and their proximity to a given region which has a high correlation and impact on a region's crime risk. while the temporally related features include traffic flow which influences urban safety. After extracting the spatial and temporal features, the next step begins which is to formulate the road network as multi-objective planning with the objective of getting the route that has the lowest crime risk index. The paper proposed a Multi-Objective Hyper-Heuristic algorithm for planning the route and, then, the CUDA framework was added to increase the speed of the planning. Finally, we get a safety-aware routing that takes the crime risk into account as well as the distance and time. Many experiments were conducted to compare the proposed algorithm with other similar algorithms. The results showed that the proposed algorithm was 173 times faster than the EMLS algorithm, 5.3 times faster than RL-MOHH, and 3.1 faster than the parallel NSGAII.

## III. CONCLUSION

Surveys with different work and efforts yielded different results with some algorithms being noticeably better than others at accomplishing specific tasks than others. In several papers [5], [10], [17], ANN seems to be the best solution for detecting junctions and straight road network sections. while DBScan combined with Grid-based Filtering [1] seem to give the best result with regard to interactivity and an application's reaction time. As a Future work, it would be extremely beneficial to add an accident Heat-map based on previous accident rates going back a certain number of years. as well as blocking off certain areas that have been known to have a high crime rate based on previous incident reports that can be retrieved from public records. And finally, using the official weather reports to avoid certain roads that do not function well under specific weather conditions.

## REFERENCES

[1] G. K. Habibullayevich, X. Chen, and H. Shin, "Efficient filtering and clustering mechanism for google maps," *Journal of Advanced Management Science*, vol. 1, no. 1, pp. 107–111, 2013.

[2] M. Ester, H.-P. Kriegel, J. Sander, X. Xu *et al.*, "A density-based algorithm for discovering clusters in large spatial databases with noise." in *Kdd*, vol. 96, no. 34, 1996, pp. 226–231.

[3] C. Zhou, D. Frankowski, P. Ludford, S. Shekhar, and L. Terveen, "Discovering personally meaningful places: An interactive clustering approach," *ACM Transactions on Information Systems (TOIS)*, vol. 25, no. 3, p. 12, 2007.

[4] K. Mumtaz and K. Duraiswamy, "An analysis on density based clustering of multi dimensional spatial data," *Indian Journal of Computer Science and Engineering*, vol. 1, no. 1, pp. 8–12, 2010.

[5] Z. Li, I. V. Kolmanovsky, E. M. Atkins, J. Lu, D. Filev, and J. Michelini, "Cloud aided safety-based route planning," *2014 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*, pp. 2495–2500, 2014.

[6] P. de Leur and T. Sayed, "Development of a road safety risk index," *Transportation Research Record*, vol. 1784, pp. 33–42, 01 2002.

[7] H. Wu and Z. Zhang, "A framework for developing road risk indices using quantile regression based crash prediction model," Tech. Rep., 2012.

[8] R. Andrášik, M. Bíl, Z. Janoška, and V. Valentová, "Identification of curves and straight sections on road networks from digital vector data."

[9] J. J. Watkins, "Mathematical omnibus: Thirty lectures on classic mathematics by dmitry fuchs and serge tabachnikov," *The Mathematical Intelligencer*, vol. 32, no. 2, pp. 71–72, Jun 2010. [Online]. Available: https://doi.org/10.1007/s00283-009-9101-7

[10] A. Barsi and C. Heipke, "Detecting road junctions by artificial neural networks," in *2003 2nd GRSS/ISPRS Joint Workshop on Remote Sensing and Data Fusion over Urban Areas*. IEEE, 2003, pp. 129–132.

[11] S. Lanser and W. Eckstein, "Eine modifikation des deriche-verfahrens zur kantendetektion," vol. 290, 01 1991, pp. 151–158.

[12] D. Habermann, C. E. Vido, F. S. Osório, and F. Ramos, "Road junction detection from 3d point clouds," in *2016 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2016, pp. 4934–4940.

[13] K. Granström, T. B. Schön, J. I. Nieto, and F. T. Ramos, "Learning to close loops from range data," *The International Journal of Robotics Research*, vol. 30, no. 14, pp. 1728–1754, 2011. [Online]. Available: https://doi.org/10.1177/0278364911405086

[14] C. M. Bishop *et al.*, *Neural networks for pattern recognition*. Oxford university press, 1995.

[15] Z. Q. John Lu, "The elements of statistical learning: Data mining, inference, and prediction," *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, vol. 173, no. 3, pp. 693–694, 2010.

[16] B. Scholkopf and A. J. Smola, *Learning with kernels: support vector machines, regularization, optimization, and beyond*. MIT press, 2001.

[17] A. Hata, D. Habermann, D. Wolf, and F. Osório, "Crossroad detection using artificial neural networks," in *Engineering Applications of Neural Networks*, L. Iliadis, H. Papadopoulos, and C. Jayne, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2013, pp. 112–121.

[18] Y. Yao, Z. Peng, and B. Xiao, "Parallel hyper-heuristic algorithm for multi-objective route planning in a smart city," *IEEE Transactions on Vehicular Technology*, vol. 67, no. 11, pp. 10 307–10 318, 2018.

# Value-Driven use Cases Triage for Embedded Systems: A Case Study of Cellular Phone

Neunghoe Kim[1]

Department of Computer
Korea University, Seoul, Republic of Korea

Younkyu Lee[2]

Multimedia Processing Lab Samsung Advanced Institute of
Technology, Suwon, Republic of Korea

Vijayan Sugumaran[3]

Department of Decision and Information Sciences
Oakland University, Rochester, United States

Soojin Park[4, *]

Graduate School of Management of Technology
Sogang University, Seoul, Republic of Korea

*Abstract*—**A well-defined and prioritized set of use cases enables the enhancement of an entire system by focusing on more important use cases identified in the previous iteration. These use cases are given more opportunities to be refined and tested. Until now, use case prioritization has been done from a user perspective, and through balanced measurement of actors/ objects usage. Lack of cost consideration for realization, however, renders it ineffective for economic purposes. Hence, this study incorporates the 'value' concept, based on cost benefit analysis, in use case prioritization for embedded systems. The use case satisfaction level is used as the surrogate for 'benefit', and the complexity of implementation for 'cost'. Based on the value, use cases are prioritized. As a proof-of-concept, we apply our value-based prioritization method to the development of a camera system in a cellular phone.**

*Keywords—Value-based software engineering; use case triage; embedded system; cost-benefit analysis*

## I. INTRODUCTION

Software development has evolved around users and there is more focus on developing use cases to cover the entire process such as inspection, requirement analysis, and testing. How to select use cases is one of the main issues in the planning process of an iteration-based software development project using the Unified Software Development Process [1]. The earlier a use case is placed in iteration, the more test opportunities it gets resulting in higher quality. Therefore, it is vital for improving the quality of the entire system to detect core use cases in the earlier stage and include them in the iteration plan.

Requirements prioritization has been an active area of research. Herrmann and Daneva [2] have conducted a systematic review of this literature and classified the existing requirements prioritization approaches based on several criteria and have identified fifteen well established methods. Based on their analysis, they point out several weaknesses among existing methods such as: a) not being able to estimate the benefits at the individual requirements level as opposed to the system level, b) lack of guidance for selecting the appropriate cost estimation technique given a specific context, and c) not taking into account the dependencies between various requirements [2]. Some of these limitations can be mitigated by

conducting cost-benefit analysis at the use case level. We contend that this process enables us to better estimate the benefits, account for complexity as well as dependencies between requirements. Thus, use case prioritization provides a systematic approach for analyzing the benefits and cost of realization of requirements. In this research, we utilize the principles from cost-benefit analysis and determine the value of use cases by combining user preferences and complexity. Specifically, our proposed value-based use case prioritization method assigns higher priority to use cases that provide maximum satisfaction to users, while consuming minimum time and cost of realization.

The rest of the paper is structured, as follows. Section 2 discusses the prior studies related to the prioritization methods of use cases. The proposed approach and case study for value-oriented prioritization of use cases are described in Section 3. In Section 4, the evaluation conducted to verify the validity and efficiency of the proposed approach earlier is discussed. Section 5 concludes the paper, and outlines the future work.

## II. RELATED WORK

Conventional studies on use case prioritization assign priority based on objective measurement of actors and objects usage metrics as well as subjective stakeholder viewpoints [3]. However, they are not very effective for actual application since they lack economic consideration of the costs of realization. Karlsson and Ryan [4] discuss a cost-value approach for prioritizing requirements, however, they do not provide a systematic way of estimating cost. They also don't consider dependency relationships and software quality attributes.

Numeral Assignment Technique, Planning Game and the Analytic Hierarchy Process ("AHP") are three representative examples of the stakeholder preference-based approach. AHP incorporates pairwise comparison and is used for computation of relative values from stakeholders and cost of individual requirements. Technical complexity-based approaches are cost estimation methods such as the Lines of Code ("LOC"), Constructive Cost Model ("COCOMO"), Function Point Method, and Use Case Points Method. Models such as LOC and COCOMO are inappropriate for cost estimation of projects

---

*Corresponding Author

where the number of lines of code is hard to estimate. The problem with Function Point Method is that application to embedded systems is almost impossible due to its indifference to the internal operations of software. Finally, with respect to Use Case Points Method, the criteria are too unclear to determine the complexity and the weight of each use case for embedded systems. Hence, propose a new method for estimating use case complexity that reflects the characteristics of embedded systems.

By using the notion of value, we introduce a new method for prioritizing use cases for an embedded system, which considers both stakeholder preferences and technical complexity.

## III. Value-Driven use Cases Triage Method and Case Study

### A. Value of a use Case

The term "value" is defined in different ways according to the needs of different fields (e.g. marketing, business management). Typically, Cost-Benefit analysis compares benefits and costs of a project or a system. In our work, we define value as the ratio between the benefits and cost of software development, which is computed using the following equation [5]:

Benefit Cost Ratio (BCR) = Benefit/Cost          (1)

In general, "cost" includes all costs ranging from capital, planning, installation, application development, to continuous maintenance, while "benefit" includes benefits from savings in labor and operation cost, and improved productivity [6].

Adapting the above equation to the software development context, the value of a use case can be determined as follows:

Value of Use Case = Level of Satisfaction of the Use Case

/Cost of Realization of the Use Case          (2)

As shown in equation 2, the benefit of a single use case that specifies a particular functionality can be substituted with the satisfaction level that a user gets from the use case. Satisfaction level is measured using pairwise comparison of all the use cases, which is part of AHP. The cost factor is the cost of realization of a use case and depends on the complexity of that use case. It is measured using the extended complexity factors of the sequence diagram associated with the use case. According to equation 2, the more satisfied a user feels about a single use case, and at the same time, the less it costs for its realization, then higher the value of that use case.

### B. Value-Driven use Cases Triage Method

The proposed method of prioritizing use cases involves cost-benefit analysis, which in turn computes the ratio between the benefit and the cost of developing that function. The proposed method consists of the following steps: 1) investigating relative satisfaction level, 2) identification of inter-component collaboration using state and sequence diagrams, 3) complexity calculation, and 4) use case value adjustment. Each of these steps is briefly described below and the computations done in each step is summarized in Table I.

The relative satisfaction level in Step 1 is obtained by means of AHP. Through pairwise comparison, we measure user satisfaction from realizing a function represented by a particular use case compared to other functions. The results form a comparison matrix, which is used to determine the relative satisfaction levels of use cases by averaging over the normalized columns, as proposed by Thomas Saaty [7].

The collaboration between components in a use case is examined in Step 2 using a state diagram and a sequence diagram. The Gray-Box based requirements specification method for embedded systems proposed in [8] is used to generate these diagrams. Among the objects constituting an embedded system such as controller, sensor, and actuator, we focus on the state diagram (top-level) for the embedded controller object. The information on how the state transition of the system interacts with internal components is represented in the sequence diagram.

In Step 3, the complexity of each use case is calculated based on the sequence diagram. The objects in the sequence diagram are classified as simple, average, and complex and weights assigned for each type. Similarly, messages are classified as synchronous or asynchronous, with weights given to each type. Since actors do not affect the complexity of software development, they are ignored. The weights are determined by an expert, since it is heavily dependent on the project and domain characteristics. The complexity of each use case is computed as follows using the equation discussed in [9] (shown in Table I). First, we count the total number of actors collaborating within the sequence diagram. Then, we count the number of objects in each category and multiply it by the corresponding weight for that category. Then, these weighted numbers are summed up. Similarly, we do the same type of computation for messages. The complexity of the sequence diagram is determined by adding up the scores for each of the three parts. The complexity of a use case is then determined by summing up the complexities of all the sequence diagrams that are associated with that use case. Overlapping in computation is avoided by not counting the actors, objects, and messages more than once if they appear in many sequence diagrams.

In Step 4, the value of each use case is determined by dividing the relative satisfaction level generated in Step 1 by its complexity computed in Step 3. Use case values are then adjusted by considering the dependency relationships among the use cases and their expected quality levels. Specifically, the value is adjusted if include and extend relationships exist among use cases, or when the sequential order of the use cases is established by the preconditions existing in them. The adjustment to reflect expected quality levels is to take into account the expectations on distinct quality attributes given for each use case. According to [10], the quality attributes relevant for embedded systems are: reliability, usability, performance, real timeliness, and purpose limitation. The expected quality level of each use case for each attribute is categorized as high, medium, or low. The value of each empirical weight is determined based on the project and domain characteristics. The use case value is multiplied by the weighted quality level scores to determine the adjusted value. These adjusted values determine the final order of development (priority) for each use case.

TABLE I.        COMPUTATIONS USED IN VALUE-DRIVEN USE CASES TRIAGE

---

**Step 1:        Investigating Relative Satisfaction Level (using AHP)**

$$\begin{bmatrix} Xij & \cdots & Xij \\ \vdots & \ddots & \vdots \\ Xij & \cdots & Xij \end{bmatrix} \bullet\bullet\bullet\bullet \begin{bmatrix} Xij & \cdots & Xij \\ \vdots & \ddots & \vdots \\ Xij & \cdots & Xij \end{bmatrix} \Rightarrow \begin{bmatrix} Xij & \cdots & Xij \\ \vdots & \ddots & \vdots \\ Xij & \cdots & Xij \end{bmatrix} \Rightarrow \begin{bmatrix} SL_1 \\ \vdots \\ SL_n \end{bmatrix} \quad \left( \begin{array}{c} \text{"Recording a Video"} \\ \text{Use Case} = 16 \end{array} \right)$$

|   Use Case Comparison Matrices from All Users   |   Comparison Matrix Created using Geometric Mean of Individual Scores   |   Satisfaction Levels of Use Cases   |   Example   |

---

**Step 2:        Identification of Inter-Component Collaboration**

The Gray-Box technique [8] is used to generate the State Diagram and the Sequence Diagram

---

**Step 3:        Complexity Calculation (adapted from [9])**

- Complexity of a Sequence Diagram = Actor Complexity+Object Complexity+Message Complexity

  = ∑No. of actors+∑(No. of objects*object weight)+∑(No. of message*message weight)

- Complexity of Use Case = ∑(Complexity of Constituent Sequence Diagram)

**Example:        Complexity Calculation from Case Study (for the "Recording a Video" Use Case)**

- Actor Complexity = User+Camsensor+MIC = 1+1+1 = 3
- Object Complexity = CamsensorIF*Complex+AudioControlIF*Average = (1*2)+(1*1.5) = 3.5
- Message Complexity = oper_set_state()*Asynchronous+get_sensordata()*Synchronous+ oper_record()*Asynchronous+oper_take_movie()*Asynchronous+set_autdio_path()*Asynchronous+get_audio_input_data()*Synchronous+check_free_space()*Asynchronous+save_recorded_data()*Asynchronous+get_temp_filename()*Asynchronous+save_file()*Asynchronous+stop_record()*Asynchronous+restartPreview()*Asynchronous = (1*1)+(1*2)+(1*1)+(1*1)+(1*1)+(1*2)+(1*1)+(1*1)+(1*1)+(1*1)+(1*1)+(1*1) = 14
- Complexity of the Sequence Diagram = 3+3.5+14 = 20.5
- Complexity of Use Case = 20.5 (This use case contained only one sequence diagram)

---

**Step 4:        Use Case Value Adjustment**

- Adjusted Value of Base Use Case under <<extend>> = value of base use case*∑extend weight
- Adjusted Value of Included Use Case under <<include>> = value of included use case*∑include weight
- In Case of Precondition, Value of Prerequisite Use Case = value of prerequisite use case*∑precondition weight
- Adjustment of Quality Attributes = value of use case*∑(No. of quality attributes*weight of quality attributes category)

**Example:        Use Case Value Adjustment ("Recording a Video" Use Case)**

- Value of Use Case = Relative Satisfaction Level/Complexity = 16/20.5 = 78 (We multiply this by 100 and round off to the nearest integer)
- Dependency (Not Related)
- Adjustment of Quality Attributes = 78*(reliability*Medium+usability*High+performance*High+real timeliness*High+purpose limitation*Low) = 78*(0.2+0.3+0.3+0.1+1) = 172 (1 is added to the sums of empirical weights to compensate for reduction of values owing to decimal values of weights)

---

*C. Case Study: A Camera System in Cellular Phone*

To demonstrate the feasibility of calculating use case values, we have conducted a case study using the camera system in a cellular phone. We implemented this system in a domestic 3G feature handset for a global electronics company. The project utilized 21 software developers and took 7 months to complete. First, use case modeling was carried out based on the requirements for the camera system. Ten use cases (Previewing, Taking a Snapshot, Recording a Video, Postviewing, Playing a Video, Album Management, Editing Photo&Video, Sending Photo&Video, Printing a Photo,

Albumview) and 9 actors (Camsensor, MIC, Speaker, User, LCD, Wallpaper Manager, MMS Manager, Bluetooth Manager, Printer) were derived. The "Recording a Video" use case is used to demonstrate the computations in our approach.

*1) STEP 1 Investigating relative satisfaction level:* As part of the AHP methodology, the use case comparison matrices from 30 users were used to generate a single comparison matrix by taking the geometric mean of the individual comparison scores. This matrix was then normalized and the relative satisfaction levels for each of the use cases were

determined. The consistency ratio of the corresponding values in the comparison matrices from users was less than 0.1, indicating that these matrices are reliable [7]. To express the relative satisfaction levels as integer values, they were multiplied by 100 and rounded off to the nearest integer. The relative satisfaction levels for the use cases are: Taking a Snapshot = 24, Recording a Video = 16, Previewing = 11, Postviewing = 11, Playing a Video = 10, Album Management = 7, Albumview = 7, Sending Photo&Video = 6, Editing Photo&Video = 5, and Printing a Photo = 4.

*2) STEP 2 identification of inter-component collaboration:* A state diagram was created for the 'CameraController' component that controls the state and transition information of the camera system in the cellular phone. As an essential object controlling the state of the entire system, the 'CameraController' object functions as the owner of the state diagram [8]. In the state diagram, the camera system in a cellular phone should maintain states such as "idle," "initialized," "preview," "postview," "recording," "albumview," "editing," "sending," "printing," "playing, " "snapshot," and "stopped." As this object is used in computing the complexity of the use case, the "Recording a Video" use case triggers the state transition upon its activation, and the transition goes through the following flow: "preview → recording → preview."

Next, the sequence diagram is created, while assigning the events and the actions shown on the state diagram and marking them chronologically. In the sequence diagram of the "Recording a Video" use case, the CameraController in the "preview" state sends commands such as Oper_set_state() and Get_sensor_data() to Camsensor_IF object, upon receipt of the StartRecord event invoked by a user. Then, while carrying out its own oper_record(), the "preview" state transitions into the "recording" state. During this transition, Camsensor_IF also sends certain messages to Camsensors to fulfill the objective(s) of the message(s) it has received. Through this analytical process, it is made clear how the goal of each use case is accomplished by understanding what messages are sent and received by each component object constituting the entire system.

*3) STEP 3 complexity calculation:* The complexity of a use case is determined based on the information contained in the sequence diagram. The sequence diagram pertaining to the "Recording a Video" use case contains three actors, two objects and 12 messages. The object weights applied in this project are Simple = 1.0, Average = 1.5, and Complex = 2.0, while message weights are set as Synchronous = 2.0 and Asynchronous = 1.0. For example, CamsensorIF was classified as "complex,' while AudioControlIF was categorized as "average." In the case of messages, get_sensordata() and get_audio_input_data() were marked as "synchronous", while the others were deemed "asynchronous." Thus, the complexity of the sequence diagram for the "Recording a Video" use case is computed to be 20.5. It is to be noted that this use case contained only one

sequence diagram. Therefore, the sequence diagram complexity also represents the complexity of the use case. The complexity values computed for the use cases in the case study are: Previewing = 28, Recording a Video = 20.5, Postviewing = 18.5, Album Management = 16.5, Sending Photo&Video = 15, Printing a Photo = 14, Playing a Video = 12.5, Taking a Snapshot = 12, Editing Photo&Video = 11, and Albumview = 10.

*4) STEP 4 use case value adjustment:* The use case values are obtained by dividing the relative satisfaction levels by the complexities. This ratio is expressed as whole number, by multiplying it by 100 and rounding off to the nearest integer. The values of the use cases computed in the case study are as follows: Taking a Snapshot = 200, Playing a Video = 80, Recording a Video = 78, Albumview = 70, Postviewing = 59, Editing Photo&Video = 45, Album Management = 42, Sending Photo&Video = 40, Previewing = 39, and Printing a Photo = 29. In the case of the "Playing a Video" use case, its relative satisfaction level is 10, or the 5th highest among the ten use cases, and its complexity is 12.5, or the 4th lowest among them. However, its value computes to 80, the second highest among the ten use cases.

Next, the use case values are adjusted based on inter use case dependencies and each case's expected quality level. The values are rounded off to the nearest integer. The value adjustment for the "Recording a Video" use case is shown in Table I. The adjusted values of the use cases in the case study are as follows: Taking a Snapshot = 420, Albumview = 189, Recording a Video = 172, Previewing = 128, Playing a Video = 104, Postviewing = 65, Album Management = 63, Editing Photo&Video = 50, Sending Photo&Video = 48, and Printing a Photo = 35. For the adjustment, a weight of 0.1 is assigned to dependency relationship, while three weights are assigned to the expected quality level (i.e. high = 0.3, medium = 0.2, low = 0.1). Playing a Video, Albumview, Postviewing, Editing Photo&Video, Sending Photo&Video, and Previewing use cases had their priority positions changed after the adjustment.

## IV. Evaluation

In the case study described in section 3, we discussed how the use case values for a cellular phone camera system were derived. To demonstrate the effectiveness of the process, we have to answer the following three questions:

- Does the use case complexity computed through our approach match the complexity experienced in realizing the use case?

- How much do the stakeholders trust the results of our proposed method after applying it to their processes?

- Are the results from our approach more useful compared to the previous use case prioritizations that were being used?

We demonstrate the validity of our complexity calculation by showing the proportional relationship between our complexity values and the actual LOC values for corresponding use cases. In addition to this quantitative

evaluation, we demonstrate the trust shown by different stakeholders in our proposed approach by administering a survey to the marketing staff and development engineers of embedded software in the case study organization. Lastly, we show the usefulness of our proposed method through comparative analysis of the results from our method and the results from previous use case prioritizations generated by the development engineers.

### A. Verification of the Complexity-Calculation Process

We counted the number of lines of executable source code upon completion of the development of the cellular phone camera system. As discussed earlier, the reason for measuring LOCs is to check whether or not the complexity-based use case priority, which had been generated prior to realization, matches the LOC size-based priority upon completion. If the two types of priority are in direct proportion to each other, our proposed method of estimating use case values is applicable to the actual development of embedded systems. To investigate the relationship between our use case complexities and the actually realized LOCs, we plotted the LOCs and the corresponding use case complexities, as shown in Fig. 1(a), (b). The use case complexity rank and the LOC rank corresponds to the rank ordering of use cases from the most complex (1) to the least complex (10) in the case study. As seen from Fig. 1(a), the complexity rank and the LOC rank for the use cases follow each other closely.

In this case study, the correlation coefficient between the complexities of the use cases computed based on our approach and the corresponding LOC was determined to be 0.96, meaning a strong relationship between them. Also, we ran a simple linear regression model with complexity as the independent variable and LOC as the dependent variable. The regression coefficients and the $R^2$ are shown in Fig. 1(b). The $R^2$ value is 0.9207, which is very significant and strongly suggests a linear relationship between the use case complexities computed through our approach and the resulting LOC. Considering these findings, it is fair to conclude that the use case complexity computation method proposed herein is a good indicator of the complexity of the actually realized code.

### B. Acceptance of the Proposed Method

A survey was administered to 40 developers of embedded systems and 10 marketing staff members from the corporation that developed the camera system for the cellular phone. On average, the marketing staff members had five years of experience, and the developers had 7 years of experience.

The survey contained questions focusing on three main aspects: a) choosing use cases based on cost, b) trustworthiness of the results from our approach, and c) usefulness of our approach. With respect to the need for choosing use cases in consideration of development costs, 90% of the marketing staff and 95% of the developers indicated that cost should be considered in selecting use cases for implementation. With respect to our model's trustworthiness, 80% of the marketing staff and 85% of the developers responded positively. In terms of usefulness of our approach, 70% of the marketers and 40% of the developers acknowledged that the method would be useful in their organization. The lower percentage value of the developers may be due to unfamiliarity with modeling, personal habits, corporate culture and internal structural/organizational issues.
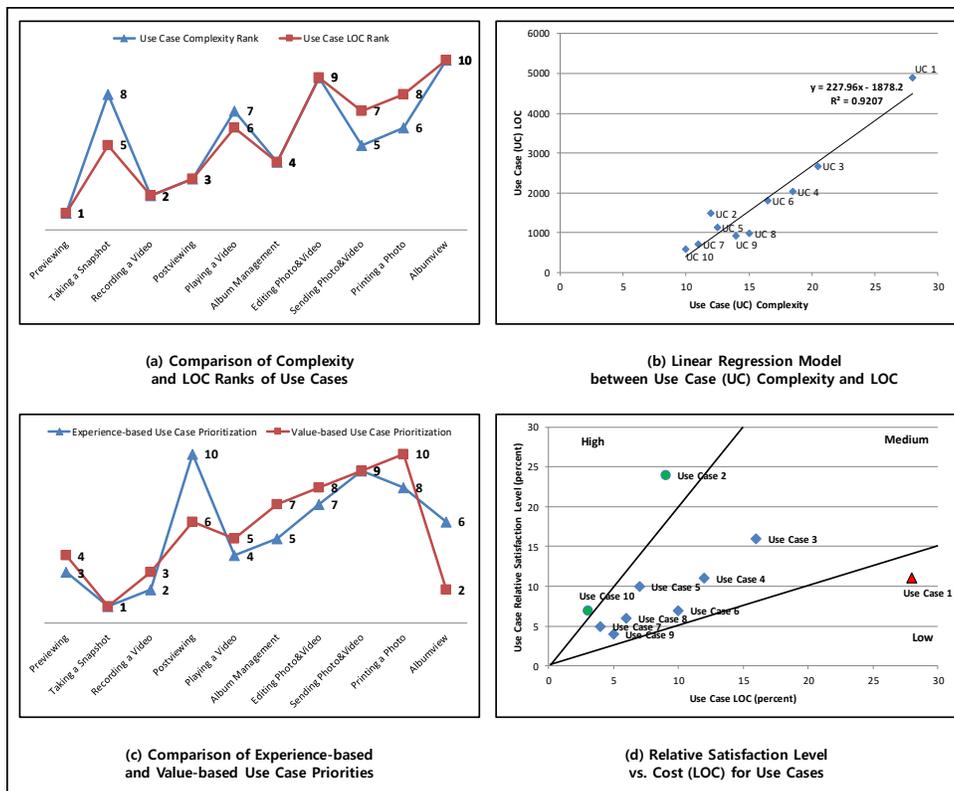


(a) Comparison of Complexity
and LOC Ranks of Use Cases

(b) Linear Regression Model
between Use Case (UC) Complexity and LOC

(c) Comparison of Experience-based
and Value-based Use Case Priorities

(d) Relative Satisfaction Level
vs. Cost (LOC) for Use Cases

Fig. 1.    Results from Case Study (Camera System in Cellular Phone).

### C. Usefulness of the Proposed Method

For the ten use cases in the case study, the developers independently estimated the priority based on their prior experience. Developers tend assign high priority to use cases corresponding to basic functions even if they have high cost of realization. For the other functions, they assign lower priority even if they are complex. Fig. 1(c) shows the experience-based use case priorities and the value-based priorities developed using our approach. While the priorities are similar for a few use cases, there is considerable difference for some of them. To further analyze the differences, we have developed a cost-value diagram similar to the one discussed in [4]. The normalized relative satisfaction levels and the normalized LOC values for the use cases are plotted, as shown in Fig. 1(d). Based on the value of use cases (ratio of satisfaction level and LOC), we group the use cases into three categories (High, Medium, Low). For high value use cases, the ratio exceeds 2, for medium value, between 0.5 and 2, and for low value below 0.5, as used in [4]. As seen in Fig. 1(d), Use Case 2 (Taking a Snapshot) and 10 (Albumview) are high value use cases and they were correctly assigned a high priority value of 1 and 2 in our approach. However, the developers assigned a priority of 6 for Use Case 10, thus failing to identify this high value use case. Use Case 1 (Previewing) is a low value use case, as shown in Fig. 1(d). Our approach assigned a priority of 4, while the developers assigned a priority of 3. Thus, our approach is better able to assign more appropriate priorities compared to the experience based use case prioritization.

## V. CONCLUSION

This study has proposed a value-based method for prioritizing use cases. This approach is used to improve quality by discerning "valuable" use cases and incorporating them in the iteration plan at an earlier stage. To demonstrate the validity and usefulness of the proposed approach, a case study and a survey were conducted.

The contributions of this study are as follows:

- In prioritizing use cases, the notion of value is defined based on the external requirement of "satisfaction level" and the internal requirement of "cost." Our prioritization process is based on value, which is a balanced metric. In this study, the cost of each use case refers to the effort required to realize that use case and it increases in direct proportion to the complexity of the use case.

- To determine the complexity of use cases tailored to the embedded system domain, it is computed based on the inter-component collaboration model.

- The validity of our model has been demonstrated by applying our complexity estimation model to an actual case and showing that the complexity estimates produced through our approach matched the actually realized LOCs.

Although we have demonstrated the feasibility of our approach, further work is needed to fully establish its efficacy. The evaluation results verify the validity of the complexity estimation of each use cases. However, further work is needed to verify the validity of users' satisfaction. As part of future work, a quantitative study will be conducted to investigate how much improvement can be achieved in the quality of the software product, when the relevant iteration planning is carried out in accordance with the prioritization results produced through our model.

### REFERENCES

[1] I. Jacobson, G. Booch, and J. Rumbaugh, The Unified Software Development Process, Addison-Wesley Professional, 1999.

[2] A. Herrmann and M. Daneva, "Requirements prioritization based on benefits and cost prediction: an agenda for future research," International Requirements Engineering Conference, pp. 125-134, September 2008.

[3] F. Moisiadis, "Prioritizing use cases and scenarios," International Conference on Technology of Object-Oriented Languages and Systems, pp. 108-119, November 2000.

[4] J. Karlsson, K. Ryan, "A cost-value approach for prioritizing requirements," IEEE Software, vol. 14, no. 5, pp. 67-74, September/October 1997.

[5] H. Erdogmus, "Cost-benefit analysis of software development techniques and practices," International Conference on Software Engineering, pp. 178-179, May 2007.

[6] T. Pisello, Return on Investment for Information Technology Providers, New Canaan Connecticut: Information Economics Press, 2001.

[7] T. L. Saaty, The Analytic Hierarchy Process, McGraw-Hill, 1980.

[8] S. Park, S. Park, "A gray-box based software requirements specification method for embedded systems," Journal of Korean Institute of Information Scientists and Engineers, vol. 38, no. 9, pp. 485-490, September 2011.

[9] A. Kanjilal, S. Sengupta, S. Bhattacharya, "Analysis of complexity of requirements: a metrics based approach," India Software Engineering Conference, pp. 131-132, February 2009.

[10] J. Lim, H. Yoon, "Extraction of quality attribute for designing the S/W architecture in weapon systems embedded software," Korean Fuzzy Logic and Intelligent Systems Society Autumn Conference, pp. 268-271, November 2006.