

In the Name of God

Journal of Information Systems & Telecommunication

Vol. 10, No.4, October-December 2022, Serial Number 40

Research Institute for Information and Communication Technology
Iranian Association of Information and Communication Technology
Affiliated to: Academic Center for Education, Culture and Research (ACECR)

Manager-in-Charge: Dr. Habibollah Asghari, ACECR, Iran

Editor-in-Chief: Dr. Masoud Shafiee, Amir Kabir University of Technology, Iran

Editorial Board

Dr. Abdolali Abdipour, Professor, Amirkabir University of Technology, Iran
Dr. Ali Akbar Jalali, Professor, Iran University of Science and Technology, Iran
Dr. Alireza Montazemi, Professor, McMaster University, Canada
Dr. Ali Mohammad-Djafari, Associate Professor, Le Centre National de la Recherche Scientifique (CNRS), France
Dr. Hamid Reza Sadegh Mohammadi, Associate Professor, ACECR, Iran
Dr. Mahmoud Moghavvemi, Professor, University of Malaya (UM), Malaysia
Dr. Mehrnoush Shamsfard, Associate Professor, Shahid Beheshti University, Iran
Dr. Omid Mahdi Ebadati, Associate Professor, Kharazmi University, Iran
Dr. Rahim Saeidi, Assistant Professor, Aalto University, Finland
Dr. Ramezan Ali Sadeghzadeh, Professor, Khajeh Nasireddin Toosi University of Technology, Iran
Dr. Sha'ban Elahi, Associate Professor, Tarbiat Modares University, Iran
Dr. Shohreh Kasaei, Professor, Sharif University of Technology, Iran
Dr. Saeed Ghazi Maghrebi, Assistant Professor, ACECR, Iran
Dr. Zabih Ghasemlooy, Professor, Northumbria University, UK

Executive Editor: Dr. Fatemeh Kheirkhah

Executive Manager: Shirin Gilaki

Executive Assistants: Mahdokht Ghahari, Ali Mokhtarani, Ali BoozarPoor

Print ISSN: 2322-1437

Online ISSN: 2345-2773

Publication License: 91/13216

Editorial Office Address: No.5, Saeedi Alley, Kalej Intersection., Enghelab Ave., Tehran, Iran,
P.O.Box: 13145-799 Tel: (+9821) 88930150 Fax: (+9821) 88930157

E-mail: info@jist.ir , infojist@gmail.com

URL: www.jist.ir

Indexed by:

- | | |
|---|-------------------------|
| - SCOPUS | www.Scopus.com |
| - Index Copernicus International | www.indexcopernicus.com |
| - Islamic World Science Citation Center (ISC) | www.isc.gov.ir |
| - Directory of open Access Journals | www.Doaj.org |
| - Scientific Information Database (SID) | www.sid.ir |
| - Regional Information Center for Science and Technology (RiCeST) | www.ricest.ac.ir |
| - Magiran | www.magiran.com |

Publisher:

Iranian Academic Center for Education, Culture and Research (ACECR)

This Journal is published under scientific support of
Advanced Information Systems (AIS) Research Group and
Telecommunication Research Group, ICTRC

Acknowledgement

JIST Editorial-Board would like to gratefully appreciate the following distinguished referees for spending their valuable time and expertise in reviewing the manuscripts and their constructive suggestions, which had a great impact on the enhancement of this issue of the JIST Journal.

(A-Z)

- Afsharirad, Majid, Kharazmi University, Tehran, Iran
- Fathi, Amir, Urmia University, Urmia, Iran
- Fortaki, Tarek, University of Batna, Batna, Algeria
- Ghaffari, Ali, Islamic Azad University, Tabriz Branch, Iran
- Gholami, Mohammad, Babol Noshirvani University of Technology, Mazandaran, Iran
- Hasan, Junayed, Universiteit van Ulsan, Ulsan, South Korea
- Lotfi, Reza, Yazd University, Yazd, Iran
- Khazaei, Mehdi, Kermanshah University of Technology, Kermanshah, Iran
- Kashef, Seyed Sadra, Urmia University, Urmia, Iran
- Kolahkaj, Maral, Islamic Azad University, Karaj Branch, Iran
- Mahdieh, Omid, University of Zanjan, Zanjan, Iran
- Mansoorizadeh, Muharram, Bu-Ali Sina University, Hamedan, Iran
- Minoofam, Seyed Amir Hadi, Qazvin Islamic Azad University, Qazvin, Iran
- Mirzaei, Abbas, Islamic Azad University, Ardabil, Iran
- Moradi, Gholamreza, Amirkabir University, Tehran, Iran
- Mirroshandel, Seyed Abolghasem, University of Guilan, Rasht, Iran
- Mohammadi, Mohammad Reza, Iran University of Science and Technology, Tehran, Iran
- Mohamed, Lashab, University of Oum, Larbi Ben M'hidi, Algeria
- Omid Mahdi, Ebadati, Kharazmi University, Tehran, Iran
- Pashazadeh, Saeid, Tabriz University, Tabriz, Iran
- Qingliang, Zhao, Northwestern University, Evanston, United State
- Qamar, Nafees, Governors State University, Illinois, south of Chicago
- Rahman, Saifur, Virginia Tech, Riva San Vitale, Switzerland,
- Rasi, Habib, Shiraz University of Technology, Shiraz, Iran
- Sadeghmohammadi, Hamidreza, ACECR, Tehran, Iran
- Saadatfar, Hamid, University of Birjand, Iran
- Samsami Khodadad, Farid, Amol University of Special Technology, Amol, Iran
- Soleimani Gharehchopogh, Farhad, Islamic Azad University Urmia, Iran
- Tanhaei, Mohammad, Ilam University, Ilam, Iran
- Tourani, Mahdi, University of Birjand, South Khorasan, Iran
- Yaghoobi, Kaebeh, Ale Taha Institute of Higher Education, Tehran, Iran
- Zayyani, Hadi, Shiraz University of Technology, Shiraz, Iran

Table of Contents

• A High Performance Dual Stage Face Detection Algorithm Implementation using FPGA Chip and DSP Processor...	241
M V Ganeswaea Rao, P Ravi Kumar and T Balaji	
• A Novel Detector based on Compressive Sensing for Uplink Massive MIMO Systems.....	249
Mojtaba Amiri and Amir Akhavan	
• A Hybrid Approach based on PSO and Boosting Technique for Data Modeling in Sensor Networks	257
Hadi Shakibian and Jalal A. Nasiri	
• Detection of Attacks and Anomalies in the Internet of Things System using Neural Networks Based on Training with PSO Algorithms, Fuzzy PSO, Comparative PSO and Mutative PSO	268
Mohammad Nazarpour, Navid Nezafati and Sajjad Shokouhyar	
• ARASP: An ASIP Processor for Automated Reversible Logic Synthesis	277
Zeinab Kalantari, Marzieh Gerami and Mohammad Eshghi	
• Propose an E-CRM Model based on Mobile Computing Technology in Pharma Distribution Industry	285
Alireza Kamanghad, Gholamreza Hashemzadeh Khorasgani , Mohammadali Afshar Kazemi and Nosratollah Shadnoosh	
• A Novel Approach for Establishing Cinnectivity in Partitioned Mobile Sensor Networks using Beamforming Techniques.....	298
Abbas Mirzaei and Shahram Zandiyan	
• Energy-Efficient User Pairing and Power Allocation for Granted Uplink-NOMA in UAV Communication Systems	310
Seyyed Hadi Mostafavi-Amjad, Vahid Solouk and Hashem Kalbkhani	

A High Performance Dual Stage Face Detection Algorithm Implementation using FPGA Chip and DSP Processor

MV Ganeswara Rao^{1*}, P Ravi Kumar¹, T Balaji²

¹. Department of Electronics and Communication Engineering, Shri Vishnu Engineering College for Women, Bhimavaram, AP, India.

². Department of Electronics and Communication Engineering PVP Siddhartha Institute of Technology, Vijayawada, AP, India

Received: 12 Oct 2021/ Revised: 01 Jan 2022/ Accepted: 01 Feb 2022

Abstract

A dual stage system architecture for face detection based on skin tone detection and Viola and Jones face detection structure is presented in this paper. The proposed architecture able to track down human faces in the image with high accuracy within time constrain. A non-linear transformation technique is introduced in the first stage to reduce the false alarms in second stage. Moreover, in the second stage pipe line technique is used to improve overall throughput of the system. The proposed system design is based on Xilinx's Virtex FPGA chip and Texas Instruments DSP processor. The dual port BRAM memory in FPGA chip and EMIF (External Memory Interface) of DSP processor are used as interface between FPGA and DSP processor. The proposed system exploits advantages of both the computational elements (FPGA and DSP) and the system level pipelining to achieve real time performance. The present system implementation focuses on high accurate and high speed face detection and this system evaluated using standard BAO image database, which include images with different poses, orientations, occlusions and illumination. The proposed system attained 16.53 FPS frame rate for the input image spatial resolution of 640X480, which is 23.4 times faster detection of faces compared to MATLAB implementation and 12.14 times faster than DSP implementation and 2.1 times faster than FPGA implementation.

Keywords: Face detection; Heterogeneous System; FPGA; DSP.

1- Introduction

The image processing algorithms with real time performance and high accuracy are idly used in diversified fields such as surveillance, surface quality inspection, Robotic vision, Assistive technology etc. [1]. The human face detection is one of the popular research areas in the field of image processing. These algorithms are used in many applications like facial recognition in security systems, human computer Interaction (HCI) and so on. In the past, many researchers proposed face detection algorithms, which are computationally simple, but not efficient. However, in the recent past, researchers proposed highly efficient methods, but demanding high computational power. Development of such a computational system becomes a challenging task. In recent years, revolutionary advancements in computational platform after merge of high performance

Digital Signal Processors (DSPs), Graphics Processing Units (GPUs), Application Specific Instruction set Processors (ASIPs) and Field Programmable Gate Array (FPGA) Chips. Each computing elements has its own advantages that make it used in associated application. Many researcher developed by hardware platforms based on single computing element to track the human faces in image.

Yang, et al. [2] Proposed face detection system based on DSP Processor, Patrick, et al. [1] developed DSP based hardware to compress video frames for wireless transmission, In [3] Nguyen et. al. implemented an optimized algorithm for video segmentation on DSP platform, Arya, et al. [4] proposed face detection (using RTC colour model) system based on Quartus II FPFA, Fekih, et al.) [5] presented new hardware architecture for face detection based on Zynq-7000 SoC, in which uses ARM CPU and FPGA as computational elements, Leung, et al. [6] proposed a FPGA platform to extract facial features from images for facial recognition and results

✉ MV Ganeswara Rao
ganesh.mgr@gmail.com

proved that detection rate and the performance are significantly high compared to Desktop implementation and Karnewar, et al. [7] proposed GPU based image processing platform to process geospatial images (CUDA application).

The state of the art Face detection algorithms are demanding a very high computing performance, which is very hard to achieve with single and homogeneous computing elements. One recent approach to meet this performance demand is to use heterogeneous systems formed by interconnecting a no of heterogeneous computing elements to build a huge computational platforms. These platforms are widely used in performance intense applications such as image processing and Medical instrumentation etc. Simple heterogeneous system architecture is presented in Fig. 1.

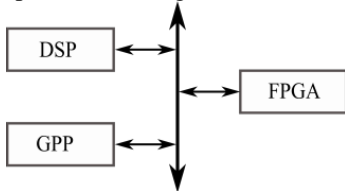


Fig. 1 Simple Heterogeneous structure (example)

In this paper, a new custom heterogeneous face detection platform based FPGA and DSP is proposed. The advantage of this system lies in its two stage systems level pipelining, which is used to attain practical performance. The paper is organized as follows: Section 2 introduces various Heterogeneous platforms proposed by the researchers. Section 3 describes about Proposed Heterogeneous platform architecture. In Section 4, the overview of face detection algorithm is presented. Section 5 deals with implementation of face detection algorithm on FPGA and DSP Processor. Section 6 provides experimental results and comparison with previous works. Section 7 presents conclusion remarks.

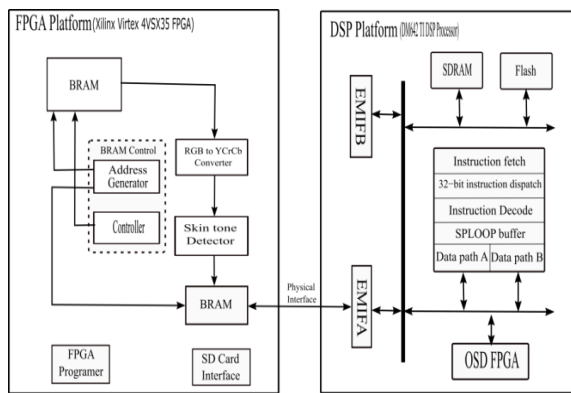


Fig. 2 Block diagram proposed face detection system

2- Allied Work

Considerable research is going on, in the past one and half decade to design and development of Heterogeneous system for real time image processing. Battle, et al. [8]. Proposed a high performance image processing system architecture based on FPGA chip and DSP processor. This system consists of array of DSP processor and the FPGA. The FPGA is used to interconnect these processors.

Liu, et al. [9] proposed multi core GPPs and one GPU (Graphic Processing unit) based heterogeneous platform to tracking human faces in images. The face tracking cannot stable with single information of the face, due to occlusion and illumination problems. Three dissimilar information, wavelet feature, colour histogram and edge orientation histogram are combined to significantly improve the face tracking performance[28][29].

Guo et al. [10] implemented video image correlation algorithm on DSP and FPGA based platform and also multimedia processing algorithms on proposed platform. This system is based on DSP (Multimedia) and FPGA chip. The functions video and audio gathering are implemented on DSP and the FPGA chip is responsible for VGA display, control logic etc.

Wei et al. [11] designed and developed Image processing platform rely on three DSPs and one FPGAs to attain real performance. The EMIF of DSPs is used as to interface DSPs and FPGA to built heterogeneous platform. In the proposed architecture DSPs are used to process core multimedia and FPGA controls data flow in the system. This system exhibits high processing power at a cost of complexity.

3- System Architecture and Overview

This section deals with architecture of proposed system, External Memory Interface (EMIF) of DSP processor and Interface design

3-1- Hardware Architecture

The proposed system architecture for face detection implemented using with Xilinx’s Virtex 4V5X35 FPGA and Texas Instruments DSP TMS320DM642 is shown in Fig. 2. In the first stage, Xilinx’s Virtex 4V5X35 FPGA is used to implement hardware architecture of skin tone detector and two BRAMs to store the image data before and after skin tone detection. There is a large Block RAM resource (3,456KB) in Virtex 4V5X35 FPGA, which is adequate to implement BRAM. This FPGA also offers 18 x 18, two’s complement, signed Multiplier and automatic programmable FIFO logic along with other advanced features [Xilinx 4V5X35 Data Sheet, 2010, p.3].

In the second stage, Texas Instruments (TI) TMS320DM642 is used as a computing element to implement Viola and Jones face detection algorithm and to set up other modules. This processor is a fixed point, high performance digital media processor and delivers the performance up to 5760 MIPS at a clock frequency of 720MHz. The TMS320DM642 provides EMIF (External/Memory Interface) services, which allows the external memory controllers connects to processors. [TMS320DM642 Technical Overview, 2017, p. 30][14]. The EMIF of TMS320DM642 are used to provide interface between the FPGA platform and DSP platform. The MTYPE (Memory Type) field of the CE3 space control register of TMS320DM642 is configured as an asynchronous RAM interface and 32 bit image data.

3-2- EMIF Overview

The External Memory Interface (EMIF) of TMS320DM642 supports interfaces to many external peripherals such as Asynchronous devices (example SRAM, and FIFOs) and Synchronous DRAM (SDRAM). The EMIF can be used as EMIFA supports data bus width 64 or 32 bits and EMIFB supports 16 bits. The signals of EMIFA are presented in Fig. 3

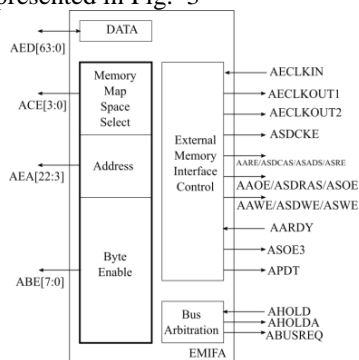


Fig. 3 Signal description of EMIFA

To configure EMIF according to requirement, some registers must be set with the required values. The TI TMS320DM642 memory address space is used to configure various interfaces of EMIF. The base address for EMIFA and EMIFB are 0x0180_0000 and 0x01A8_0000 respectively. The 4 bit CE control Registers CECTL corresponds to the four memory spaces of EMIF. The MTYPE is key filed in the CECTL register, which defines the type of memory interfacing to corresponding memory space.

3-3- Interface Design

There are two ways to interface FPGA to DSP, one way is to use dual port BRAM, and other one is FIFO. In this

design, dual port synchronous BRAM used to interface with ports of EMIF

This interface utilizes the Virtex -4 IOB, which is configured as simple input and tri-state buffer. This memory based interface between FPGA and DSP significantly reduce the interface logic to obtain maximum image data throughput. The EMIF of TMS320DM642 uses the memory configured in FPGA as a memory system with 32 bit data word length and 31kb memory depth for image size of 640X480. The Virtex 4 Device offers a large number of Block RAMs with size of 18Kb. However, these blocks can be interconnected to build wider and deeper memory systems. The 18Kb BRAM is a dual port RAM with 18Kb memory space and two ports A and B are completely independent. Data can be written and read on both ports simultaneously and each port has its own data lines, address lines and control lines (See Fig. 4).

This interface uses BRAM as a memory and FPGA I/O pins as a physical connection between EMIF of DSP and FPGA. To implement 32kb X 32 bit dual port BRAM (see Fig. 4), a set of 8 BRAMs blocks are configured as a 32 bit wide and 32KB deep True dual port RAM memory. Port A is used as the access port for EMIF of DSP and Port B configured to act as contact port for FPGA.

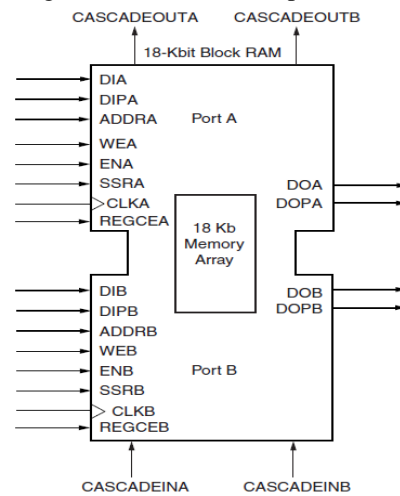


Fig. 4 18KB BRAM in Vertex 4 Device

4- Face Detection Theory & Algorithm

A hybrid face detection algorithm based on skin tone detection and Viola and Jones face detection structure implemented on Heterogeneous platform, which is on discussed in the previous section.

Face Detection Overview

In Recent days, Human face recognition plays critical role in the automation of various processes in this technologically advanced world. The fundamental step in

the human face recognition algorithm is, to detect whether image consists of face or not. If detected, the region of the human face is estimated (see Fig. 5). The difficulty with face detection greatly related with pretence, the existence of structural items, Facial expressions, Image orientation, Occlusion etc.

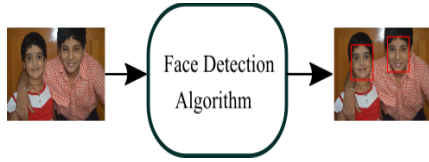


Fig. 5 Face Detection System

The still image face detection algorithms broadly divided into four distinct methods. a. Knowledge based methods: In this methods human knowledge about the face and what a typical face consists are used to define a convention to obtain relation between facial features (Yang, et al. 1994). b. Feature invariant methods: in these methods various structural features such as colour, shape, texture and other local feature, which are invariant even with changes in illumination condition, pose and viewpoint are used to detect faces in images (Leung, et al. 1995; Dai and Nakano, 1996; McKenna, et al. 1998; Kjeldsen and Kender, 1996). c. Template matching methods: un like other methods, in this various regular arrangements of human face collected and stored in the template database. These templates are used to correlate with input image to detect faces in the image (Craw, et al. 1992; Lanitis, et al. 1995). d. Appearance based methods: these type of algorithms relies on the large image database, which comprise a huge range of human faces with numerous variations. Support Vector Machine (SVM) (Osuna and Girosi, 1997) and Neural Networks (Rowley, et al. 1998) are the most commonly used techniques in this category

4-1- Face detection Algorithm

The proposed two stage face detection algorithm based on skin tone and Viola and Jones face detection structure is shown in Fig. 6. In the first stage, the input image, which is in the RGB colour model converted into YCrCb model and skin patches are segmented in the input image. The second stage, extract facial features from skin segmented image and detect faces in the image using Viola and Jones face detection algorithm.

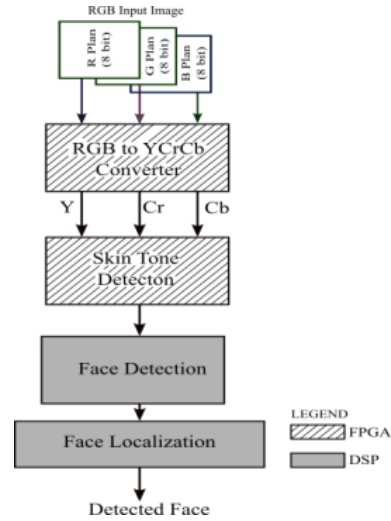


Fig. 6 Face Detection Algorithms

4-1-1 RGB to YCrCb Colour Model Conversion

Colour model is a mathematical representation of colour in terms of three or four components. The different colour models are used based on applications such as processing of digital image data, Display, transmission and TV broadcasting. There are several colour models are proposed and some most popular colour models are RGB, HSI, HSV, HSL, YIQ, YCrCb and YUV.

RGB (Red Green Blue) Colour model is most commonly used colour model to represent digital images. In this any colour is represented by three primary colours Red, Green and Blue based on how much percentage taken from each component. The skin tone detection based on the RGB colour model not preferred because of the high correlation between chrominance values and illumination value (Jones and Rehg, 2002) [12]. The normalized RGB can obtain from eq. [1-3].

$$r = \frac{R}{R + G + B} \tag{1}$$

$$g = \frac{G}{R + G + B} \tag{2}$$

$$b = \frac{B}{R + G + B} \tag{3}$$

In YCrCb colour model, Y components represent the luminance information and Cr and Cb represent chrominance information of image pixels. This colour model is most commonly used model because luminance and chrominance components are highly independent. YCrCb value can be obtained from the RGB model according to eq. 4.

$$\begin{bmatrix} Y \\ Cb \\ Cr \end{bmatrix} = \begin{bmatrix} 16 \\ 128 \\ 128 \end{bmatrix} + \begin{bmatrix} 0.279 & 0.504 & 0.098 \\ -0.148 & -0.291 & 0.439 \\ 0.439 & -0.368 & -0.071 \end{bmatrix} \begin{bmatrix} R \\ G \\ B \end{bmatrix} \quad (4)$$

4-1-2 Skin Tone Detection

Skin tone detection efficiency largely depends on choices of appropriate colour model and clustering of skin pixels in that colour model. We selected YCrCb colour model since it seems identical space and also adopted in very popular video compression standards such as MPEG and JPEG. Many researchers have reported that Cr and Cb values of skin pixels are uncorrelated with the Y value of the pixel. But, in practical, skin tone is nonlinearly dependent on the luma component. Some researchers demonstrated that detecting skin tone in CrCb and Cb/Y -Cr/Y sample subspaces results in many false positives and false negatives respectively. Therefore, the Rein-Lien et al., (2002) proposed a nonlinear transformation of YCrCb in order to make the skin colour space independent of the luma component (eq 5-10).

$$\begin{aligned} \text{MeanCr} &= 142.53 - Y * 0.091 & Y < 128 \\ &= 66 + Y * 0.468 & 129 \leq Y < 188 \\ &= 0 & \text{otherwise} \end{aligned} \quad (5)$$

$$\begin{aligned} \text{MeanCb} &= 119.46 - Y * 0.091 & Y < 128 \\ &= 68 + Y * 0.212 & 129 \leq Y < 188 \\ &= 0 & \text{otherwise} \end{aligned} \quad (6)$$

$$\begin{aligned} \text{WidthCr} &= 17.24 - Y * 0.172 & Y < 128 \\ &= 153.76 + Y * 0.611 & 129 \leq Y < 188 \\ &= 0 & \text{otherwise} \end{aligned} \quad (7)$$

$$\begin{aligned} \text{WidthCb} &= 19.48 - Y * 0.21 & Y < 128 \\ &= 178.85 + Y * 0.70 & 129 \leq Y < 188 \\ &= 0 & \text{otherwise} \end{aligned} \quad (8)$$

$$\begin{aligned} \text{TransCr} &= Cr & \text{if } 125 \leq Y \leq 188 \\ &= [Cr - \text{MeanCr}] * \left[\frac{38.76}{\text{WidthCr}} \right] + \text{MeanCr}(188) \end{aligned} \quad (9)$$

$$\begin{aligned} \text{TransCb} &= Cb & \text{if } 125 \leq Y \leq 188 \\ &= [Cb - \text{MeanCb}] * \left[\frac{38.76}{\text{WidthCb}} \right] + \text{MeanCb}(188) \end{aligned} \quad (10)$$

The transformed Chroma components [TransCr, TransCb] space that represent the elliptical model is described by eq. 11-12.

$$\begin{bmatrix} K \\ L \end{bmatrix} = \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix} \begin{bmatrix} \text{transCb} - c_x \\ \text{transCr} - c_y \end{bmatrix} \quad (11)$$

4-1-3 Face Detection using Adaboost

The basic structure proposed by Viola and Jones [13] is used to solve the face detection problem. In this framework, facial features are used instead of pixel-based operations. Many researchers reported very detailed versions of this algorithm; hence we are presenting very little information about this method. This approach is based on Haar-like features. Given a set of features and training based on face images and non-face images, the Adaboost algorithm can choose the best single Haar-like feature which isolates face and non-face images and strong classifiers are formed by combining weak classifiers. These strong classifiers are cascaded to detect faces in images. The number of stages required for face detection depends on accuracy and speed requirements. In the learning stage, after each round, weights of images which were correctly judged by the preceding weak classifier are enhanced.

5- Implementation

The proposed heterogeneous system consists of an FPGA and a DSP processor. The skin tone detection algorithm is implemented on the FPGA and the face detection algorithm on the DSP processor. The hardware connections of the proposed system are shown in Fig. 7.

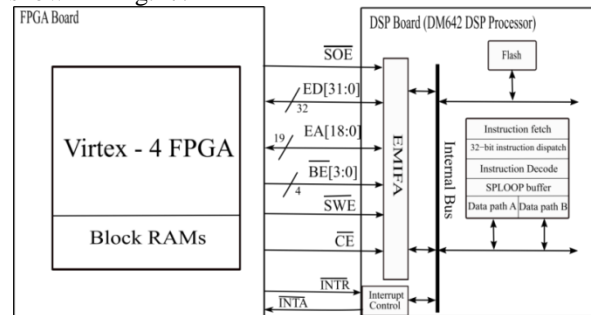


Fig. 7 Proposed System Architecture

5-1- Implementation of Skin Tone Detection

The proposed skin tone detection algorithm has been implemented on Xilinx's Virtex 4VSX35 FPGA. The input image is loaded into the Block RAM of the FPGA, and a Block RAM controller is realized to read the Block RAMs. The transformed Chroma values, transCr and transCb, are used to find the skin score for each pixel, as shown in Fig. 8. After processing, the skin-segmented image is loaded into dual-port RAM, which is later read by the Texas Instruments (TI) TMS320DM642 DSP processor using EMIFA to detect human faces in the skin-segmented image.

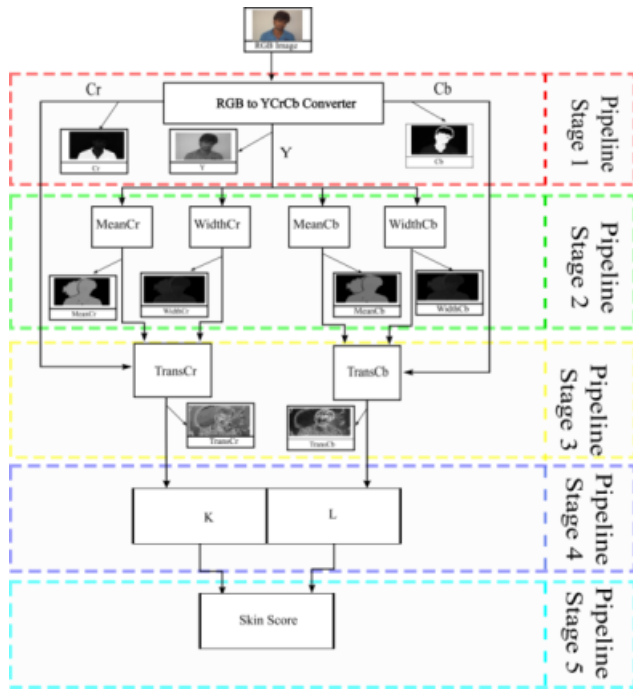


Fig. 8 Skin Tone Detection Hardware flow

5-2- DSP Implementation of Face Detection

We have chosen TMS320DM642 DSP platform to implement face detection algorithm, which offers a low cost solution for high performance requirements and it offers speed up to 4800 MIPS (Million Instruction per Second) at a clock frequency of 600 MHz.

In the heterogeneous system, TMS320DM642 DSP processor reads skin segmented image from dual port BRAM configured in FPGA using EMIF facility provided in TMS320DM642 DSP processor. The face detection algorithm based on Viola and Jones face detection structure is successfully implemented on the target TMS320DM642 DSP processor. The Face detection algorithm is implemented in MATLAB R2015a and converted to C code using the MATLAB coder facility and CCS (Code Composer Studio) project is implemented using generated code.

6- Experiment Results

The proposed system heterogeneous face detection system architecture consists of two modules. The first module read the RGB images from Block RAM and generates skin segmented images. The second module, read the skin segmented image from first module and localize the face in the skin regions of the image.

In feature classification stage, it is need to generate sum of pixels in a various rectangular area for haar feature

classifications. Different sizes of rectangular area require different time to compute the sum of the pixels in the rectangular area. An integral image used to reduce the computational time of sum of the pixels under rectangle. By using integral image, area under all sized rectangle are computed at constant time with two adders and one subtractor. The state of the art implementations convert entire image (640X420) it's required 270MB of RAM to hold the 640X420 size integral image. However, in the proposed system integral image generated only for sub-image (24X24) and it required 580 bytes of RAM to hold 24X24 sized integral image. It intern reduced the memory requirement of the overall system.

The performance of the complete system is tested by using Bio Database [Image Data Base, 2017], which includes images with single and multiple faces with different pose, orientation, occlusions and illumination.

The database images are resized to 640X480 and pixel size of 24 bits (8 bits for each RGB component). We have implemented proposed skin tone detector architecture on Xilinx's Virtex 4 VSX35 FPGA and results are presented . Table 1

Table 1 Synthesis Results of the proposed architecture implemented on Xilinx's Virtex 4VSX35 FPGA

Parameter	Value
No of clock cycles	264354
Execution time at 5Mhz clock	4.8ms
Execution time at 120Mhz clock	1.8ms
Hardware utilization	982 Logical elements
Core Power Dissipation	90.2mW

The detection performances of the proposed system, are presented in **Error! Reference source not found.** The proposed dual stage face detection system achieved a very high detection rate of 94.5% and performance of 13.1 FPS for image resolution of 640X480. The performance comparisons of our system with some existing systems are reported in Table 3

Table 2 face detection rates of proposed system

Type of the image (With variations)	No of images	No of faces in the image	Accuracy	
			No of detected faces	Detection rate
Single	150	150	139	92.66%
Group	220	1230	1173	95.36%

Table 3 Performance Comparison of proposed Architecture

Platform	Detection time(ms)	Frame Per Second (FPS)
FPGA+DSP (Proposed system)	112	13.0
FPGA (Fekih, et. al., 2015) [5]	160	6.18
DSP Processor (Zhao, et.al., 2009) [14]	940	1.07
PC (MATLAB)	1438	0.69

7- Conclusion

In this paper new hardware architecture for real face detection is presented. This implementation allows the user to choose image resolution and speed with available resources in the FPGA. The current implementation is based The Xilinx Starter Kit Virtex-4 SX35 Starter Kit and XEVM642 Development Kit powered by a TI TMS320DM642 DSP processor. The proposed system achieved 13.7 FPS average frame rate, when tested with an images with a spatial resolution of 640X480. This system exhibits performance improvement of 2.12 times compared with equivalent FPGA implementation, 12.3 times compared to DSP implementation and 18.98 times compared to PC implementation to solve the real time performance problem. The proposed hardware architecture achieved average detection accuracy of 94.5%, which low compared to the implementation on PC (97%), since the low accuracy pixel data are used in FPGA hardware architecture.

Future Scope

In the future, this hybrid architecture can be extended to design high performance facial recognition system by modifying the second stage of the proposed system.

Acknowledgments

Supported by VLSI Lab, Dept. of ECE, Shri Vishnu Engineering College for women

References

- [1] Y. Lei, Z. Gang, R. Si-Heon, Lee Choon-Young, Lee Sang-Ryong and K. -M. Bae, "The Platform of Image Acquisition and Processing System Based on DSP and FPGA," 2008 International Conference on Smart Manufacturing Application, 2008, pp. 470-473, doi: 10.1109/ICSMA.2008.4505567.
- [2] C. Kotropoulos and I. Pitas, "Rule-based face detection in frontal views," 1997 IEEE International Conference on Acoustics, Speech, and Signal Processing, 1997, pp. 2537-2540 vol.4, doi: 10.1109/ICASSP.1997.595305.
- [3] D. Nguyen, D. Halupka, P. Aarabi and A. Sheikholeslami, "Real-time face detection and lip feature extraction using field-programmable gate arrays," in IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics), vol. 36, no. 4, pp. 902-912, Aug. 2006, doi: 10.1109/TSMCB.2005.862728.
- [4] D. N. Arya, K. L. V. Sivanji, R. Reddy, S. Sivanantham, and K. Sivasankaran, "A face detection system implemented on FPGA based on RCT colour segmentation," *Proc. 2016 Online Int. Conf. Green Eng. Technol. IC-GET 2016*, 2017, doi: 10.1109/GET.2016.7916781.
- [5] H. Ben Fekih, A. E. B, and B. Juurlink, "An Efficient and Flexible FPGA Implementation of a Face Detection System," pp. 243-254, 2015, doi: 10.1007/978-3-319-16214-0.
- [6] H.-Y. Leung, L.-M. Cheng, and X. Y. Li, "A FPGA implementation of facial feature extraction," *J. Real-Time Image Process.*, vol. 10, no. 1, pp. 135-149, 2015, doi: 10.1007/s11554-012-0263-8.
- [7] A. S. Kamewar, "Processing geospatial images using GPU," 2017 International Conference on Emerging Trends & Innovation in ICT (ICEI), 2017, pp. 27-32, doi: 10.1109/ETIICT.2017.7977005.
- [8] J. Battle, "A New FPGA/DSP-Based Parallel Architecture for Real-Time Image Processing," *Real-Time Imaging*, vol. 8, no. 5, pp. 345-356, 2002, doi: 10.1006/rtim.2001.0273.
- [9] K. L. Y. Li *et al.*, "A new parallel particle filter face tracking method based on heterogeneous system," *J. Real-Time Image Process.*, vol. 7, no. 3, pp. 153-163, 2012, doi: 10.1007/s11554-011-0225-6.
- [10] L. Guo, "An embedded multimedia communication terminal based on DSP+FPGA," *Multimed. Tools Appl.*, vol. 76, no. 16, pp. 16949-16961, 2017, doi: 10.1007/s11042-016-3597-6.
- [11] Z. Ding, F. Zhao, T. Wang, W. Shu, and M.-Y. Wu, "Hecto-Scale Frame Rate Face Detection System for SVGA Source on FPGA Board," *2011 IEEE 19th Annu. Int. Symp. Field-Programmable Cust. Comput. Mach.*, pp. 37-40, 2011, doi: 10.1109/FCCM.2011.16.
- [12] Rein-Lien Hsu, M. Abdel-Mottaleb and A. K. Jain, "Face detection in color images," in IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 24, no. 5, pp. 696-706, May 2002, doi: 10.1109/34.1000242.
- [13] P. Viola and M. Jones, "Robust real-time face detection," *Int. J. Comput. Vis.*, vol. 57, no. 2, pp. 137-154, 2004, doi: 10.1023/B:VISI.0000013087.49260.fb.
- [14] F. Zhao, L. Yang, Y. Zhu, and P. Liao, "Enhancing the implementation of Adaboost algorithm on a DSP-based

- platform," *Int. Conf. Scalable Comput. Commun. - 8th Int. Conf. Embed. Comput. ScalCom-EmbeddedCom 2009*, pp. 393–395, 2009, doi: 10.1109/EmbeddedCom-ScalCom.2009.77.
- [15] Ganeswara Rao M.V., Panakala R.K., Mallikarjuna Prasad A. (2018) A New VLSI Architecture for Skin Tone Detection in an Uncontrolled Background. In: Anguera J., Satapathy S., Bhateja V., Sunitha K. (eds) *Microelectronics, Electromagnetics and Telecommunications. Lecture Notes in Electrical Engineering*, vol 471. Springer, Singapore
- [16] Fekih H.B., Elhossini A., Juurlink B. (2015) An Efficient and Flexible FPGA Implementation of a Face Detection System. In: Sano K., Soudris D., Hübner M., Diniz P. (eds) *Applied Reconfigurable Computing. ARC 2015. Lecture Notes in Computer Science*, vol 9040. Springer, Cham
- [17] Dong Zhang, S. Z. Li and D. Gatica-Perez, "Real-time face detection using boosting in hierarchical feature spaces," *Proceedings of the 17th International Conference on Pattern Recognition*, 2004. ICPR 2004., Cambridge, 2004, pp. 411-414 Vol.2.
- [18] Y. N. Chae, T. Han, Y.-H. Seo, and H. S. Yang, "An efficient face detection based on color-filtering and its application to smart devices," *Multimed. Tools Appl.*, vol. 75, no. 9, pp. 4867–4886, 2016, doi: 10.1007/s11042-013-1786-0.
- [19] C. Kumar and M. S. Azam, "A multi-processing architecture for accelerating Haar-based face detection on FPGA," *9th Int. Conf. Ind. Inf. Syst. ICIIS 2014*, 2015, doi: 10.1109/ICIINFS.2014.7036525.
- [20] S. Liao, A. K. Jain, and S. Z. Li, "A Fast and Accurate Unconstrained Face Detector," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 2, pp. 211–223, 2016, doi: 10.1109/TPAMI.2015.2448075.
- [21] A. N. Rajagopalan, K. S. Kumar, J. Karlekar, R. M. M. Patil, U. B. Desai, and P. G. P. S. Chaudhuri, "Finding Faces in Photographs," *IEEE Int. Conf. Comput. Vis.*, no. 1, pp. 640–645, 1998, doi: 10.1109/ICCV.1998.710785.
- [22] M. S. Lew, "Information theoretic view-based and modular face detection," *Proceedings of the Second International Conference on Automatic Face and Gesture Recognition*, Killington, VT, USA, 1996, pp. 198-203.
- [23] A. J. Colmenarez and T. S. Huang, "Face Detection With Information- Based Maximum Discrimination," *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, pp. 782–787, 1997, doi: <http://dx.doi.org/10.1109/CVPR.1997.609415>.
- [24] K. S. Park, R. H. Park, and Y. G. Kim, "Face detection using the 3x3 block rank patterns of gradient magnitude images and a geometrical face model," *Dig. Tech. Pap. - IEEE Int. Conf. Consum. Electron.*, no. c, pp. 793–794, 2011, doi: 10.1109/ICCE.2011.5722867.
- [25] P. P. Paul and M. Gavrilova, "PCA based geometric modeling for automatic face detection," *Proc. - 2011 Int. Conf. Comput. Sci. Its Appl. ICCSA 2011*, pp. 33–38, 2011.
- [26] A. Majumder, L. Behera, and V. K. Subramanian, "Automatic and robust detection of facial features in frontal face images," *Proc. - 2011 UKSim 13th Int. Conf. Model. Simulation, UKSim 2011*, pp. 331–336, 2011, doi: 10.1109/UKSIM.2011.69.
- [27] J. Guo, C. Lin, M. Wu, C. Chang and H. Lee, "Complexity Reduced Face Detection Using Probability-Based Face Mask Prefiltering and Pixel-Based Hierarchical-Feature Adaboosting," in *IEEE Signal Processing Letters*, vol. 18, no. 8, pp. 447-450, Aug. 2011.
- [28] Katkooari Arun Kumar and Ravi Boda, "A Threshold-based Brain Tumour Segmentation from MR Images using Multi-Objective Particle Swarm Optimization," *Journal of Information Systems and Telecommunication*, Vol. 9, No. 4, 2021, pp. 218–225.
- [29] Hamed Agahi and Kimia Rezaei, "An Automatic Thresholding Approach to Gravitation-Based Edge Detection in Grey-Scale Images," *Journal of Information Systems and Telecommunication*, Vol. 9, No. 4, 2021, pp. 285–296.
- [30] K. Li, Y. Tian, B. Wang, Z. Qi, and Q. Wang, "Bi-Directional Pyramid Network for Edge Detection," *Electronics*, vol. 10, no. 3, 2021, p. 329-333.
- [31] D. Wang, J. Yin, C. Tang, X. Cheng, and B. Ge, "Color edge detection using the normalization anisotropic Gaussian kernel and multichannel fusion," *IEEE Access*, vol. 8, 2020, pp. 228277-228288,.
- [32] Azamossadat Nourbakhsh, Mohammad-Shahram Moin and Arash Sharifi, "Facial Images Quality Assessment based on ISO/ICA0 Standard Compliance Estimation by HMAX Model," *Journal of Information Systems and Telecommunication*, Vol. 7, No. 27, 2009, pp. 225–237.
- [33] Azar Mahmoodzadeh, "Human Activity Recognition based on Deep Belief Network Classifier and Combination of Local and Global Features," *Journal of Information Systems and Telecommunication*, Vol. 9, No. 36, 2021, pp. 45–52.

A Novel Detector based on Compressive Sensing for Uplink Massive MIMO Systems

Mojtaba Amiri¹, Amir Akhavan^{2*}

¹. School of Electrical and Computer Engineering, College of Engineering, University of Tehran, Tehran, Iran

². Department of Electrical and Computer Engineering, Isfahan University of Technology, Isfahan 84156-83111, Iran

Received: 19 Feb 2022/ Revised: 07 Apr 2022/ Accepted: 13 Apr 2022

Abstract

Massive multiple-input multiple-output is a promising technology in future communication networks where a large number of antennas are used. It provides huge advantages to the future communication systems in data rate, the quality of services, energy efficiency, and spectral efficiency. Linear detection algorithms can achieve a near-optimal performance in large-scale MIMO systems, due to the asymptotic orthogonal channel property. But, the performance of linear MIMO detectors degrades when the number of transmit antennas is close to the number of receive antennas (loaded scenario). Therefore, this paper proposes a series of detectors for large MIMO systems, which is capable of achieving promising performance in loaded scenarios. The main idea is to improve the performance of the detector by finding the hidden sparsity in the residual error of the received signal. At the first step, the conventional MIMO model is converted into the sparse model via the symbol error vector obtained from a linear detector. With the aid of the compressive sensing methods, the incorrectly detected symbols are recovered and performance improvement in the detector output is obtained. Different sparse recovery algorithms have been considered to reconstruct the sparse error signal. This study reveals that error recovery by imposing sparse constraint would decrease the bit error rate of the MIMO detector. Simulation results show that the iteratively reweighted least squares method achieves the best performance among other sparse recovery methods.

Keywords: Massive MIMO; MMSE Detector; Error Recovery; Compressive Sensing; Iteratively Reweighted Least Squares (IRLS) Method.

1- Introduction

The number of cellular phones and mobile data traffic are extremely growing each year. Telecommunication companies are asked to provide higher data rates, further spectral efficiency, and larger capacity. The fifth-generation (5G) wireless communication systems are being designed to answer excessive data rate demands. Massive MIMO technology is a good candidate for the next-generation of the wireless communication. The base station (BS) in massive MIMO systems equipped with hundreds of antennas are used to serve tens of users simultaneously [1, 2]. But there are some challenges such as pilot contamination, detection performance, channel estimation and detection complexity [3-5].

The purpose of each detection algorithm is to obtain an estimate of the transmit signal, given knowledge of the received signal and the channel state information (CSI). The maximum a posteriori (MAP) and the maximum

likelihood (ML) algorithms provide the optimal detectors but they are not practically feasible for the massive MIMO systems since their computational complexity increase exponentially with the number of antennas. Linear MIMO detectors such as zero forcing (ZF) and minimum mean square error (MMSE) receivers can achieve near optimal performance when the number of users is much lower than the number of the antennas in BS [6]. Many methods have been proposed to achieve the performance of MMSE detector with low complexity such as the optimized coordinate descent (OCD) [7], Gauss-Seidel (GS) method [8], parallelizable Chebyshev iteration (PCI) [9] and alternating minimization method (Alt-Min) [10]. The performance of the linear detectors and also the previously mentioned methods degrade when the number of transmitters is close to the number of receiver antennas [11]. Therefore, new and efficient detectors with a low error rate are highly needed to solve this problem. This paper focuses on developing a detector which achieves favorable performance in loaded scenarios.

Recently, compressive sensing (CS) and sparse signal

recovery techniques have received much attention in different signal processing applications. Compressive sensing has emerged as a promising approach for use in large MIMO systems [12, 13].

It is noteworthy that the original signals in massive MIMO systems are not intrinsically sparse, but it is expected that the detector output contains an error only for a few number of users. Thus, the error vector resulting from a primary estimator is likely to be sparse, especially in high SNR regime. The motivation of this paper is to improve the performance of the detector by using the sparsity in the residual error of large MIMO systems. In order to exploit the sparsity of the detection errors, the conventional model is converted into a sparse model via the symbol error vector [13, 14]. After that, the error recovery algorithm can be performed to improve the detection performance by recovering the non-zero entries of the error vector.

Sparse signal recovery is basically an optimization problem with l_0 -norm and is NP-hard. Therefore, different approaches are proposed to solve this problem. Greedy methods [15-17], l_1 -relaxation based optimization [18, 19] and Bayesian methods [20, 21] are the main approaches to estimate the sparse vector. Many different algorithms have been proposed for sparse signal reconstruction. The contribution of this paper is to address the effectiveness of the l_1 -relaxation-based sparse recovery methods in massive MIMO detectors for the first time.

The rest of the paper is organized as follows. Section II introduces the system model of the massive MIMO system. Section III presents the conversion of the conventional MIMO system model into the sparse error domain. Sparse recovery algorithms are introduced in section IV. The simulation results and discussions about the performance of the proposed algorithm are presented in section V and finally the paper is concluded in Section VI.

Notation: Boldface capital letters and lowercase letters represent matrices and vectors, respectively. \mathbf{I}_K denotes the $K \times K$ identity matrix; $(\cdot)^H$, $(\cdot)^{-1}$, $(\cdot)^T$ denote the conjugate transposition, the inversion and the transposition, respectively. $\mathbb{C}^{m \times n}$ denotes the $m \times n$ complex matrix. The p -norm (also called l_p -norm) of vector $\mathbf{v} = (v_1, v_2, \dots, v_n)$ is $\|\mathbf{v}\|_p = (\sum_{i=1}^n |v_i|^p)^{1/p}$.

2- System Model

Consider a multi-user MIMO model with n_t users and n_r receivers in the BS. The received signal can be described as:

$$\mathbf{y} = \mathbf{H}\mathbf{x} + \mathbf{n} \quad (1)$$

where $\mathbf{H} \in \mathbb{C}^{n_r \times n_t}$ is the channel matrix between the BS and the n_t users whose entries are modeled as independent and identically distributed (iid) complex Gaussian with zero mean and unit variance. $\mathbf{x} \in \mathbb{C}^{n_r \times 1}$, is the complex-valued information vector, and \mathbf{n} is a white Gaussian noise vector with zero mean and correlation matrix $E(\mathbf{n}\mathbf{n}^H) = \sigma_n^2 \mathbf{I}_{n_r}$. It is assumed that the channel matrix is known perfectly at the BS but it is unknown at the transmitter side.

2-1- Linear Detection

The maximum likelihood (ML) detector is not suitable for solving large dimensional problems due to the high computational complexity. Therefore, suboptimal detectors such as Minimum Mean Square Error with low complexity are beneficial in operational conditions.

The MMSE detector can be obtained by the solution of the following minimization problem.

$$\min_{\hat{\mathbf{x}}} E[|\mathbf{x} - \hat{\mathbf{x}}|^2] \quad (2)$$

where $\hat{\mathbf{x}} = \mathbf{G}\mathbf{y}$ is the estimation of the user's data.

3- Massive MIMO Detection in Sparse Domain

This section, presents a class of detectors based on error recovery technique for detection of the transmitted symbols in uplink massive MIMO system. This method iteratively achieves near-optimal performance in terms of bit error rate. In the following the error domain sparse model is introduced.

3-1- Sparse Model

At the first step, the conventional system model (1) should be converted into a sparse model via the symbol error vector obtained from a linear detector.

The error vector \mathbf{x}_e , is defined as the difference between the original signal and the recovered one. Therefore, the system model in error domain can be formulated as follows:

$$\mathbf{y}_e = \mathbf{y} - \mathbf{H}\hat{\mathbf{x}}_{\text{MMSE}} = \mathbf{H}(\mathbf{x} - \hat{\mathbf{x}}_{\text{MMSE}}) + \mathbf{n} = \mathbf{H}\mathbf{x}_e + \mathbf{n} \quad (3)$$

where $\mathbf{x}_e = \mathbf{x} - \hat{\mathbf{x}}_{\text{MMSE}}$ is the error vector of the primary estimated symbols and its nonzero values correspond to the incorrectly detected symbols.

It is noteworthy that, by recovering the incorrectly detected symbols, the performance of the detector can be improved. Since it is expected that only a few symbols are

incorrectly detected, \mathbf{x}_e , is a sparse vector. Therefore, the detection operation is equivalent to recover the sparse error vector, \mathbf{x}_e , from the difference signal, \mathbf{y}_e .

Once the sparse error vector is recovered, the estimation of the transmitted signal, $\hat{\mathbf{x}}$, is obtained by adding the error vector to the initial estimate, $\hat{\mathbf{x}}_{\text{MMSE}}$. Thus, the final estimation of the user's data is obtained by

$$\hat{\mathbf{x}} = \hat{\mathbf{x}}_{\text{MMSE}} + \hat{\mathbf{x}}_e \quad (4)$$

3-2- Error Recovery

The problem of sparse representation in the MIMO detection is to find the vector \mathbf{x}_e . Therefore, we are looking for the sparsest solution which can be done by solving the following optimization problem:

$$P_0: \quad \min \|\mathbf{x}_e\|_0 \quad s.t. \quad \mathbf{y}_e = \mathbf{H}\mathbf{x}_e \quad (5)$$

Where $\|\mathbf{x}_e\|_0$ denotes the l_0 -norm of \mathbf{x}_e and gives the total number of non-zero elements in the vector. Since (P_0) is NP-hard, the optimization problem is relaxed with convex l_1 -norm. Taking into account the effect of the noise component, the problem (P_0) can be converted to the following optimization problem:

$$P_1: \quad \min \|\mathbf{x}_e\|_1 \quad s.t. \quad \|\mathbf{y}_e - \mathbf{H}\mathbf{x}_e\|_2 \leq \varepsilon \quad (6)$$

It is assumed that the noise has bounded entries, i.e. $\|\mathbf{n}\|_2 \leq \varepsilon$ for some sufficiently small ε . Additionally, according to the Lagrange multiplier theorem, there exists an appropriate constant λ such that the problem (P_1) is equivalent to the following unconstrained minimization problem.

$$P_2: \quad \min \lambda \|\mathbf{x}_e\|_1 + \frac{1}{2} \|\mathbf{y}_e - \mathbf{H}\mathbf{x}_e\|_2^2 \quad (7)$$

Where the Lagrange multiplier λ depends on \mathbf{y}_e and ε . Note that the cost function in (P_2) is not differentiable with respect to \mathbf{x}_e and specific optimization algorithms are required to solve (P_2) . The following section addresses three well-known sparse coding algorithms to estimate the error vector \mathbf{x}_e . The minimization function in (P_2) is composed of two parts. The first term with l_1 -norm induces sparsity to the estimated error vector, while the second term, $\frac{1}{2} \|\mathbf{y}_e - \mathbf{H}\mathbf{x}_e\|_2^2$, makes the estimated vector consistent with \mathbf{y}_e . In order to investigate the effectiveness of the sparsity promoting term in (P_2) , the results of the MIMO detection with the following minimization problem are also considered.

$$P_3: \quad \min \lambda \|\mathbf{x}_e\|_2 + \frac{1}{2} \|\mathbf{y}_e - \mathbf{H}\mathbf{x}_e\|_2^2 \quad (8)$$

where the l_1 -norm in (P_2) is replaced with the l_2 -norm. The closed-form solution of the convex minimization problem (P_3) can be formulated as follows:

$$\hat{\mathbf{x}}_e = (2\lambda\mathbf{I} + \mathbf{H}^H\mathbf{H})^{-1}\mathbf{H}^H\mathbf{y}_e \quad (9)$$

In the simulation result section, the solution to (P_3) is called the regularized least square (RLS) estimation.

4- Sparse Error Reconstruction

The most significant stage in error recovery-based MIMO detection is the sparse error reconstruction. In this study, three algorithms are considered for the sparse coding step. More explicitly, Iterative Re-weighted Least Squares (IRLS), Alternating Direction Method of Multipliers (ADMM), and Iterative Shrinkage-Thresholding Algorithm (ISTA) [22-25] have been applied to reconstruct the error vector. In the following, these three algorithms are introduced briefly.

4-1- IRLS Algorithm

The Iterative Re-weighted Least Squares algorithm is one of the strategies which is able to recover sparse signals. In this algorithm, the l_1 -norm in (P_2) is replaced by a weighted l_2 -norm [26]:

$$P_3: \quad \min \lambda \mathbf{x}_e^T \mathbf{E}^{-1} \mathbf{x}_e + \frac{1}{2} \|\mathbf{y}_e - \mathbf{H}\mathbf{x}_e\|_2^2 \quad (10)$$

Where \mathbf{E} is a diagonal weight matrix and it is updated from the current iterate $(\mathbf{x}_e)_k$.

The minimization in $(P3)$ is a quadratic optimization problem, soluble using linear algebra. The pseudo-code for the IRLS error recovery-based MIMO detector has been shown in Algorithm 1.

4-2- ADMM Algorithm

The alternating direction method of multipliers is an alternative algorithm for sparse coding. This algorithm uses the augmented Lagrangian to splits the main optimization problem into two quadratic and separable minimization problems.

In this method, the augmented Lagrangian is defined as [27]

$$L_\mu(\mathbf{x}_e, \mathbf{z}, \lambda_a) = \lambda \|\mathbf{x}_e\|_1 + \frac{1}{2} \|\mathbf{z}\|_2^2 - \langle \lambda_a, \mathbf{y}_e - \mathbf{H}\mathbf{x}_e - \mathbf{z} \rangle + \frac{\mu}{2} \|\mathbf{y}_e - \mathbf{H}\mathbf{x}_e - \mathbf{z}\|_2^2 \quad (11)$$

Algorithm 1: IRLS error recovery-based detector**Input:** \mathbf{y} , \mathbf{H} , and σ_n^2 **Parameters:** maximum iteration number (K), threshold (η)**Output:** The estimation of the transmitted symbols: $\hat{\mathbf{x}}$ **initialization:**

- 1: $\mathbf{A} = \mathbf{H}^H \mathbf{H} + \sigma_n^2 \mathbf{I}_{n_t}$, $\mathbf{D} = \text{diag}(\mathbf{A})$ and $\mathbf{y}_{MF} = \mathbf{H}^H \mathbf{y}$
- 2: $\hat{\mathbf{x}}_{MMSE} = \mathbf{A}^{-1} \mathbf{y}_{MF}$ ‘Primary Estimation’
- 3: $\mathbf{y}_e = \mathbf{y} - \mathbf{H} \hat{\mathbf{x}}_{MMSE}$
- 4: The initial weight matrix $\mathbf{E} = \mathbf{D}$

Iteration: Increase k

- 5: **Regularized Least-Squares:** approximately solve the linear system

$$(2\lambda \mathbf{E}^{-1} + \mathbf{H}^H \mathbf{H})(\mathbf{x}_e)_k = \mathbf{H}^H \mathbf{y}_e$$

- 6: **Weight Update:** Update the diagonal weight matrix \mathbf{E}

$$\mathbf{E} = \text{diag}(|(\mathbf{x}_e)_k| + \varepsilon)$$

- 7: **Stopping Rule:** if $\|(\mathbf{x}_e)_k - (\mathbf{x}_e)_{k-1}\| < \eta$ break else go back to step 5

8: **Output:** $(\mathbf{x}_e)_k$.9: **return** $\hat{\mathbf{x}} = \hat{\mathbf{x}}_{MMSE} + (\mathbf{x}_e)_k$

where λ_a is the Lagrangian multiplier and $\mu > 0$ is a penalty parameter. The pseudo-code of the MIMO detection using the ADMM algorithm has been shown in Algorithm 2.

4-3- ISTA Algorithm

Another algorithm which can be used for solving problem (P2) is the iterative shrinkage-thresholding algorithms (ISTA). The solution based on the ISTA algorithm can be written as [23, 28]

$$\boldsymbol{\psi}_k = (\mathbf{x}_e)_{k-1} - 2t\mathbf{H}^H(\mathbf{H}(\mathbf{x}_e)_{k-1} - \mathbf{y}_e) \quad (12)$$

where t is the step size and the error vector is updated as

$$(\mathbf{x}_e)_k = \mathcal{H}(\boldsymbol{\psi}_k) \quad (13)$$

where $\mathcal{H}(\cdot)$ is the shrinkage operator and is described by

$$\mathcal{H}(\boldsymbol{\psi}_k) = \max(\mathbf{0}, |\boldsymbol{\psi}_k| - \alpha) \circ \text{sgn}(\boldsymbol{\psi}_k) \quad (14)$$

where \circ and $\text{sgn}(\cdot)$ are the Schur product and the sign function respectively. The parameter $\alpha = \lambda t$ represents the threshold value and λ is the proper scale hyperparameter.

Algorithm 2: ADMM error recovery-based detector**Input:** \mathbf{y} , \mathbf{H} , and σ_n^2 **Parameters:** maximum iteration number (K), threshold (η), penalty parameter (μ)**Output:** The estimation of the transmitted symbols: $\hat{\mathbf{x}}$ **initialization:**

- 1: $\mathbf{A} = \mathbf{H}^H \mathbf{H} + \sigma_n^2 \mathbf{I}_{n_t}$, $\mathbf{D} = \text{diag}(\mathbf{A})$ and $\mathbf{y}_{MF} = \mathbf{H}^H \mathbf{y}$
- 2: $\hat{\mathbf{x}}_{MMSE} = \mathbf{A}^{-1} \mathbf{y}_{MF}$ ‘Primary Estimation’
- 3: $\mathbf{y}_e = \mathbf{y} - \mathbf{H} \hat{\mathbf{x}}_{MMSE}$

Iteration: Increase k

- 4: **Update error vector, \mathbf{x}_e :**

$$(2\mathbf{H}^H \mathbf{H} + \mu \mathbf{I})(\mathbf{x}_e)_k = 2\mathbf{H}^H \mathbf{y}_e + \mu \mathbf{z}^{k-1} + \lambda_a^{k-1}$$

- 5: **Update \mathbf{z}^k :** compute \mathbf{z}^k via soft shrinkage

$$\mathbf{z}^k = \arg \min_{\mathbf{z}} \lambda \|\mathbf{z}\|_1 - \langle \lambda_a^{k-1}, (\mathbf{x}_e)_k - \mathbf{z} \rangle + \frac{\mu}{2} \|(\mathbf{x}_e)_k - \mathbf{z}\|_F^2$$

- 6: **Update Lagrangian multiplier, λ_a :** $\lambda_a^k = \lambda_a^{k-1} - \mu((\mathbf{x}_e)_k - \mathbf{z}^k)$

- 7: **Stopping Rule:** if $\|(\mathbf{x}_e)_k - (\mathbf{x}_e)_{k-1}\| < \eta$ break else go back to step 4

8: **Output:** $(\mathbf{x}_e)_k$.9: **return** $\hat{\mathbf{x}} = \hat{\mathbf{x}}_{MMSE} + (\mathbf{x}_e)_k$

The pseudo-code of the MIMO detection using the ISTA algorithm has been shown in Algorithm 3.

5- Simulation Result

In this section, numerical simulation results and complexity of detectors are presented to demonstrate the performance of the proposed methods. The simulations are conducted for $n_r \times n_t$ MIMO system, where n_r and n_t are the number of receive and transmit antennas, respectively. In the simulations, the massive MIMO system with 4-QAM and 16-QAM modulations are considered. Each entry of the channel matrix \mathbf{H} is an i.i.d. circularly symmetric complex Gaussian random variable (i.e., $\mathbf{H} \sim \mathcal{N}(0,1)$) and the channel statistics information is available for the BS and satisfy

$$\lim_{n_r \rightarrow \infty} \frac{1}{n_r} \mathbf{h}_n^H \mathbf{h}_n = 1 \quad (15)$$

where \mathbf{h}_n is the n th column of the matrix \mathbf{H} .

Algorithm 3: ISTA error recovery based detector**Input:** \mathbf{y} , \mathbf{H} , and σ_n^2 **Parameters:** maximum iteration number (K), threshold (η),**Output:** The estimation of the transmitted symbols: $\hat{\mathbf{x}}$ **initialization:**1: $\mathbf{A} = \mathbf{H}^H \mathbf{H} + \sigma_n^2 \mathbf{I}_{n_t}$, $\mathbf{D} = \text{diag}(\mathbf{A})$ and $\mathbf{y}_{MF} = \mathbf{H}^H \mathbf{y}$ 2: $\hat{\mathbf{x}}_{MMSE} = \mathbf{A}^{-1} \mathbf{y}_{MF}$ ‘Primary Estimation’3: $\mathbf{y}_e = \mathbf{y} - \mathbf{H} \hat{\mathbf{x}}_{MMSE}$ **Iteration:** Increase k 4: **Update** $\boldsymbol{\psi}_k$:

$$\boldsymbol{\psi}_k = (\mathbf{x}_e)_{k-1} - 2t\mathbf{H}^H(\mathbf{H}(\mathbf{x}_e)_{k-1} - \mathbf{y}_e)$$

5: **Update error vector**, \mathbf{x}_e :

$$(\mathbf{x}_e)_k = \mathcal{H}(\boldsymbol{\psi}_k)$$

6: **Stopping Rule:** if $\|(\mathbf{x}_e)_k - (\mathbf{x}_e)_{k-1}\| < \eta$ break else go back to step 47: **Output:** $(\mathbf{x}_e)_k$.8: **return** $\hat{\mathbf{x}} = \hat{\mathbf{x}}_{MMSE} + (\mathbf{x}_e)_k$

All simulations are carried out in Matlab 2015b on a processor Intel(R) Core (TM) i5-6200U CPU at 2.30 GHz and 8GB RAM and all results are averaged over 10000 iterations.

Prior to apply the minimization problems (P_2) or (P_3) for the MIMO detection, the coefficient λ should be adjusted. Fig. 1 shows the BER of the MIMO detectors for ISTA, ADMM, and IRLS algorithms with respect to different values of λ . In this simulation, the SNR has been fixed at 15 dB and the parameter λ varies from 0 to 50. The simulations are conducted with $n_r = n_t = 64$ for 4-QAM modulation. According to this figure, the values of the parameters λ in the following simulations are set to $\lambda_{ADMM} = 30$, $\lambda_{IRLS} = 5$ And $\lambda_{ISTA} = 17$.

Fig. 2 (a) shows the error of the primary detector in an uplink massive MIMO system with 16-QAM modulation with $n_r = n_t = 64$. This simulation shows that the error of the estimated user symbols is sparse. Fig. 2 (b), (c) and (d) illustrate the recovered error vector using the ADMM, ISTA, and IRLS respectively. It can be seen that only the error corresponding to the 40th user is not completely recovered. To further investigate the performance of different error recovery methods, various detection scenarios are simulated.

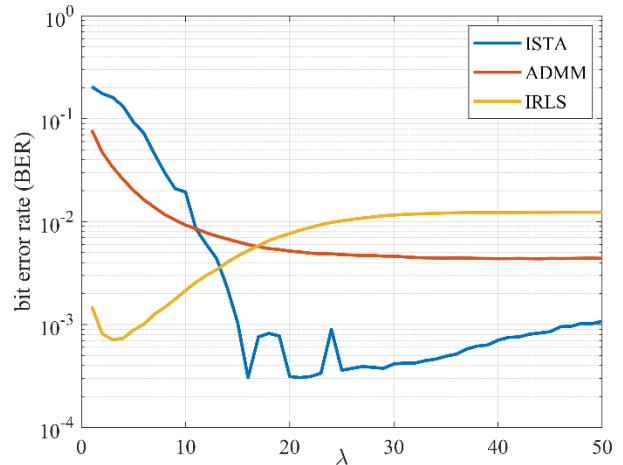


Fig. 1 BER performance versus λ in the uplink massive MIMO for 4-QAM modulation with SNR = 15 dB.

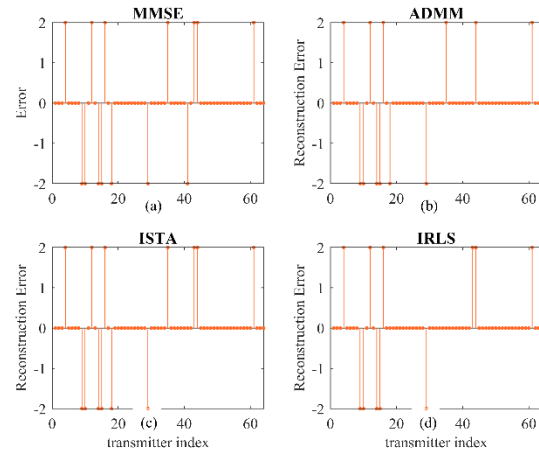


Fig. 2 (a) the error of the primary detector and (b), (c), (d) the recovered error vector using the ADMM, ISTA, and IRLS algorithms in the uplink massive MIMO system for 16-QAM modulation for $n_r = n_t = 64$ with SNR = 15 dB.

Fig.3-Fig. 6 shows the bit error rate (BER) of the MIMO detection for $n_r = n_t \in \{32, 64\}$ and $\{4, 16\}$ -QAM modulations. In Fig. 3 and Fig. 4, 4-QAM constellation with $n_r = n_t = 32$ and $n_r = n_t = 64$ are considered respectively. In comparison to the MMSE detector, performance improvement of the error recovery methods are markedly evident. It can be seen that the IRLS method has the best performance among other error recovery methods. In addition, all sparsity-based error recovery methods lead to lower BER in comparison to the RLS.

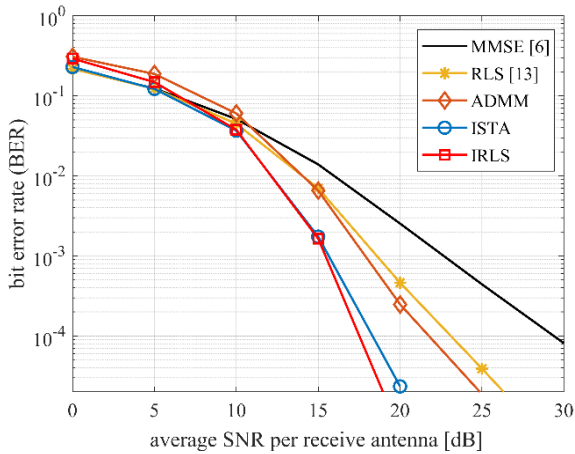


Fig. 3 BER performance comparison of the error recovery algorithms in the uplink massive MIMO for 4-QAM modulation for $n_r = n_t = 32$.

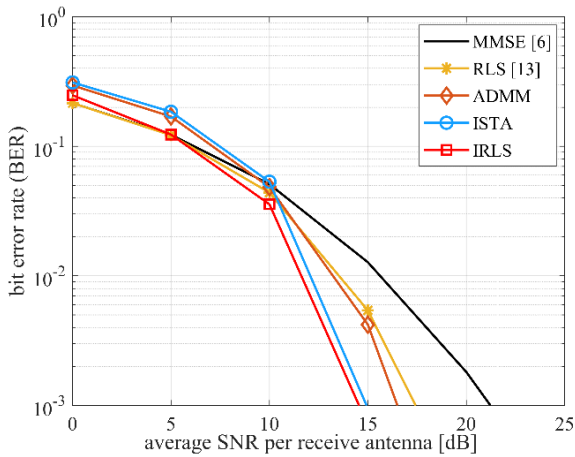


Fig. 4 BER performance comparison of the error recovery algorithms in the uplink massive MIMO for 4-QAM modulation for $n_r = n_t = 64$.

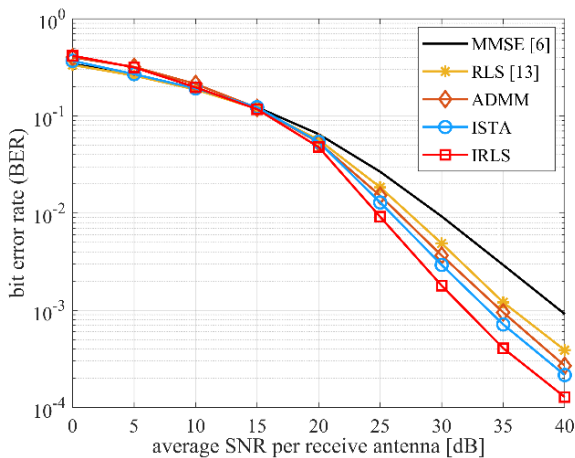


Fig. 5 BER performance comparison of the error recovery algorithms in the uplink massive MIMO for 16-QAM modulation for $n_r = n_t = 32$.

Fig. 5 and Fig. 6 show the detection performance for 16-QAM modulation with $n_r = n_t = 32$ and $n_r = n_t = 64$ respectively. Although the detection improvement is decreased but still all error recovery detectors achieved better performance than the MMSE detector.

Fig.7 compares the run time of the previously mentioned error recovery algorithms for different number of transmitters and $n_r = 64$.

The times are averaged over 10000 iterations. It can be seen that the run time of the all methods increase with the system dimensions. Generally, the run time of the IRLS method is less than that of the ADMM and ISTA algorithms. Since the RLS method has a close form solution, it leads to the least run time.

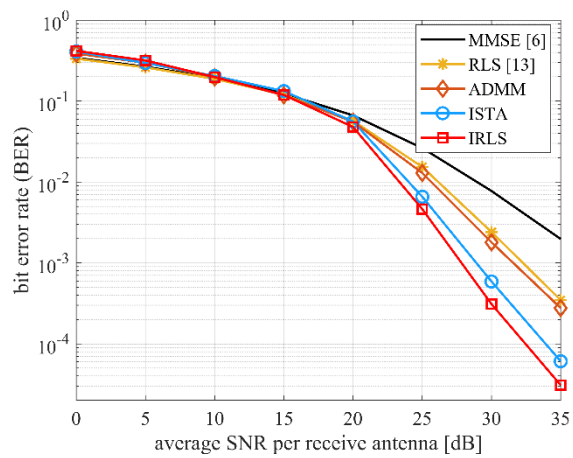


Fig. 6 BER performance comparison of the error recovery algorithms in the uplink massive MIMO for 16-QAM modulation for $n_r = n_t = 64$.

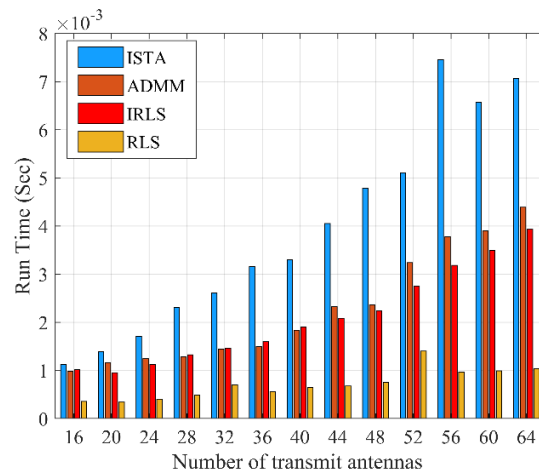


Fig. 7 Run time evaluation for the error recovery algorithms versus number of transmit antennas in the uplink massive MIMO for 64-QAM modulation with $n_r = 64$ and SNR = 10 dB.

The total computational complexity of the methods can be analyzed with respect to the number of multiplications in the Big-O notation. Since in the simulations n_t is close to n_r , it can be easily shown that the computational complexity of all methods is of order $O(n_t^3)$ which is similar to that of the MMSE MIMO detector. In order to summarize the results, it was demonstrated that the IRLS method leads to the best MIMO detection performance. Note that, since the IRLS method is an iterative algorithm and also it requires the matrix inversion operation in each iteration, the run time of the proposed algorithm is more than that of the MMSE detector. Applying the approximation methods in matrix inversion computation such as Gauss-Seidel, Chebyshev, and conjugate gradient methods would decrease the run time of the IRLS sparse recovery method.

The performance of the large-scale MIMO systems depends on the accuracy of the channel state information (CSI). In future works, an algorithm for joint channel estimation and signal detection in sparse error domain would be considered.

6- Conclusions

This paper focused on the problem of detection in massive MIMO systems. The main idea of this algorithm is to improve the performance of the detector by finding the hidden sparsity in the residual error of the received signal. In this paper, three sparse recovery algorithms, i.e. Iterative Re-weighted Least Squares (IRLS), Alternating Direction Method of Multipliers (ADMM), and Iterative Shrinkage-Thresholding Algorithm (ISTA) have been applied to reconstruct the error of the primary detector. It is noteworthy that the iteratively reweighted least-squares (IRLS) method achieved the best performance among other sparse recovery methods. The proposed methods outperform the MMSE detector but it is obvious that the complexity of the sparse error recovery-based MIMO detectors is more than that of the MMSE detector. Consequently, more efforts are needed to decrease the computational burden of the sparse error recovery algorithms.

References

- [1] E. G. Larsson, O. Edfors, F. Tufvesson, and T. L. Marzetta, "Massive MIMO for next generation wireless systems," *IEEE communications magazine*, vol. 52, no. 2, pp. 186-195, 2014.
- [2] E. Björnson, E. G. Larsson, and T. L. Marzetta, "Massive MIMO: Ten myths and one critical question," *IEEE Communications Magazine*, vol. 54, no. 2, pp. 114-123, 2016.
- [3] E. Björnson, J. Hoydis, M. Kountouris, and M. Debbah, "Massive MIMO systems with non-ideal hardware: Energy efficiency, estimation, and capacity limits," *IEEE Transactions on Information Theory*, vol. 60, no. 11, pp. 7112-7139, 2014.
- [4] M. A. Albreem, M. Juntti, and S. Shahabuddin, "Massive MIMO detection techniques: a survey," *IEEE Communications Surveys & Tutorials*, vol. 21, no. 4, pp. 3109-3132, 2019.
- [5] M. A. Albreem et al., "Low complexity linear detectors for massive MIMO: A comparative study," *IEEE Access*, vol. 9, pp. 45740-45753, 2021.
- [6] L. Fang, L. Xu, and D. D. Huang, "Low Complexity Iterative MMSE-PIC Detection for Medium-Size Massive MIMO," *IEEE Wireless Communications Letters*, vol. 5, no. 1, pp. 108-111, 2016.
- [7] M. Wu, C. Dick, J. R. Cavallaro, and C. Studer, "High-Throughput Data Detection for Massive MU-MIMO-OFDM Using Coordinate Descent," *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 63, no. 12, pp. 2357-2367, 2016.
- [8] L. Dai, X. Gao, X. Su, S. Han, I. C. L., and Z. Wang, "Low-Complexity Soft-Output Signal Detection Based on Gauss&Side Method for Uplink Multiuser Large-Scale MIMO Systems," *IEEE Transactions on Vehicular Technology*, vol. 64, no. 10, pp. 4839-4845, 2015.
- [9] G. Peng, L. Liu, P. Zhang, S. Yin, and S. Wei, "Low-Computing-Load, High-Parallelism Detection Method based on Chebyshev Iteration for Massive MIMO Systems with VLSI Architecture," *IEEE Transactions on Signal Processing*, vol. PP, no. 99, pp. 1-1, 2017.
- [10] A. Elgabli, A. Elghariani, V. Aggarwal, and M. R. Bell, "A low-complexity detection algorithm for uplink massive MIMO systems based on alternating minimization," *IEEE Wireless Communications Letters*, vol. 8, no. 3, pp. 917-920, 2019.
- [11] M. Amiri and M. F. Naeiny, "Low-Complexity Iterative Detection for Uplink Multiuser Large-Scale MIMO," *Journal of Information Systems and Telecommunication (JIST)*, vol. 1, no. 29, p. 25, 2020.
- [12] Z. Gao, L. Dai, S. Han, I. Chih-Lin, Z. Wang, and L. Hanzo, "Compressive sensing techniques for next-generation wireless communications," *IEEE Wireless Communications*, vol. 25, no. 3, pp. 144-153, 2018.
- [13] M. Amiri and A. Akhavan, "An iterative detector based on sparse bayesian error recovery for uplink large-scale MIMO systems," *AEU-International Journal of Electronics and Communications*, vol. 138, p. 153848, 2021.
- [14] R. Ran, J. Wang, S. K. Oh, and S. N. Hong, "Sparse-aware minimum mean square error detector for MIMO systems," *IEEE Communications Letters*, vol. 21, no. 10, pp. 2214-2217, 2017.
- [15] S. Kwon, J. Wang, and B. Shim, "Multipath matching pursuit," *IEEE Transactions on Information Theory*, vol. 60, no. 5, pp. 2986-3001, 2014.
- [16] D. L. Donoho, Y. Tsaig, I. Drori, and J.-L. Starck, "Sparse solution of underdetermined systems of linear equations by stagewise orthogonal matching pursuit," *IEEE transactions on Information Theory*, vol. 58, no. 2, pp. 1094-1121, 2012.

- [17] Z. Zhang, Y. Xu, J. Yang, X. Li, and D. Zhang, "A survey of sparse representation: algorithms and applications," *IEEE access*, vol. 3, pp. 490-530, 2015.
- [18] H. Xu, C. Caramanis, and S. Mannor, "Robust regression and lasso," *IEEE Transactions on Information Theory*, vol. 56, no. 7, pp. 3561-3574, 2010.
- [19] M. Tan, I. W. Tsang, and L. Wang, "Matching pursuit LASSO part I: Sparse recovery over big dictionary," *IEEE Transactions on Signal Processing*, vol. 63, no. 3, pp. 727-741, 2014.
- [20] R. Torkamani and R. A. Sadeghzadeh, "Wavelet-based Bayesian Algorithm for Distributed Compressed Sensing," *Information Systems & Telecommunication*, p. 87, 2019.
- [21] W.-C. Chang and Y. T. Su, "Sparse Bayesian Learning Based Tensor Dictionary Learning and Signal Recovery with Application to MIMO Channel Estimation," *IEEE Journal of Selected Topics in Signal Processing*, vol. 15, no. 3, pp. 847-859, 2021.
- [22] I. Daubechies, M. Defrise, and C. De Mol, "An iterative thresholding algorithm for linear inverse problems with a sparsity constraint," *Communications on Pure and Applied Mathematics: A Journal Issued by the Courant Institute of Mathematical Sciences*, vol. 57, no. 11, pp. 1413-1457, 2004.
- [23] T. Blumensath and M. E. Davies, "Iterative thresholding for sparse approximations," *Journal of Fourier analysis and Applications*, vol. 14, no. 5-6, pp. 629-654, 2008.
- [24] C. J. Miosso, R. von Borries, M. Arguez, L. Velázquez, C. Quintero, and C. Potes, "Compressive sensing reconstruction with prior information by iteratively reweighted least-squares," *IEEE Transactions on Signal Processing*, vol. 57, no. 6, pp. 2424-2431, 2009.
- [25] J. Yang and Y. Zhang, "Alternating direction algorithms for l_1 -problems in compressive sensing," *SIAM journal on scientific computing*, vol. 33, no. 1, pp. 250-278, 2011.
- [26] M. Elad, *Sparse and redundant representations: from theory to applications in signal and image processing*. Springer Science & Business Media, 2010.
- [27] S. Boyd, N. Parikh, and E. Chu, *Distributed optimization and statistical learning via the alternating direction method of multipliers*. Now Publishers Inc, 2011.
- [28] A. Beck and M. Teboulle, "A fast iterative shrinkage-thresholding algorithm for linear inverse problems," *SIAM journal on imaging sciences*, vol. 2, no. 1, pp. 183-202, 2009.

A Hybrid Approach based on PSO and Boosting Technique for Data Modeling in Sensor Networks

Hadi Shakibian^{1*}, Jalal A. Nasiri²

¹. Department of Computer Engineering, Faculty of Engineering, Alzahra University, Tehran, Iran

². Department of Mathematical Sciences, Ferdowsi University of Mashhad, Mashhad, Iran

Received: 14 Jul 2021/ Revised: 04 Dec 2021/ Accepted: 29 Jan 2022

Abstract

An efficient data aggregation approach in wireless sensor networks (WSNs) is to abstract the network data into a model. In this regard, regression modeling has been addressed in many studies recently. If the limited characteristics of the sensor nodes are omitted from consideration, a common regression technique could be employed after transmitting all the network data from the sensor nodes to the fusion center. However, it is not practical nor efficient. To overcome this issue, several distributed methods have been proposed in WSNs where the regression problem has been formulated as an optimization based data modeling problem. Although they are more energy efficient than the centralized method, the latency and prediction accuracy needs to be improved even further. In this paper, a new approach is proposed based on the particle swarm optimization (PSO) algorithm. Assuming a clustered network, firstly, the PSO algorithm is employed asynchronously to learn the network model of each cluster. In this step, every cluster model is learnt based on the size and data pattern of the cluster. Afterwards, the boosting technique is applied to achieve a better accuracy. The experimental results show that the proposed asynchronous distributed PSO brings up to 48% reduction in energy consumption. Moreover, the boosted model improves the prediction accuracy about 9% on the average.

Keywords: Wireless sensor network; Distributed optimization; Particle swarm optimization; Regression; Boosting.

1- Introduction

In wireless sensor networks (WSNs), keeping massive ongoing data is an expensive task due to the limited power supply and capacity of the sensor nodes. Moreover, this data is expected to be analyzed in order to extract more useful information about the phenomenon of interest. In this regard, regression modelling has been addressed as an efficient approach for abstracting [1], [2] and analyzing the network data [3], [4].

Distributed data and limited characteristics of the sensor nodes impose major challenges on performing regression over WSNs. A naive simple solution is to gather all the network data in the fusion center and obtain the network regressor using a well-known technique [5], [6]. Although a high accuracy is achieved, a huge data transmission from the sensor nodes to the fusion center is needed which makes this solution inapplicable, especially when the network grows in size.

To overcome both the communication and the computation constraints of the sensor nodes, several Learning/optimization algorithms have been proposed in many research papers.

A distributed sub-gradient algorithm with uncoordinated dynamic step sizes has been proposed for multi-agent convex optimization problems [7]. In this algorithm, each agent i can utilize its estimation of the local function value. Theoretical analysis show all the agents reach a consensus on the optimal solution. The gradient methods have also been studied by [8], [9], [10] over a network with communication constraints.

In [11], the information theoretic optimality of the distributed learning algorithms has been addressed in which each node is given i.i.d. samples and sends an abstracted function of the observed samples to a central node for decision making.

The use of machine learning algorithms in clustered WSNs has been studied by [12] in order to decrease data communications and make use of the features of WSNs. Different applications of machine learning algorithms in

✉ Hadi Shakibian
h.shakibian@alzahra.ac.ir

the context of WSNs has been recently reviewed by [13], [14], [15].

A kernel regression algorithm has been introduced in [16] to predict a signal y_t defined over the N network nodes with a series of T regularly sampled data points. A Laplace approximation is proposed to provide a lower bound for the marginal out-of-sample prediction uncertainty to address the large problems.

Logistic regression fusion rule (LRFR) has been proposed in [17] in which the coefficients of the LRFR is learnt at first, and then, it is used to make a global decision about the presence/absence of the target.

In [18], a quantized communication based distributed online regression algorithm has been proposed. Also, a distributed quantile regression algorithm has been proposed by [19], where, each node estimates the global parameter vector of a linear regression model by employing its local data as well as collaboration with the other nodes. Due to the sparsity of numerous natural and artificial systems, they have introduced l_1 – distributed quantile regression algorithm to exploit the sparsity and consequently to improve the performance of the method.

An energy-efficient distributed learning framework has been proposed using the quantized signals in the context of IoT networks [20], [21]. This is a recursive least-squares algorithm that learns the parameters using low-bit quantized signals and requires low computational cost.

Some distributed learning algorithms have also been suggested based on linear and polynomial regression models [22], [23].

On the other hand, several distributed regression models have been proposed in WSNs in which the learning problem is formulated as an optimization task [24]. To solve it, Incremental Gradient (IG) algorithm has been proposed in which the parameter to be estimated is circulated through the network. Along the way, each sensor node adjusts the parameter by performing a sub-gradient [25] based on its own local data set. Increasing the network cycles, the accuracy might be improved. In [26], IG has been proposed with the addition of quantization technique which can be used in the presence of low bandwidth to reduce the bits of transmitted data. In [27], a cluster-based version of IG has been developed. It brings a better energy efficiency and robustness. Incremental Nelder-Mead Simplex (IS) has been proposed in [28] and [29] with the addition of boosting and re-sampling techniques, respectively. They introduce a better accuracy and convergence rate.

In [30] a new evolutionary based approach has been proposed based on the PSO algorithm, denoted as Distributed PSO (DP). In DP, the network is partitioned into a number of clusters, dedicating a swarm of particles for which. Then the regressor of each cluster is trained by employing PSO algorithm distributively within the cluster. The final model is obtained after combining the clusters

models by the fusion center. This approach obtains a model closer to the centralized case, and decreases the latency significantly. However, its synchronous processes are in contrast with autonomous nature of WSNs. In addition, different clusters have their own cluster size and data pattern which are not taken into account by DP.

IVeP [31] is another PSO based distributed approach that learns the network regression model using a multi-objective optimization technique. They employ VEPSO model to perform the optimization task through inter- and intra-cluster cycles. The results show high prediction accuracy with moderate energy consumption.

In this paper, a modified version of DP algorithm is proposed that can simultaneously decrease the communication overheads as well as improves the final prediction accuracy. Firstly, Asynchronous DP (ADP), has been proposed by defining a diversity threshold for the particles within each cluster swarm. As a result, each cluster regressor is learned regardless of the status of the other clusters. Defining diversity thresholds, the number of transmissions is reduced. However the final accuracy might be decreased on the other hand. In this regard, Boosted ADP (BADP) has been introduced which boosts the clusters regressors and keeps the overall accuracy in high. The proposed algorithms have been compared with IG- and IS-based algorithms as well as IVEP and centralized approaches in terms of the accuracy, latency, and communication cost. The results show that ADP and BADP bring the lowest latency. Moreover, thanks to the boosting technique, BADP learns a model closer to the centralized approach while the communication cost still remains considerably acceptable. The contributions of this paper are:

- Asynchronous DP algorithm is proposed in which in-cluster optimization is performed asynchronously based on the size and data patterns of the cluster. While this is in accordance to the autonomous operations of the sensor networks, it brings more energy efficiency.
- The obtained model by ADP is boosted to improve the overall accuracy even further. Accordingly, Boosted ADP algorithm is proposed that obtains a high accurate network model and closer to the centralized approach with quite acceptable communication requirements.

The rest of this paper is organized as follows. Distributed regression problem is formally stated in section 2. The proposed approach is introduced in section 3. Evaluation and experimental results are discussed in section 4 and the last section is concluding remarks.

2- Distributed Regression in WSNs

Consider a sensor network with n nodes and m measurements per node spatially distributed in an

area. Every sensor node is expected to capture the phenomenon of interest in pre-defined time intervals [32]. Each measurement is stored as a record as:

$$\langle locx_i, locy_i, time_{i,t}, l_{i,t} \rangle$$

in which $(locx_i, locy_i)$ denotes to the i -th node's location, $time_{i,t}$ is epoch number, and $l_{i,t}$ is the captured measurement. Now, considering

$$A = \{locx_i, locy_i, time_{i,t}\}_{i=1, \dots, n}^{t=1, \dots, m}$$

as the feature space and:

$$B = \{l_{i,t}\}_{i=1, \dots, n}^{t=1, \dots, m}$$

as the labels, the aim of the parametric regression is to learn the coefficients of the mapping function $g: A \rightarrow B$, i.e. θ , such that the RMS error be minimized:

$$RMSE(g(A|\theta)) = \sqrt{\frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m [g(locx_i, locy_i, time_{i,t}|\theta) - l_{i,t}]^2} \quad (1)$$

Throughout this paper the following assumptions will be held:

- The learning process starts by disseminating a query from the fusion center to cluster heads.
- Every sensor node can localize itself by executing a well-known localization algorithm [33], [34].
- Since clustering is not the subject of this paper, it is assumed that the network is partitioned into C clusters via a well know clustering algorithm [35], [36], designating a cluster head for each cluster, CH_1, \dots, CH_C .
- The member nodes belonging to the cluster j are denotes as $\{sn_1^{(j)}, \dots, sn_{n_j}^{(j)}\}$ where n_j is the size of the cluster.
- The local data set of $sn_i^{(j)}$, cluster data j , and global network data are denoted as $LD_i^{(j)}$, $CD_j = \cup_{i=1}^{n_j} \{LD_i^{(j)}\}$, and $GD = \cup_{j=1}^C \{CD_j\}$, respectively.
- The L denotes the size of the parameter under estimate.

Table 1 shows the Nomenclature used in this study.

Table 1. Nomenclature used in this study.

Symbol	Definition
n	Number of sensor nodes
m	Number of sensor measurements
C	Number of clusters
CH_j	Cluster head j
$sn_i^{(j)}$	Sensor node i in cluster j
$LD_i^{(j)}$	The local data of $sn_i^{(j)}$
CD_j	The cluster data j
GD	The global(network) data
g_j	The cluster regression model j
g_{net}	The network regression model
A	The feature space
N_s	The swarm size
N_d	The problem dimensionality
$S^{(j)}$	The diversity of swarm j
$p_{i,d}$	The dimension d of particle i
τ_i^j	The weight of the local repressor of $sn_i^{(j)}$

3- The Proposed Approach

In Algorithm1, the basic idea of DP algorithm [13] has been recalled. The ADP is introduced afterwards. In summary, DP has the following steps:

1. Inside each cluster, every sensor node is given a swarm of particles to learn the cluster model. To do this, each cluster node obtains the model of its local data and sends it to other cluster nodes. Then, every cluster node employs the received local models to regenerate the whole cluster data. Now, every sensor is ready to start its local PSO to learn a candidate cluster regressor.
2. During the in-cluster optimizations, in order to guarantee the convergence of different swarms of the cluster nodes, the best particles are exchanged inside the cluster.
3. After completing the in-cluster optimization process, the cluster models are transmitted from the cluster heads to the fusion center.
4. The final network model is obtained by a weighted combination technique.

Algorithm 1: Distributed PSO (DP) [13]

```

Fusion Center disseminates the desired model
for each cluster  $j$  do
  data_view_unification()
  parameters_initialization()
  for  $i$  in range  $1:N_{migrations}$ 
    for each cluster node  $i$  do
       $sn_i^j$  runs a local PSO
       $sn_i^j$  sends its best particle to the  $CH_j$ 
    end for
     $CH_j$  sends the best of the best particles to its members
  end for
   $CH_j$  sends  $g_j$  and its RMSE to the fusion center
end for
The fusion center obtains  $G_{net}$  by weighted averaging

```

3-1- Asynchronous DP (ADP)

The major drawback of DP is that the migration steps should be synchronized for all clusters. In more words, the particles of a particular cluster might be converged before the final migration, while more migration steps might be required in another cluster. This is because different clusters have different data patterns and cluster size. By eliminating the extra migrations inside the converged clusters, the energy consumption is reduced. Furthermore, the synchronized clusters are in contrast with the autonomous nature of WSNs. To resolve these issues, asynchronous DP (ADP) is introduced in this section.

Attractive and Repulsive PSO, called as ARPSO, is a variant of PSO model in which the particles can switch between two phases [37], [38]. This approach is based on the diversity guided evolutionary algorithm (DGEA) developed by [39]. In ARPSO, the particles obey from the diversity of the swarm to alternate between an attraction and repulsion phases to make a proper exploitation-exploration tradeoff. Accordingly, the swarm diversity is defined as:

$$diversity(S(z)) = \frac{1}{N_s} \sum_{i=1}^{N_s} \sqrt{\sum_{d=1}^{N_d} (p_{i,d}(z) - \bar{p}_d(z))^2} \quad (2)$$

where N_s is the swarm size, N_d is the dimensionality of the problem, and \bar{p}_d is the average of the dimension d over all the particles, i.e.

$$\bar{p}_d(z) = \frac{\sum_{i=1}^{N_s} p_{i,d}(z)}{N_s} \quad (3)$$

Although ARPSO was originally applied to one swarm, nothing prevents its application to sub-swarms [40].

The diversity equation in ARPSO has been adopted in ADP for measuring the diversity of the clusters swarms. In cluster j , the diversity is calculated using only the best particles received from the cluster nodes:

$$diversity(S^j(z)) = \frac{1}{n_j} \sum_{i=1}^{N_j} \sqrt{\sum_{d=1}^{N_d} (gbest_i^j(z) - \overline{gbest}_d(z))^2} \quad (4)$$

where:

$$\overline{gbest}_d(z) = \frac{\sum_{i=1}^{n_j} gbest_{i,d}^j(z)}{n_j} \quad (5)$$

If the diversity (Eq. 4) be greater than a threshold φ , the in-cluster optimization is stopped, and the cluster regressor is transmitted to the fusion center as well as the corresponding RMS error. The final model is obtained by the fusion center similar to the idea proposed in DP algorithm. The steps of ADP is shown in Algorithm 2.

3-2- Boosted ADP (BADP)

Defining smaller thresholds, the quality of the clusters models are expected to be increased in ADP algorithm. However, it brings more communication cost. In this regard, in order to keep both energy efficiency and high accuracy, a boosting technique is applied on ADP inspiring from [28]. In Boosted ADP (BADP) algorithm, firstly, the clusters regressors are obtained using a diversity threshold, as explained in ADP. Then, each cluster model is boosted before transmitting to the fusion center. To do this, within the cluster j , the cluster head broadcasts the final obtained regressor and the size of the

Algorithm 2: Asynchronous PSO (ADP)

Fusion Center disseminates the desired model
for each cluster j **do**
 data_view_unification()
 parameters_initialization()
 for each cluster node i **do**
 sn_i^j runs a local PSO
 sn_i^j sends its best particle to the CH_j
 end for
 CH_j calculates the cluster diversity, i.e. Eq. 4
 if the diversity is larger than φ **then**
 CH_j sends the best of the best particles to its members
 else
 CH_j sends g_j and its RMSE to the fusion center
 end if
end for
The fusion center obtains G_{net} by weighted averaging

cluster data to its member nodes. Each member node, e.g. s_i^j , tests the cluster model on its own local data set and calculates a partial weight for it, ω_i^j . Afterwards, a new learner, v_i^j , is trained over data points labeled incorrectly by the cluster model [28]. Similarly, the new learner is test over the local data set and a local weight is computed [11]:

$$\tau_i^j = \frac{\# \text{ of truly labeled data points}}{|CD_j|} \quad (6)$$

The new learner acts as a weak learner when applying on the cluster data, as it has been trained over a small data set. So, it should be combined with the new learners obtained by the other cluster nodes to build a second stronger regressor. For this pupose, s_i^j sends $\tau_i^j \times v_i^j$ as well as ω_i^j to the cluster head. The cluster head aggregates the received partial weights to compute the weight of its regressor, ω_j . Now, a new boosted cluster model is obtained as:

$$g_j^{boosted} = w_j \times g_j + \sum_{i=1}^{n_j} \tau_i^j \times v_i^j \quad (7)$$

The last step is to calculate the in-cluster RMS error of the new boosted model using the cluster data as explained in ADP algorithm. Finally, the boosted model and its RMS error are sent to the fusion center, and the global model is obtained. Algorithm 3 describes the steps of BADP algorithm. Although the computational complexity has not been found as a major concern in WSNs, we can provide an estimation of the computational complexity for a single

Algorithm 3: Boosted ADP (BADP)

Fusion Center disseminates the desired model
for each cluster j **do**
 data_view_unification()
 parameters_initialization()
 for each cluster node i **do**
 sn_i^j runs a local PSO
 sn_i^j sends its best particle to the CH_j
 end for
 CH_j calculates the cluster diversity, i.e. Eq. 4
 if the diversity is larger than φ **then**
 CH_j sends the best of the best particles to its members
 else
 CH_j sends g_j and the cluster data size to its members
 for each cluster node i **do**
 sn_i^j tests g_j^{ADP} on LD_i^j and obtains two data partitions as $LD_{true,i}^j$ and $LD_{false,i}^j$
 sn_i^j computes ω_i^j , the partial weight of g_j^{ADP}
 sn_i^j runs a local PSO over $LD_{false,i}^j$ to learn v_i^j
 sn_i^j computes τ_i^j , the weight of v_i^j as Eq. 6.
 sn_i^j sends $\tau_i^j \times v_i^j$ and ω_i^j to CH_j
 end for
 CH_j computes ω_j using $\{\omega_i^j\}_{i=1}^{n_j}$
 CH_j computes its final regressor $g_j^{boosted}$ as Eq. 7
 CH_j sends g_j and its RMSE to the fusion center
 end if
end for
The fusion center obtains G_{net} by weighted averaging

sensor node belonging to the cluster j in DP, ADP, and BADP algorithms. For DP algorithm we have:

$$T(DP) = T(\text{data_view_unification step}) + T(\text{optimization})$$

where:

$$T(\text{data_view_unification}) = t(\text{local PSO}) + O(mk)$$

in which $t(\text{local PSO})$ denotes the computational time of running the PSO algorithm over the cluster data by the sensor node and $O(mk)$ is the required time for resampling of m measurements using a k -parameters data model and:

$$T(\text{optimization}) = O(M \times t(\text{local PSO}))$$

where M denotes the number of migration steps. Similarly, for ADP and BADP algorithms we have:

$$T(\text{ADP}) = O(\text{local PSO}) + O(mk) + O(M_j \times t(\text{local PSO}))$$

where M_j is the required migration steps for the corresponding cluster j and:

$$T(\text{BADP}) = t(\text{local PSO}) + O(mk) + O(M_j \times t(\text{local PSO})) + t(\text{local PSO for the boosting task})$$

Totally, the computational time complexity of a sensor node in DP algorithm could be simplified as:

$$T(\text{DP}) = O(M \times t(\text{local PSO})) + O(mk)$$

and for ADP and BADP we would have:

$$T(\text{ADP}) = T(\text{BADP}) = O(M_j \times t(\text{local PSO})) + O(mk)$$

4- Evaluation and Results

The proposed algorithms have been compared with their distributed counterparts, IG, IS, BIS, IS-Resampling, IVEP and the centralized approach. In all of these algorithms, the learning problem is formulated as an optimization task, as discussed in Section 2. Two datasets have been used for comparison. In the first one, Berkeley Intel Lab network [41], there are 54 sensor nodes with two corrupted ones. Mica2Dot sensors with weather boards capture humidity, light, voltage, and temperature in every 31 seconds. As mentioned before, regression modeling has been performed only over temperature readings. Two portions of the network data, named as DS1 and DS2, have been chosen such that each sensor node has 100 and 2880 (measuring for one day) data points, respectively. So, the global data, GD , has 5200 and 149760 data points in total, respectively. The second network, denoted as DS3, is an artificial network with 100 sensors distributed uniformly over a square of 100 m^2 . Each sensor has collected a dataset of size 200. The phenomenon under study, temperature, is sensed in each epoch from 1 to 200. This data is generated using Eq. (8) with an additive Gaussian noise of mean 0 and variance 1, $N(0,1)$. The coefficients of the model are randomly chosen in the range of $(-10, 10)$.

Ten Fold Cross Validation method (10-CV) has been adopted for each approach. In 10-CV, the data set is divided into 10 partitions for 10 times. Each time, one of the partitions is used as the test data and the learning will be executed with the remaining parts. Finally, the average results will be found. In [42], some polynomial models have been suggested for the Berkeley network data set. It was reported that a linear space and quadratic time model can be fitted more accurately. Accordingly, in order to fit a model on the network data, a spatiotemporal model

with linear in space and quadratic in time has been chosen as:

$$G_{net}(\text{locx}, \text{locy}, \text{time} | \theta) = \theta_1 \text{locx} + \theta_2 \text{locy} + \theta_3 \text{time}^2 + \theta_4 \text{time} + \theta_5 \quad (8)$$

where the node location (locx, locy) and epoch number (time) introduce the feature set while the captured temperature per epoch is the label. Therefore, the proposed approach aims to learn the coefficients of a set of basis functions as $\{\text{locx}, \text{locy}, \text{time}^2, \text{time}, 1\}$. Accordingly, the RMS error could be calculated as:

$$RMSE(G_{net}(A | \theta)) = \sqrt{\frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m [\theta_1 \text{locx}_{ij} + \theta_2 \text{locy}_{ij} + \theta_3 \text{time}_{ij}^2 + \theta_4 \text{time}_{ij} + \theta_5 - T_{ij}]^2} \quad (9)$$

In order to have a good exploration-exploitation tradeoff, the inertia weight, w , is usually decreased during the time as [43]:

$$w(z) = w_{min} + \frac{m_i - z}{m_i} (w_{max} - w_{min}) \quad (10)$$

where m_i is the maximum number of iterations and z denotes to the current iteration number. The particles starts with a maximum value w_{max} and linearly decrease their inertia weights to a pre-defined minimum value w_{min} .

As the problem addressed in this study is a data-centric application, the discrete-event simulators are not required. Accordingly, all the algorithms have been implemented with Java using Eclipse IDE and the experiments were performed on an Intel dual core processor with 4 GB RAM memory.

4-1- Prediction Accuracy

The prediction accuracy of different approaches have been shown in Table 2. As all the data points are available for the centralized approach, a good accuracy is achieved at the end of the learning process. In practice, IG suffers from a low convergence rate and requires to pass several cycles in order to obtain an average accuracy. In our experiments, the accuracy of IG has been obtained through 40 cycles. While IS based approaches obtain better results within one network cycle.

In ADP, thanks to (i) learning several candidate models and (ii) the high accuracy of each cluster model, a good accuracy is achieved. However, integrating the boosting technique with the ADP algorithm leads the final accuracy becomes much closer to the centralized case and consequently BADP outperforms its distributed counterparts in most cases. Moreover, BADP

shows more stable accuracies in both networks rather than the other methods. This indicates how the boosted

Table 2. The final RMS error of different approaches based on each data set.

Approach	DS1	DS2	DS3
IG	17.481	21.549	110.678
IS	8.206	5.059	13.452
BIS	6.268	5.011	9.647
IS-Res.	5.806	3.104	11.714
IVeP	2.219	3.009	3.348
ADP	2.060	4.311	3.448
BADP	1.892	2.917	3.017
Central	0.835	2.536	1.005

Table 3. The RMSE comparison of BADP algorithm with [44].

Algorithm	RMSE
BADP	2.917
[44]-LG	$\cong 2.35$
[44]-PV	$\cong 2.75$
[44]-UV	$\cong 3.02$
[44]-RA	$\cong 2.91$

regressor can accurately predict those parts of the phenomena that labeled inaccurately by the first learner. The prediction accuracy of the BADP has also been compared with the reported results of [44] where a multi-objective sensor placement algorithm has been proposed. The performance of the state estimation of the temperature measurements has been evaluated based on the RMSE. As shown in Table 3, the proposed BADP could obtain better or completely close prediction accuracy compared to [44] which is not a distributed algorithm.

In Figure 1, the convergence rate of different cluster swarms have been depicted based on three data sets. As it is expected, some cluster swarms converge faster than the other ones due to their size and data patterns. As a result, lesser in-cluster communications would be required. To show how the swarms diversity could decrease the energy consumption, 5 diversity thresholds have been defined for each dataset and shown in Table 4.

The prediction accuracy using each defined diversity threshold has been demonstrated in Figure 2 along with the corresponding amount of energy saving in comparison with the DP algorithm. It is concluded from the Figure 2 that the energy consumption of DP algorithm could be decreased up to %48 by the proposed approach while the final RMSE is quite high. As it was shown, by using more tight diversity thresholds, the final prediction accuracy is increased. However, a trade-off should be made between the energy consumption and the prediction accuracy using diversity thresholds. It should be noticed that the reported results in Table 1 are based on φ_4 .

Table 4. Swarms diversity thresholds

Dataset	φ_1	φ_2	φ_3	φ_4	φ_5
DS1	3.50	2.25	1.25	0.75	0.50
DS2	4.75	2.50	2.25	1.50	1.00
DS3	3.75	2.50	1.75	1.25	0.75

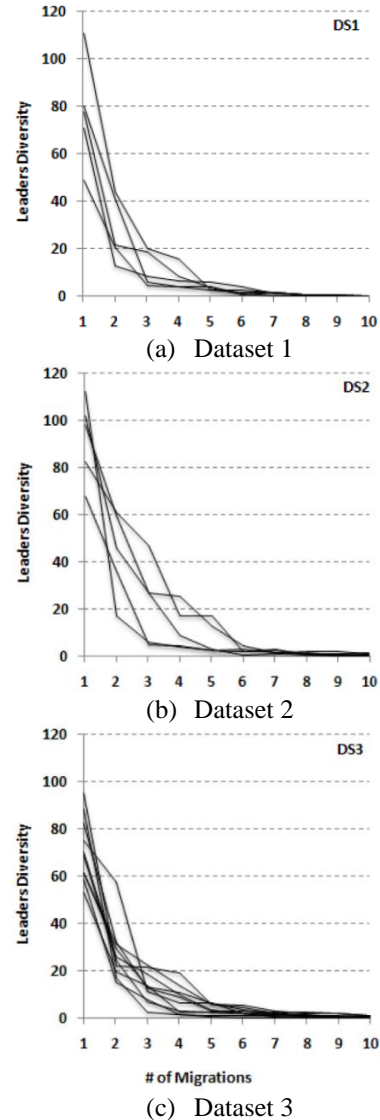


Figure 1. The convergence rate of different clusters' swarm

4-2- Latency

Regarding to the ongoing sensor measurements, the network model is valuable for some pre-specified periods of time. Thus, when the measurements is refreshed, it is required to train the regression model with the new network data. In this regard, the required time to rebuild the model is important, known as the latency metric: *the*

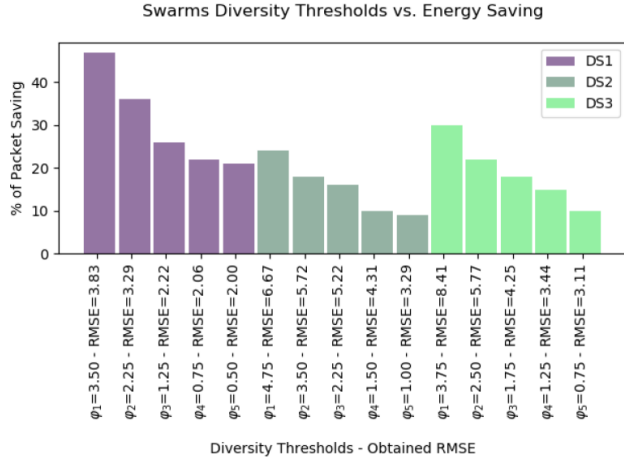


Figure 2. The impact of varying swarms diversity thresholds on the energy consumption.

number of iterations to visit all the network data for the first time [27]. The latency of different approaches has been compared in Table 5.

The centralized approach just needs one iteration to visit the whole network data, and so, its latency is $O(1)$. In IG and IS based approaches, the latency is $O(n)$, as they need one network cycle to get access to the network data. In IVEP algorithm, one in-cluster cycle is required and thus its latency is $O(n^*)$ where n^* is the size of the largest cluster. If we let the number of clusters equals to \sqrt{n} , then the latency of IVEP would be $O(\sqrt{n})$. In ADP and BADP, the training process is started synchronously in all clusters. Consequently the latency would be the same as the centralized method, i.e. $O(1)$.

4-3- Communication Cost

Bit/meter metric has been used for comparing the communication requirements of the approaches. According to [24], assume the network has been deployed in a unit square area. Having the size and average distance of each transmission, the communication requirements of each approach can be evaluated. Every transmission in all the studied approaches falls in one of the following transmission types and the corresponding average distance can be achieved similar to [24]:

cluster node - cluster node:

Algorithm	Latency
IG	$O(n)$
IS	$O(n)$
BIS	$O(n)$
IS-Resampling	$O(n)$
IVeP	$O(\sqrt{n})$
ADP	$O(1)$
BADP	$O(1)$
Centralized	$O(1)$

$$d_1 = O(\sqrt{\log^2 n/n})$$

cluster head - cluster head:

$$d_2 = O(\sqrt{\log^2 \sqrt{n}/\sqrt{n}})$$

sensor node - fusion center:

$$d_3 = O(1)$$

cluster head - cluster node:

$$d_4 = O(1/\sqrt{n})$$

cluster head - fusion center:

$$d_5 = O(1/\sqrt{n})$$

The size of the parameter(s) transmitted between two consecutive sensors in IG and IS based approaches are as follows:

- In IG algorithm, a double vector of size L is transmitted [24].
- In IS algorithm and the first pass of BIS, $|LD_k|$ (an integer of size 1) and a double vector of size L are transmitted [28].
- In the second pass of BIS, three parameter transmissions are happened: the partial weight of the learned regressor (a double value), the size of the global data, $|GD|$ (an integer of size 1), and the partial weighted combination of local regressors which is a double vector of size L .
- In IS-Resampling algorithm, a double vector of size L and a double vector of size 2 ($locx_k, locy_k$) are transmitted.

In the centralized approach, m data values are transmitted between a sensor node and the fusion center. Thus, n transmissions of size v would be required. As each data point in our experiments contains three features with a label, we have $v = 4$.

The communication requirements of the ADP algorithm is similar to DP. Firstly, a parameter of size L is transmitted from the cluster nodes to cluster

Table 6. Types and size of transmissions of different approaches

	Sensor-Sensor	Sensor-Fusion
IG	$C_{IG} \times (n-1)L$	L
IS	$C_{IS} \times (n-1)(L+1)$	L
BIS	$cost(IS) + 2(n-1)(L+1)$	$2L$
IS-R	$C_{IS-R} \times (n-1)(L+2)$	L
Centralized		nmv
	Sensor-CHead	CHead-Fusion
ADP	$3nL + n - LC + (2L+1) \sum_{j=1}^C l_j n_j$	$C(L+1)$
BADP	$4nL + 2n - 2LC - C + (2L+1) \sum_{j=1}^C l_j n_j$	$C(L+1)$

Table 7. Comparing the communication order. Without loss of generality, we follow [27] and let $C = \sqrt{n}$ and K denotes the average number of iterations in ADP and BADP algorithms.

Approach	Communication requirement	Rank
IG	$\mathcal{O}((C_{IG})(L)(\sqrt{n} \cdot \log n + 1))$	7
IS	$\mathcal{O}((C_{IS})(L)(\sqrt{n} \cdot \log n + 1) + L)$	1
BIS	$\mathcal{O}(IS) + \mathcal{O}(L)(\sqrt{n} \cdot \log n + 1)$	3
IS-Resampling	$\mathcal{O}((C_{IS-Res})(L+2) \cdot \sqrt{n} \cdot \log n) + L$	2
IVeP	$\mathcal{O}(n(L\sqrt{n}) + 3\sqrt{n}L + \sqrt{\sqrt{n}L} \log n)$	6
ADP	$\mathcal{O}((K + 2LK + 3L + 1)\sqrt{n} + 1)$	4
BADP	$\mathcal{O}((K + 2LK + 4L + 2)\sqrt{n} - L)$	5
Central	$\mathcal{O}(n)(m)(v)(1)$	8

head, and vice versa. Then, a driver message by size of $2L + 1$ is transmitted from the cluster head to the members. Afterwards, during the in-cluster optimization, the best particle of each cluster node by size of $2L$ with an RMS error of size 1, and the best of the best particles by size of $2L$ are transmitted between the cluster head and the cluster members at each migration step. Then, each cluster node sends its regressor with the corresponding RMS error to the cluster head by size of $L+1$. Finally, each cluster head sends the obtained cluster model with its RMS error by size of $L + 1$ to the fusion center. In BADP algorithm, every cluster node transmits an extra parameter of size $L + 1$, new learner plus its partial weight, to the cluster head. In Table 6, the parameter transmissions of all approaches have been summarized based on transmission size and type. For IVEP algorithm, the communication cost analysis is recalled from [31]. Accordingly, the total communication cost of different approaches are obtained as shown in Table 7.

In this regard, the communication cost of the centralized approach is the highest due to a huge data transmission. As mentioned before, IG practically needs to meet a large number of cycles to obtain an average accuracy. As a result, its energy

consumption is higher than the other distributed approaches while IS has the lowest transmissions. From Table 6, it can be understood that in-cluster communications consume less energy due to a smaller average distance. On the other hand, the main part of the transmissions in ADP as well as BADP has been spent in clusters. Thanks to this property, ADP and BADP both work moderately in terms of the energy consumption, as shown in Table 7.

5- Conclusion

A novel distributed data modeling approach has been proposed based on multi-swarm PSO algorithm. In the proposed approach, the task of learning the regression model of a cluster is assigned to a swarm of particles. Each swarm executes an in-cluster optimization process to learn the cluster regressor asynchronously. The most important feature of this approach is that the optimization of each swarm is terminated according to the size and data pattern of its cluster. This property leads to save up to 48% of energy consumption by eliminating extra migration steps while the accuracy is high. In order to improve the prediction accuracy even further, a boosting technique is also employed in a distributed manner. The proposed approach has been evaluated against two real and artificial network data and compared to common distributed regression modeling techniques as well as the centralized approach. The results show the boosted model improves the prediction accuracy about 9% on the average. Due to the recent advances in the sensor nodes technologies, more complex machine learning algorithms, such as Deep Learning, could be employed to achieve higher prediction accuracies in the context of WSNs [45], [46].

References

- [1] Sharma, Himanshu, Ahteshamul Haque, and Frede Blaabjerg. "Machine Learning in Wireless Sensor Networks for Smart Cities: A Survey." *Electronics* 10.9 (2021): 1012.
- [2] Liu, Longgeng, et al. "An algorithm based on logistic regression with data fusion in wireless sensor networks." *EURASIP Journal on Wireless Communications and Networking* 2017.1 (2017): 1-9.
- [3] Deng, Yulong, et al. "Temporal and spatial nearest neighbor values based missing data imputation in wireless sensor networks." *Sensors* 21.5 (2021): 1782.
- [4] Zuhairy, Ruwaida M., and Mohammed GH Al Zamil. "Energy-efficient load balancing in wireless sensor network: An application of multinomial regression analysis." *International Journal of Distributed Sensor Networks* 14.3 (2018): 1550147718764641.

- [5] Kumar, D. Praveen, Tarachand Amgoth, and Chandra Sekhara Rao Annavarapu. "Machine learning algorithms for wireless sensor networks: A survey." *Information Fusion* 49 (2019): 1-25.
- [6] Ghate, Vasundhara V., and Vaidehi Vijayakumar. "Machine learning for data aggregation in WSN: A survey." *International Journal of Pure and Applied Mathematics* 118.24 (2018): 1-12.
- [7] Ren, Xiaoxing, et al. "Distributed Subgradient Algorithm for Multi-Agent Optimization With Dynamic Stepsize." *IEEE/CAA Journal of Automatica Sinica* 8.8 (2021): 1451-1464.
- [8] Doan, Think T., Siva Theja Maguluri, and Justin Romberg. "Convergence rates of distributed gradient methods under random quantization: A stochastic approximation approach." *IEEE Transactions on Automatic Control* (2020).
- [9] Zhang, Peng, and Gejun Bao. "An incremental subgradient method on Riemannian manifolds." *Journal of Optimization Theory and Applications* 176.3 (2018): 711-727.
- [10] Berahas, Albert S., Charikleia Iakovidou, and Ermin Wei. "Nested distributed gradient methods with adaptive quantized communication." 2019 IEEE 58th Conference on Decision and Control (CDC). IEEE, 2019.
- [11] Xu, Xiangxiang, and Shao-Lun Huang. "An information theoretic framework for distributed learning algorithms." 2021 IEEE International Symposium on Information Theory (ISIT). IEEE, 2021.
- [12] Perumal, T. Sudarson Rama, V. Muthumanikandan, and S. Mohanalakshmi. "Energy Efficiency Optimization in Clustered Wireless Sensor Networks via Machine Learning Algorithms." *Machine Learning and Deep Learning Techniques in Wireless and Mobile Networking Systems*. CRC Press, 2021. 59-77.
- [13] Kumar, D. Praveen, Tarachand Amgoth, and Chandra Sekhara Rao Annavarapu. "Machine learning algorithms for wireless sensor networks: A survey." *Information Fusion* 49 (2019): 1-25.
- [14] Mohanty, Lipika, et al. "Machine Learning-Based Wireless Sensor Networks." *Machine Learning: Theoretical Foundations and Practical Applications*. Springer, Singapore, 2021. 109-122.
- [15] Pundir, Meena, and Jasminder Kaur Sandhu. "A systematic review of Quality of Service in Wireless Sensor Networks using Machine Learning: Recent trend and future vision." *Journal of Network and Computer Applications* (2021): 103084.
- [16] Antonian, Edward and Peters, Gareth and Peters, Gareth and Chantler, Michael John and Yan, Hongxuan, GLS Kernel Regression for Network-Structured Data (August 9, 2021). Available at SSRN: <https://ssrn.com/abstract=3901694>
- [17] Liu, Longgeng, et al. "An algorithm based on logistic regression with data fusion in wireless sensor networks." *EURASIP Journal on Wireless Communications and Networking* 2017.1 (2017): 1-9.
- [18] Wang, Heyu, Lei Xia, and Chunguang Li. "Distributed online quantile regression over networks with quantized communication." *Signal Processing* 157 (2019): 141-150.54532343WE332
- [19] Wang, Heyu, and Chunguang Li. "Distributed quantile regression over sensor networks." *IEEE Transactions on Signal and Information Processing over Networks* 4.2 (2017): 338-348.
- [20] Danaee, Alireza, Rodrigo C. de Lamare, and Vitor H. Nascimento. "Energy-Efficient Distributed Recursive Least Squares Learning with Coarsely Quantized Signals." 2020 54th Asilomar Conference on Signals, Systems, and Computers. IEEE, 2020.
- [21] Danaee, Alireza, Rodrigo C. de Lamare, and Vitor H. Nascimento. "Energy-efficient distributed learning with coarsely quantized signals." *IEEE Signal Processing Letters* 28 (2021): 329-333.
- [22] Hellkvist, Martin, Ayça Özçelikkale, and Anders Ahlén. "Generalization error for linear regression under distributed learning." 2020 IEEE 21st International Workshop on Signal Processing Advances in Wireless Communications (SPAWC). IEEE, 2020.
- [23] Shuman, David I., et al. "Distributed signal processing via Chebyshev polynomial approximation." *IEEE Transactions on Signal and Information Processing over Networks* 4.4 (2018): 736-751.
- [24] M. Rabbat and R. Nowak, "Distributed Optimization in Sensor Networks," in Proceedings of the 3rd international symposium on Information processing in sensor networks, Berkeley, California, USA, (2004), pp. 20-27.
- [25] Bertsekas, Dimitri P., "Incremental gradient, subgradient, and proximal methods for convex optimization: a survey," *Optimization for Machine Learning*, No. 85, pp. 1-38, 2011.
- [26] M. Rabbat, and Nowak, R. "Quantized Incremental Algorithms for Distributed Optimization," *IEEE Journal on Selected Areas in Communications*, 23 (4) (2006), pp. 798-808.
- [27] S.H. Son, M. Chiang, S. R. Kulkarni, and S. C. Schwartz, "The Value of Clustering in Distributed Estimation for Sensor Networks," in proceedings of International Conference on Wireless Networks, Communications and Mobile Computing, Maui, Hawaii, 2 (2005), pp. 969-974.
- [28] P.J. Marandi, N.M. Charkari, "Boosted Incremental Nelder-Mead Simplex Algorithm: Distributed Regression in Wireless Sensor Networks," *IFIP Joint Conference on Mobile and Wireless Communications Networks*, France, (2008), pp. 199-212.
- [29] P.J. Marandi, M. Mansourizadeh, N. M. Charkari, "The Effect of Resampling on Incremental Nelder-Mead Simplex Algorithm: Distributed Regression over Wireless Sensor Network," in Proceedings of the Third International Conference on Wireless Algorithms, Systems, and Applications, LNCS, 5258 (2008), Dallas, Texas, pp. 420-431.
- [30] H. Shakibian and N. Moghadam Charkari, "D-PSO for Distributed Regression over Wireless Sensor Networks," *Iranian Journal of Electrical and Computer Engineering*, Vol. 11, No. 1, pp. 43-50, 2012.
- [31] Shakibian, Hadi, and Nasrollah Moghadam Charkari. "In-cluster vector evaluated particle swarm optimization for distributed regression in WSNs." *Journal of network and computer applications* 42 (2014): 80-91.
- [32] Zhao, Jijun, Hao Liu, Zhihua Li, and Wei Li., "Periodic Data Prediction Algorithm in Wireless Sensor Networks,"

- In *Advances in Wireless Sensor Networks*, pp. 695-701, 2013.
- [33] Cheng, Long, et al. "An Indoor Localization Algorithm based on Modified Joint Probabilistic Data Association for Wireless Sensor Network." *IEEE Transactions on Industrial Informatics* (2020).
- [34] Shahbazian, Reza, and Seyed Ali Ghorashi. "Distributed cooperative target detection and localization in decentralized wireless sensor networks." *The Journal of Supercomputing* 73.4 (2017): 1715-1732.
- [35] Zandhessami, Hessam, Mahmood Alborzi, and Mohammadsadegh Khayyatian. "Energy Efficient Routing-Based Clustering Protocol Using Computational Intelligence Algorithms in Sensor-Based IoT." *Journal of Information Systems and Telecommunication (JIST)* 1.33 (2021): 55.
- [36] Daanoune, Ikram, Baghdad Abdennaceur, and Abdelhakim Ballouk. "A comprehensive survey on LEACH-based clustering routing protocols in Wireless Sensor Networks." *Ad Hoc Networks* (2021): 102409.
- [37] Qi, Hong, et al. "Inversion of particle size distribution by spectral extinction technique using the attractive and repulsive particle swarm optimization algorithm." *Thermal Science* 19.6 (2015): 2151-2160.
- [38] Mo, Simin, Jianchao Zeng, and Weibin Xu. "Attractive and repulsive fully informed particle swarm optimization based on the modified fitness model." *Soft Computing* 20.3 (2016): 863-884.
- [39] Ursem, Rasmus K. "Diversity-guided evolutionary algorithms." *International Conference on Parallel Problem Solving from Nature*. Springer, Berlin, Heidelberg, 2002.
- [40] A.P. Engelbrecht, *Computational Intelligence: An introduction*, 2ed., Wiley, 2007.
- [41] Madden, S., 2003. Intel Berkeley research lab data. USA: Intel Corporation, 2004 [2004-06-08]. <http://berkeley.intel-research.net/labdata>, html.
- [42] C. Guestrin, P. Bodi, R. Thibau, M. Paskin, and S. Madde, "Distributed Regression: An Efficient Framework for Modeling Sensor Network data," in *Proceedings of third international symposium on Information processing in sensor networks*, Berkeley, California, USA, (2004), pp. 1-10.
- [43] Y. Shi, R.C. Eberhart, "Empirical study of particle swarm optimization," in *Proceedings of the IEEE International Congress on Evolutionary Computation*, 3 (1999), pp. 101-106.
- [44] Xu, Zhaoyi, Yanjie Guo, and Joseph Homer Saleh. "Multi-objective optimization for sensor placement: An integrated combinatorial approach with reduced order model and Gaussian process." *Measurement* 187 (2022): 110370.
- [45] Premkumar, M., and T. V. P. Sundararajan. "DLDM: Deep learning-based defense mechanism for denial of service attacks in wireless sensor networks." *Microprocessors and Microsystems* 79 (2020): 103278.
- [46] Mohanty, Sachi Nandan, et al. "Deep learning with LSTM based distributed data mining model for energy efficient wireless sensor networks." *Physical Communication* 40 (2020): 101097.

Detection of Attacks and Anomalies in the Internet of Things System using Neural Networks Based on Training with PSO Algorithms, Fuzzy PSO, Comparative PSO and Mutative PSO

Mohammad Nazarpour¹, Navid Nezafati^{2*}, Sajjad Shokouhyar²

¹. Department of Information Technology Management, Central Tehran Branch, Islamic Azad University, Tehran, Iran

². Department of Management, Shahid Beheshti University, Tehran, Iran

Received: 13 Jun 2021/ Revised: 04 Nov 2021/ Accepted: 13 Dec 2021

Abstract

Integration and diversity of IOT terminals and their applicable programs make them more vulnerable to many intrusive attacks. Thus, designing an intrusion detection model that ensures the security, integrity, and reliability of IOT is vital. Traditional intrusion detection technology has the disadvantages of low detection rates and weak scalability that cannot adapt to the complicated and changing environment of the Internet of Things. Hence, one of the most widely used traditional methods is the use of neural networks and also the use of evolutionary optimization algorithms to train neural networks can be an efficient and interesting method. Therefore, in this paper, we use the PSO algorithm to train the neural network and detect attacks and abnormalities of the IOT system. Although the PSO algorithm has many benefits, in some cases it may reduce population diversity, resulting in early convergence. Therefore, in order to solve this problem, we use the modified PSO algorithm with a new mutation operator, fuzzy systems and comparative equations. The proposed method was tested with CUP-KDD data set. The simulation results of the proposed model of this article show better performance and 99% detection accuracy in detecting different malicious attacks, such as DOS, R2L, U2R, and PROB.

Keywords: Attack Detection; Internet of Things (IOT); Neural Network; PSO Algorithm; Fuzzy Rule; Adaptive Formulation.

1-Introduction

With the advancement of information technology, IT-related issues have also developed rapidly. The Internet of Things is a new model that integrates the Internet and physical objects belonging to different fields such as home automation, industrial process, human health and environmental monitoring. Having Internet-connected devices deepens our day-to-day operations, in addition to having many benefits, brings with many security challenges. For more than two decades, intrusion detection systems have been an important tool for protecting networks and information systems. However, it is difficult to apply the former IDS techniques to the Internet of Things because of its special features such as limited resources, special protocol stacks, and certain standards. The proliferation of IOT has led to new challenges such as increased power consumption, more complex management due to increased data volume, more bandwidth demands to transmit IOT data, and use more powerful processors for

information analysis. Moreover, protecting the privacy of individuals by protecting and safeguarding the information of individuals is very important and vital to achieve the commercialization of this industry [1-3]. Today, of course, the use of technologies such as optical fibers in the transmission of information and optical integrated circuits with Nano dimensions in fast processing and reducing energy consumption has greatly contributed to the commercialization of the Internet of Things. In contrast, the use of cloud storage, computing for data storage, processor and the use of SDN-based software pose a serious threat to attackers of the IOT infrastructure. Threats and anomalies created in the Internet of Things can be divided into four general categories: Dos attacks, R2L attacks, U2R attacks and Probing attacks. In Dos attacks, a large number of requests are sent to a system to disable it. In the U2R attacks, the intruder enters as the system administrator and destroys the system radically. In the R2L type of attack, the attacker enters the system as a local user and then takes control of the system by designing attacks. In the Probing attack, the intruder tries

to obtain information from the system such as passwords, user numbers, important files and types of system services. One of the most important and popular tools in the field of attack prevention is the use of machine learning systems [4-5]. In this system, the system is modeled using artificial intelligence and based on existing experiences to prepare for predicting new conditions. Therefore, the system should be trained using training data that is the result of past experiences. One of the most powerful and efficient modeling tools in the field of machine learning is the use of artificial neural networks [6-7]. In simpler terms, neural networks are modern systems and computational methods for machine learning, knowledge display, and finally the application of knowledge gained to maximize the output responses of complex systems. The main idea of such networks is partly inspired by the way the biological neural system works to process data and information to learn and create knowledge. The key element of this idea is to create new structures for the information processing system. The system is made up of a large number of extremely interconnected processing elements called neurons that work together to solve a problem and transmit information through synapses (electromagnetic communications). In these networks, if one cell is damaged, other cells can make up for its absence and contribute to its regeneration. These networks are able to learn. For example, by injecting tactile nerve cells, the cells learn not to go to the hot body, and with this algorithm, the system learns to correct its error. Learning in these systems is adaptive, that is, using examples, the weight of the synapses changes in such a way that the system produces the correct response if new inputs are given. The main philosophy of the artificial neural network is to model the processing properties of the human brain to approximate conventional computational methods with the biological processing method. In other words, the artificial neural network is a method that learns the knowledge of the communication between several data sets through training and stores it to use in similar cases. This processor works in two ways similar to the human brain: Neural network learning is done through education. Weighting similar to the information storage system takes place in the neural network of the human brain.

An artificial neural network consists of three layers: input, output and processing. Each layer contains a group of nerve cells (neurons) that normally communicate with all the neurons in the other layers unless the user restricts communication between neurons; but the neurons in each layer have no connection with other neurons in the same layer. A neuron is the smallest unit of information processing that forms the basis of the function of neural networks. A neural network is a collection of neurons that, being located in different layers, form a special architecture based on the connections between neurons in different layers. Neurons can be a nonlinear mathematical

function, so a neural network made up of a community of these neurons and can also be a completely complex, nonlinear system. In the neural network, each neuron operates independently, and the overall behavior of the network is the result of the behavior of multiple neurons. In other words, neurons correct each other in a process of cooperation. Figure 1 shows an artificial neural network versus the neural network of the human body. The system inputs, called X_1, X_2, X_n , enter in the input neurons and transfer to the hidden layers via $W_1, W_2 \dots W_n$. This transfer is done by multiple inputs on the W coefficients. Now apply the nonlinear function to the output layer to enhance the modelling application for nonlinear samples and data collections. In the learning procedure, the w coefficient is determined by machine learning algorithms. Now we want to focus on the weight coefficient determination. One of the interesting and attractive methods is the backpropagation algorithm for determination the W coefficients [8-9]. This method is based on the slope of error and has a good speed to response determination. Instead, it is trapped in the local optimum point and unable to find the global optimum [10-11]. One solution is to use meta-heuristic algorithms. A metaheuristic optimization algorithm is an innovative method that can be applied to various optimization problems with minimal modifications. Metamorphism algorithms significantly increase the ability to find high-quality solutions to difficult optimization problems.

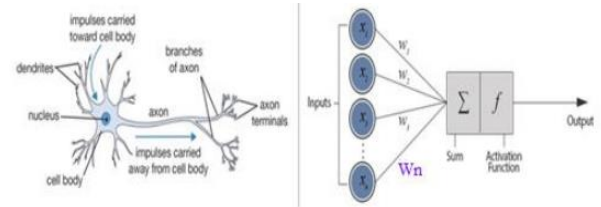


Fig. 1. schematic biological neuron (Left) versus artificial neural network (Right)

One of the evolutionary optimization algorithms that have a good performance speed is the PSO algorithm. Additionally, this algorithm, like other evolutionary algorithms (genetics and colonial competition and so on) has simpler calculations. In this article, we used the PSO algorithm to train the neural network. Then, we will show that although training by the PSO algorithm gives a much more accurate answer than the BP training method, it is still possible to reach much more accurate answers by changing the PSO algorithm. For this purpose, we used a combination of fuzzy, comparative and mutation methods to alter this algorithm and showed that we get very acceptable results by training the neural network by the altered PSO algorithm.

2-Methodology

This research work aims to propose a Neural Network Model of KDD-Data set for intrusion detection in IOT devices. This part of the paper describes the proposed work methodology, i.e., proposed attack detection framework, proposed network model, data set description, and preprocessing.

a. The Framework of the attack detection

The proposed procedure is illustrated in Figure 2. As you can see in this flowchart, the data must first be collected for training the neural network. To collect the data, we use the kdd-cup data set. In the continuation of the data preprocessing operation is done, it includes deleting similar data, extracting more effective data, and normalizing the data. We then classify the data into two categories: training data and test data, so that test data makes up 20% of the data and training data makes up 80% of the data. In the next step we architected the neural network and trained it based on PSO and Modified PSO algorithms and training data. Finally, the evaluation of the model created by the neural network is performed based on test data.

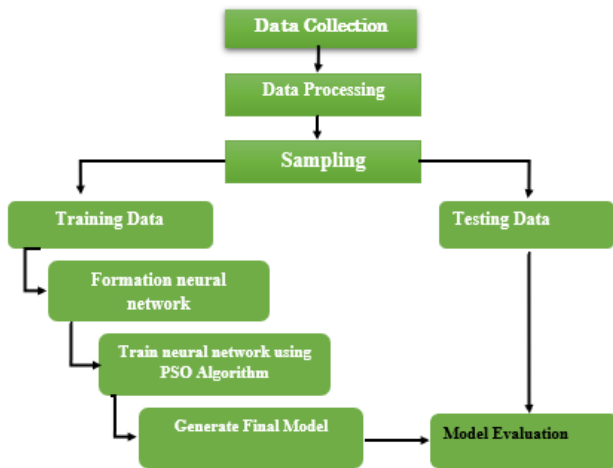


Fig. 2. overall framework of the attack detection using neural-network based PSO algorithm

b. Neural Network

An artificial neural network, also called a simulated neural network or neural network, is an interconnected group of artificial neurons that uses a mathematical or computational model to process information and based on the connection approach. One of the classic types of the artificial neural networks is the perceptron network. The following figure shows a perceptron neural network:

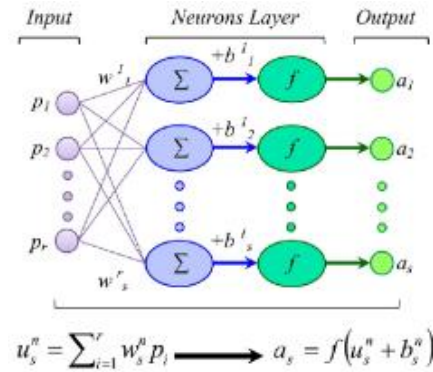


Fig. 3. Multi-layer Perceptron neural network.

A multilayered (deep) perceptron neural network will result from the stacking of several perceptron's. That is, we will have multiple layers of neurons in such a network. Here we have an output layer and an input layer. There are also several layers of neurons between the input and output layers. The layers between the input and output layers are called the Hidden Layer. Layers that are close to the input layer are usually called bottom layers. Layers that are close to the output layer are also called top layers. Except for the output, each layer has a bias. A network that has a large number of hidden layers is called a Deep Neural Network. As mentioned, there are several classical methods for determining weight and bias coefficients. But all of these methods are caught in local optimal points and are not able to determine the global optimal point. To solve this problem, in this paper, we use the training method based on the PSO optimization algorithm and extract these coefficients for system modeling. In addition, despite the high speed of the bird algorithm, it does not have enough accuracy and to improve the system, we use fuzzy, mutation and adaptive models to increase the accuracy of neural network performance in addition to speed.

c. Classical PSO

Particle swarm optimization algorithms are one of the heuristic optimization algorithms. The most significant benefit of these algorithms over the other optimization algorithms is that they do not postulate intricate operations and mathematical relationships such as integrals and derivatives [12-13]. These algorithms are either modeled on the foundation of the biological processes and exchanges of organisms (such as ants, particles, genetics, etc.) or human socio-political exchanges and treatments (such as colonial competition algorithms or teacher-learning based optimization) [14-15]. The PSO algorithm is also modeled based on the search for appropriate lodging by particles. This algorithm was suggested and developed in 1995 by a common study of Eberhart and Kennedy based on the motion of fish and particles on the basis of the two axioms of artificial life and evolution. In

similar other evolutionary algorithms, algorithm begins with a collection of particles of a matrix with a completely random position. Any particle in this matrix is called a particle, and these particles can jump in the n th - perspective space (n is the number of variables in the optimization problem). And at each step, their new situation is updated based on the previous personal experiences and the situation of their proximities. The strength of each particle of this set of particles is defined by the following vector [16-18]:

$$\mathbf{X}_i = [\mathbf{X}_{i1}, \mathbf{X}_{i2}, \dots, \mathbf{X}_{in}]^T \in S \quad (1)$$

In this regard, S is the search space and X_i is the position of each particle in the iteration i algorithm. Each particle has a velocity at any step. Therefore, the velocity vector of all particles is defined by relationship 2 [16-18]:

$$\mathbf{V}_i = [\mathbf{V}_{i1}, \mathbf{V}_{i2}, \dots, \mathbf{V}_{in}]^T \in S \quad (2)$$

The best personal position that each particle has from the beginning to step is called the best personal position and is defined for all particles by the following vector in each step [16-18]:

$$\mathbf{P}_i = [\mathbf{P}_{i1}, \mathbf{P}_{i2}, \dots, \mathbf{P}_{in}]^T \in S \quad (3)$$

Based on the relationships and definitions described above, the rate and speed of each particle at each step of repetition is calculated and updated by the following relationship [16-18]:

$$\vec{v}_i^{k+1} = w\vec{v}_i^k + c_1 r_1 \times (\vec{p}_i - \vec{x}_i^k) + c_2 r_2 \times (\vec{p}_g - \vec{x}_i^k) \quad (4)$$

$$\vec{X}_i^{k+1} = \vec{V}_i^{k+1} + \vec{X}_i^k \quad (5)$$

In this regard, the updated speed of the particle is in the iteration of $k + 1$ and the previous velocity and location of the particle respectively. It is also the best i -th particle location ever as well as the location of particle that has the p -best between particles. Here c_1 and c_2 are fixed coefficients and are usually 2. If the quantity of c_1 increases, the particle tends to follow the search around its best personal location. However, if c_2 is higher than c_1 , the tendency of the particle is to probe around the global location. Hence, it is better to assimilate the procedure of choice between these two parameters. The coefficient w is known as the inertial weight coefficient. This coefficient specifies the impact of the previous velocity on the new velocity. If the low w coefficient is c , the search step is short and consequently, the search space is small and of course, the search accuracy is increased. However, if the selected number is large, the search step and the search space for each particle will be longer but the search accuracy will be lower. r_1 and r_2 are two random numbers between zero and one that gives a random nature to the search pattern. In many cases, the w coefficient is fixed and about 0.9. However, in some cases it is linear and a function of program repetition. So first, a large search is selected to enlarge the search space at the beginning of

the search. Then, with increasing iteration pattern, its value decreases so that the further we go, the more accurate the search accuracy. Although this method gives a more accurate answer than the choice of w with a constant value, it still cannot be applicable in all engineering issues. Therefore, it is then selected by fuzzy rules in a comparative manner. If the target function is close to the optimal value, the coefficient w is small and if it is far away, the coefficient w is selected. In addition to the coefficient w , the coefficients c_1 and c_2 will be selected by comparative relationships according to what will be mentioned in the next section. As mentioned above in the particle cluster algorithm, particles are inclined to follow a search pattern that can obtain the best personal and global location at each stage. This causes premature convergence of the algorithm because the actual main optimal point may be far from these two points. To overcome this problem, in this article we use the mutation operator to prevent particles from getting entangled in the optimum local point.

d. Modification of the Classic PSO Algorithm with Mutation Operator:

As mentioned before, the particle aggregation algorithm, despite the genetic algorithm, does not have a mutation operator and always tries at every step to find the search around the two points of the best personal and overall position to continue the same step. This phenomenon can lead to two demerits. First, the algorithm may experience premature convergence. This means that it is entangled in an optimal local location, in other words, the population loses its diversity. Second, the response varies from program to program since the final response depends almost on the randomly selected primary population. Hence, in this article, we use the mutation operator to overcome these two problems. Since the mutation is a powerful tool in improving particle population diversity. In this article, we will use a new mutation operator. First, five vectors are randomly selected from the previous population in each repetition of the program (H_4, H_5, H_1, H_2, H_3), so that $H_4 \neq H_5 \neq H_1 \neq H_2 \neq H_3$. Now the jump operator selects the new position of the particle as follows:

$$X_{mut} = X_{H1} + \beta_1 (X_{H3} - X_{H2}) + \beta_2 (X_{H5} - X_{H4}) \quad (6)$$

Here, the coefficients β_1 and β_2 are supposed as mutation coefficients, the value of which should be selected experimentally in the range $0 < \beta < 1$. In the following, the value of the position of each particle of the following relation is calculated:

$$X_{new,i} = \begin{cases} X_{mut,i}, & \text{if } (rand < crossover) \\ X_i & \text{otherwise} \end{cases} \quad i = 1, 2, \dots, n \quad (7)$$

In this article, the crossover value is calculated as 0.2

e. Determining the Coefficients of c1 and c2 in a Comparative Manner:

As mentioned before, the coefficients C1 and C2 in the classical PSO algorithm are considered constant and equal to the value of 2. In some papers, these coefficients change linearly over different iterations. Increasing or decreasing these coefficients, in addition to directing the search around a particular point, can reduce or increase search space. As the matter of increasing or decreasing the weighting coefficient of inertia w. Experimental results show that such a choice for these coefficients prevents accurate response. For this reason, in this article, these coefficients are determined comparatively using the following relationship.

$$c_1, c_2 = 1 + [1 + \exp(-\frac{G_Best_Valve}{G_0})^n]^{-1} \tag{8}$$

In this case, n = 2 and G0 are equivalent to Gbest in the first repetition. Notice that the smaller the Gbest in the current repetition, the closer we get to the answer. So its value is reduced to increase the accuracy of the search. However, if Gbest is a large number, the answer is far from the optimal global answer and makes the search space bigger.

f. Fuzzy Rules for Diagnosing the Inertial Coefficient W:

The weight factor W has a huge impact on the velocity of each particle at the current stage, so increasing this factor increases the velocity. Since it is supposed that in relationship number 5, the amount of each displacement is considered one second, so the higher the velocity, the higher the particle displacement in one step, and consequently, the search space is large and its accuracy is decreased. The opposite is true. Hence, an appropriate balance must be taken into account in selecting this particle. In this article, this equilibrium is performed using fuzzy rules and ifs. The best choice is to match the w coefficient to whether Gbest is close to or far from each step of the desired Gbest using fuzzy logic. Here, the values of w and NFV, which are defined below, are the inputs of the fuzzy inference motor and its output is Δw [19-20].

$$NFV = \frac{(FV - FV_{min})}{(FV_{max} - FV_{min})} \tag{9}$$

Here FV is the Gbest level in the current step and FVmin is the Gbest level in the first repetition and FVmax is a very large number. Usually, the W coefficient must be between 0.9 and 0.4. Since the correction of the W factor during the implementation of the program may be increasing or decreasing, both positive and negative corrections are essential for this coefficient. In this research, a small number with a value of 0.1 is regarded, which is added and subtracted by the W factor.

$$\omega^{k+1} = \omega^k + \Delta\omega \tag{10}$$

Here Δw is asimilar correction value and is equal to ± 1. Of course, sometimes the value is zero and its status is

suggested according to Table 1. Notice that Gbest values must be expressed as membership functions to attain an optimal value for the weight factor W. In this article, it is recommended that triangular membership functions be selected so that they have three states:

Large or L, small or M, and medium or M. Also, the fuzzy model outputs, as shown in Table 1, have three values of PE ((+0.1, NE ((-0.1) or, ZE (0). As shown in Table 1, the 9 states may be based on different values of NFV and W occur. If both NFV and W are small, there is no need to change w because on the one hand, Gbest has reached the optimal level and on the other, hand it is not possible to decrease W so much that it excels the permutable limit. If the NFV is low and the W is medium, you can still reduce the W by 0.1 to increase the search accuracy. If the NFV is low and the W is high, you can reduce the w by 0.1 as much as before. Here the relationship between inputs and outputs is shown in Table 1. Also, the triangular membership functions are represented in Figure 4. These functions are used to get the input and output variables.

Table 1: Fuzzy rules of the input and output variables

ΔW		W		
		S	M	L
NFV	S	ZE	NE	NE
	M	PE	ZE	NE
	L	PE	ZE	NE

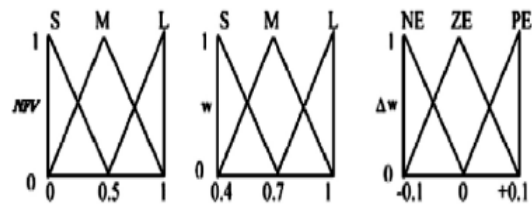


Fig. 4. The membership functions.

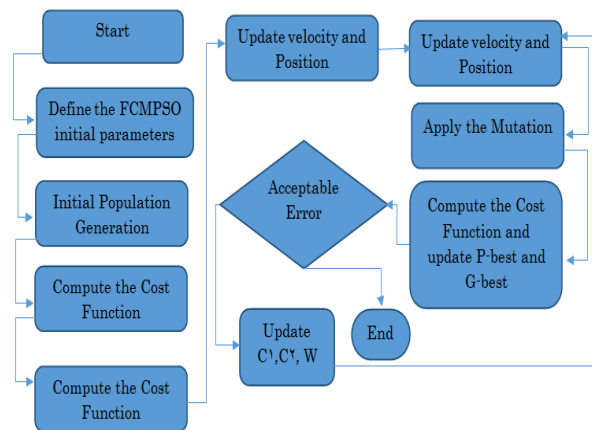


Fig. 5. Flowchart of FCMPPO algorithm

g. Experimental Data

KDDCUP99 [24] and NSL-KDD are the most commonly used datasets in intrusion detection research. We used the NSL-KDD intrusion dataset which is available in CSV format for model validation and evaluations. The dataset composes of the attacks shown in Tables 2 and 3 and identified as a key attack in IOT computing. Sherasiya and Upadhyay [25] pointed out that IOT objects are also exposed to such types of attacks, and the data that IoT objects exchange are of the same value and importance, or occasionally more important than a non-IoT counterpart.

h. The Objective Function:

In this research, to model the attack detection system and anomalies, we used the multilayer perceptron neural network structure as ML. additionally we trained the neural network using BP algorithms, classical particle algorithms, modified particle algorithms with FPSO (fuzzy PSO), FCPSO (Fuzzy comparative PSO) and FCMPPO (Fuzzy combinations. Moreover, we used the sigmoid function as the last layer of the neural network according to the following formula.

$$a(z) = \frac{1}{1 + \exp(-z)} \quad (11)$$

The accuracy of the suggested model is calculated based on the correct detection of the model attained by the neural network and by the following relationship:

$$Accuracy = \frac{\text{True Positive} + \text{True Negative}}{\text{True Positive} + \text{True Negative} + \text{False Positive} + \text{False Negative}} \quad (12)$$

Since the particle algorithm inherently minimizes the target function, the following function should be defined to increase the accuracy of the target function:

$$\text{Cost Function} = - \text{accuracy} \quad (13)$$

3- Result and Discussion

As referred to in the previous section, the PSO algorithm is a powerful algorithm for finding optimal points in complex and multi-purpose problems. Hence, in this article, the neural network has one hidden layer with 15 neurons and training is done by the PSO algorithm. However, the classic model of this algorithm has a number of coefficients that if selected consistently decrease particle diversity and premature convergence, resulting in localized optimal locations. So, in this paper, these coefficients c_1 , c_2 are comparatively diagnosed using exponential relationships. Additionally, the weighted coefficient of inertia is determined using rolls and fuzzy logic rules. Also, since this algorithm, unlike the genetic algorithm, did not have a mutation operator, it led to the search for the best personal position or the best global

position at any stage, so we suggested adding a new mutation operator to the algorithm's function. This operator is expected to curb the algorithm from getting trapped in the optimal local locations. So, we used the combination of the above methods and taught them the neural network and compared the outputs. Figure 5 shows the accuracy level for different neural network training methods. Here we suppose that the maximum repetition is equal to 50 and also the number of particles is equal to 40. As represented in the figure, neural network training by classical PSO algorithm is much more optimal than training by BP algorithm. Moreover, as expected, the classic PSO algorithm was entangled at the local optimal point, and the combination of FPSO, FCPSO, and FCMPPO gave more accurate responses. Also, the combination of the mutation operator with the classic PSO algorithm gives good results. Figure 7 shows the convergence speed of different algorithms drowned on the iteration of the algorithm for diagnosing different attacks. As shown in the figure, the FCMPPO algorithm, in addition to being much more accurate, has a better convergence pace. So, this algorithm is a very optimal algorithm to increase the accuracy and speed of attack detection.

From the Dos detection picture, we can see that when the number of trainings exceed 35 times, the Classic ANN curve is basically stable, and with the increasing of the number of trainings, the accuracy rate no longer increases significantly. In this method, the performance accuracy of the algorithm does not exceed 74%. In contrast to this method is the ANN-FCMPPO algorithm. This method has higher accuracy (99%) and achieves faster response. As shown in this figure convergence point is 26 and the point of this sentence is that FCMPPO algorithm is faster than the previous algorithm. Moreover figure 7 shows that by applying any corrective methodology in PSO algorithms such as ANN FPSO and ANN FCPSO, accuracy and convergence speed improved simultaneously. Another matter is that among the four attack type, these methods give the best performance to the Dos attack.

Lastly, in Figure 8, we show the accuracy of the PSO and FCMPPO algorithms after running the program 20 times to detect Dos attacks. As shown in the figure, the FCMPPO algorithm is more dependable than the PSO algorithm. Since in different performances, the program represents relatively the same answers.

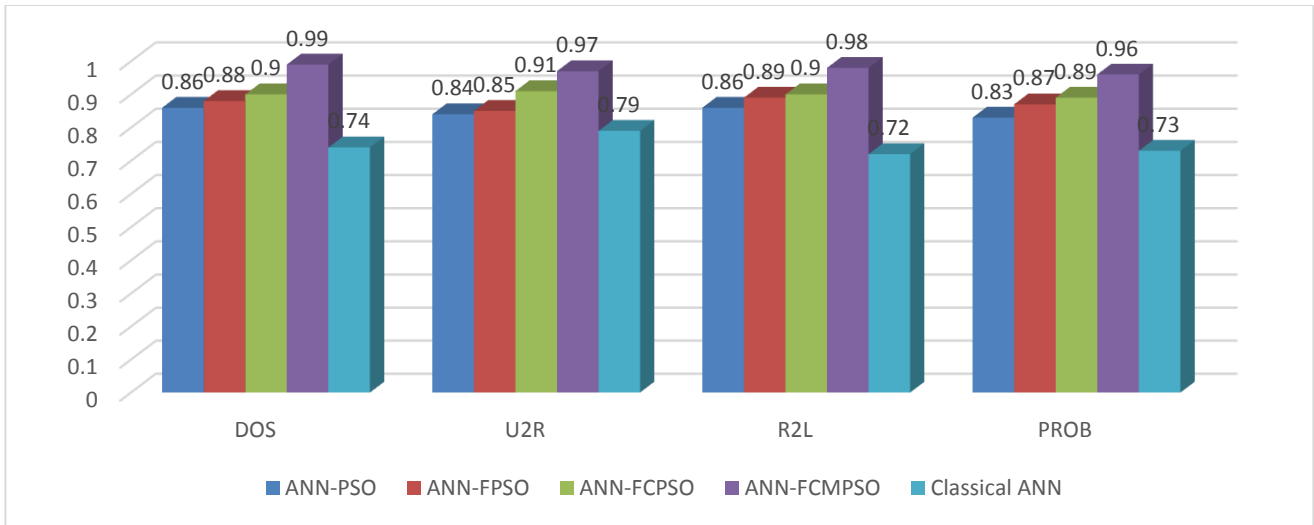


Fig. 6: accuracy for different machine learning algorithm

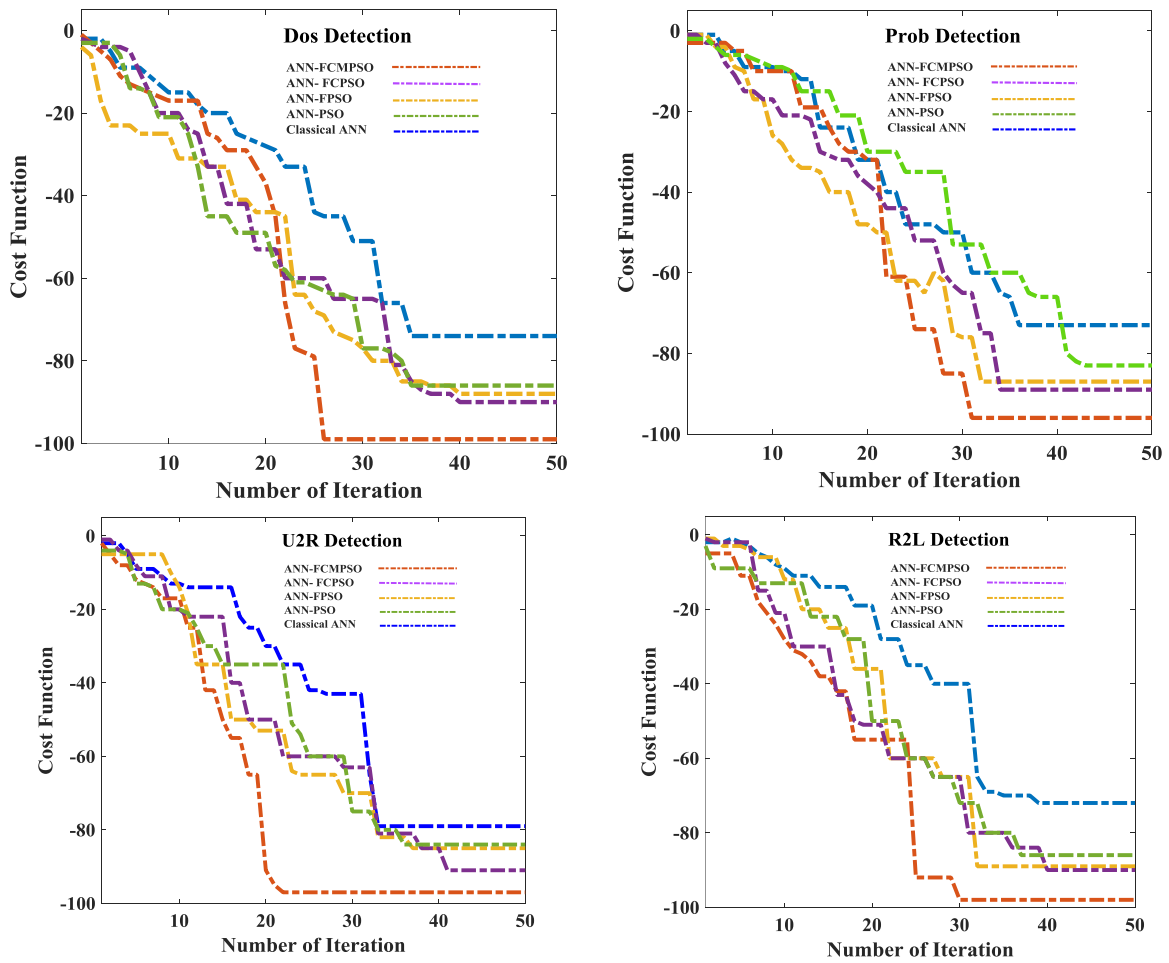


Fig. 7. convergence characteristic of proposed method in different attack detection

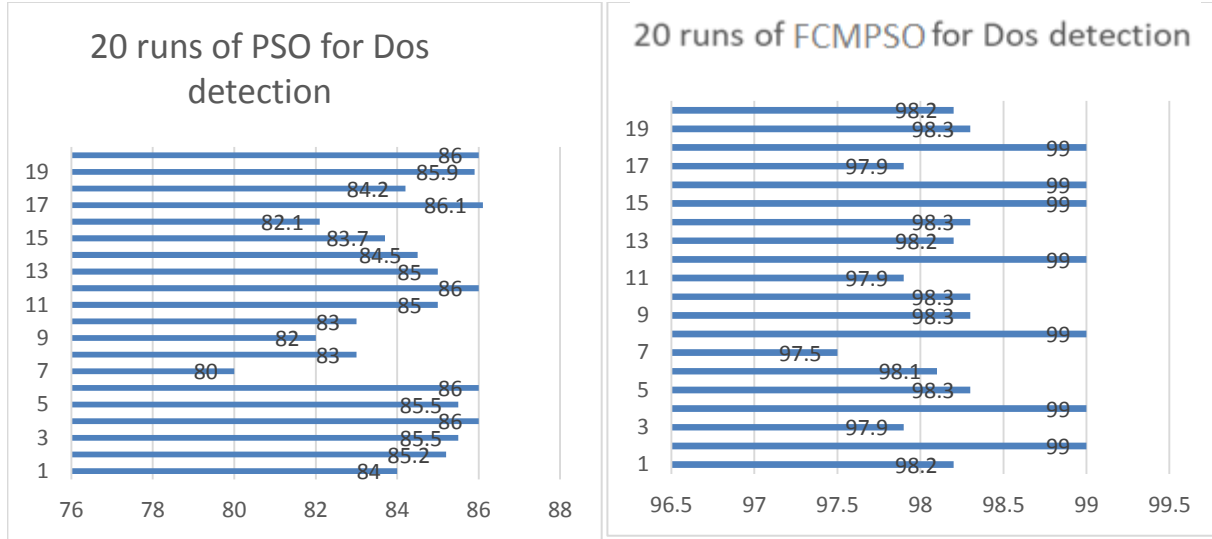


Fig. 8. Accuracy for 20 runs of the left algorithm: ANN-PSO right: ANN FCMPs

Table 2: Input parameters of Neural Network

S/N	Name	Type	S/N	Name	Type
1	duration	Continuous	25	error_rate	Continuous
2	protocol_type	Symbolic	26	srv_error_rate	Continuous
3	service	Symbolic	27	error_rate	Continuous
4	flag	Symbolic	28	srv_error_rate	Continuous
5	src_bytes	Continuous	29	same_srv_rate	Continuous
6	dst_bytes	Continuous	30	diff_srv_rate	Continuous
7	land	Symbolic	31	srv_diff_host_rate	Continuous
8	wrong_fragment	Continuous	32	dst_host_count	Continuous
9	urgent	Continuous	33	dst_host_srv_count	Continuous
10	hot	Continuous	34	dst_host_same_srv_rate	Continuous
11	num_failed_logins	Continuous	35	dst_host_diff_srv_rate	Continuous
12	logged_in	Symbolic	36	dst_host_same_src_port_rate	Continuous
13	num_compromised	Continuous	37	dst_host_srv_diff_host_rate	Continuous
14	root_shell	Continuous	38	dst_host_error_rate	Continuous
15	su_attempted	Continuous	39	dst_host_srv_error_rate	Continuous
16	num_root	Continuous	40	dst_host_error_rate	Continuous
17	num_file_creations	Continuous	41	dst_host_srv_error_rate	Continuous
18	num_shells	Continuous			
19	num_access_files	Continuous			
20	num_outbound_cmds	Continuous			

Table 3: Output Parameters of Neural Network (Attack Type)

S/N	Name	Type
1.	Back	dos
2.	buffer_overflow	u2r
3.	ftp_write	r2l
4.	guess_passwd	r2l
5.	imap	r2l
6.	ipsweep	probe
7.	land	dos
8.	loadmodule	u2r
9.	multihop	r2l
10.	neptune	dos

S/N	Name	Type
11.	nmap	probe
12.	perl	u2r
13.	phf	r2l
14.	pod	dos
15.	portsweep	probe
16.	rootkit	u2r
17.	satant	probe
18.	smurf	dos
19.	spy	r2l
20.	teardrop	dos
21.	warezclient	r2l
22.	warezmaster	r2l

4-Conclusion

In this study, we used modified PSO and PSO algorithms to train the neural network to model the IOT network attack detection. We showed that meta-heuristic algorithms can be a more effective method than classical education systems. In addition, we have shown that the PSO algorithm has coefficients that, if not properly adjusted, lose their efficiency and cannot be suitable for neural network training methods. The correction model proposed in this paper is the simultaneous combination of a PSO algorithm with a fuzzy system and a mutational and adaptive operator. The suggested ANN-FCMPSO algorithm is about 97% (99% for Dos type attack, 97% for U2R, 98% for R2L and 96% for PROB), and the accuracy for the PSO-ANN algorithm is about 86%.

References

- 1- Haji, Saad Hikmat, and Siddeeq Y. Ameen. "Attack and anomaly detection in IOT networks using machine learning techniques: A review." *Asian Journal of Research in Computer Science* (2021): 30-46.
- 2- Aversano, Lerina, et al. "Effective Anomaly Detection Using Deep Learning in IoT Systems." *Wireless Communications and Mobile Computing 2021* (2021). (+)
- 3- Khan, Arshiya, and Chase Cotton. "Detecting Attacks on IoT Devices using Featureless 1D-CNN." 2021 IEEE International Conference on Cyber Security and Resilience (CSR). IEEE, 2021.
- 4- Bello, Ibrahim, et al. "Detecting ransomware attacks using intelligent algorithms: recent development and next direction from deep learning and big data perspectives." *Journal of Ambient Intelligence and Humanized Computing* 12.9 (2021): 8699-8717.
- 5- Foley, John, Naghmeh Moradpoor, and Henry Ochen. "Employing a Machine Learning Approach to Detect Combined Internet of Things Attacks against Two Objective Functions Using a Novel Dataset." *Security and Communication Networks 2020* (2020).
- 6- Ullah, Imtiaz, and Qusay H. Mahmoud. "Design and development of a deep learning-based model for anomaly detection in IoT networks." *IEEE Access* 9 (2021): 103906-103926.
- 7- Syed, Naeem Firdous, et al. "Denial of service attack detection through machine learning for the IOT." *Journal of Information and Telecommunication* (2020): 1-22.
- 8- Manimurugan, S., et al. "Effective Attack Detection in Internet of Medical Things Smart Environment Using a Deep Belief Neural Network." *IEEE Access* 8 (2020): 77396-77404.
- 9- Churcher, Andrew, et al. "An experimental analysis of attack classification using machine learning in iot networks." *Sensors* 21.2 (2021): 446. (+)
- 10- Latif, Shahid, et al. "A Novel Attack Detection Scheme for the Industrial Internet of Things Using a Lightweight Random Neural Network." *IEEE Access* 8 (2020): 89337-89350.
- 11- Alkronz, EyadSameh, et al. "Prediction of Whether Mushroom is Edible or Poisonous Using Back-propagation Neural Network." (2019).
- 12- Wang, Weilin, et al. "Estimation of PM2.5 concentrations in China using a spatial back propagation neural network." *Scientific reports* 9.1 (2019): 1-10.
- 13- Mohammadi, Farzaneh, et al. "Modelling and optimizing pyrene removal from the soil by phytoremediation using response surface methodology, artificial neural networks, and genetic algorithm." *Chemosphere* 237 (2019): 124486.
- 14- Azimi, Yousef, Seyed Hasan Khoshrou, and MortezaOsanloo. "Prediction of blast induced ground vibration (BIGV) of quarry mining using hybrid genetic algorithm optimized artificial neural network." *Measurement* 147 (2019): 106874.
- 15- Cai, Jianghui, et al. "A Novel Clustering Algorithm Based on DPC and PSO." *IEEE Access* 8 (2020): 88200-88214.
- 16- Singh, Shakti, Prachi Chauhan, and NirbhawJap Singh. "Capacity optimization of grid connected solar/fuel cell energy system using hybrid ABC-PSO algorithm." *International Journal of Hydrogen Energy* (2020).
- 17- Devarasiddappa, D., M. Chandrasekaran, and R. Arunachalam. "Experimental investigation and parametric optimization for minimizing surface roughness during WEDM of Ti6Al4V alloy using modified TLBO algorithm." *Journal of the Brazilian Society of Mechanical Sciences and Engineering* 42.3 (2020): 1-18.
- 18- Qiao, Weibiao, Hossein Moayedi, and Loke KokFoong. "Nature-inspired hybrid techniques of IWO, DA, ES, GA, and ICA, validated through a k-fold validation process predicting monthly natural gas consumption." *Energy and Buildings* (2020): 110023.
- 19- Prithi, S., and S. Sumathi. "LD2FA-PSO: A novel Learning Dynamic Deterministic Finite Automata with PSO algorithm for secured energy efficient routing in Wireless Sensor Network." *Ad Hoc Networks* 97 (2020): 102024.
- 20- Kacimi, MohandAkli, et al. "New mixed-coding PSO algorithm for a self-adaptive and automatic learning of Mamdani fuzzy rules." *Engineering Applications of Artificial Intelligence* 89 (2020): 103417.
- 21- Jallal, Mohammed Ali, Samira Chabaa, and AbdelouhabZeroual. "A novel deep neural network based on randomly occurring distributed delayed PSO algorithm for monitoring the energy produced by four dual-axis solar trackers." *Renewable Energy* 149 (2020): 1182-1196.
- 22- Niknam, Taher, Ehsan Azadfarsani, and Masoud Jabbari. "A new hybrid evolutionary algorithm based on new fuzzy adaptive PSO and NM algorithms for distribution feeder reconfiguration." *Energy Conversion and Management* 54.1 (2012): 7-16.
- 23- Niknam, Taher, Hassan DoagouMojarrad, and Majid Nayeripour. "A new fuzzy adaptive particle swarm optimization for non-smooth economic dispatch." *Energy* 35.4 (2010): 1764-1778.
- 24- M. Tavallae, E. Bagheri, W. Lu, and A. Ghorbani, "A detailed analysis of the KDD CUP 99 data set," Proc. 2nd IEEE Symp. Comput. Intell. Secur. Defense Appl. (CISDA), Ottawa, ON, Canada, Jul. 2009, pp. 1-6.
- 25- A, Alghuried, "A model for anomalies detection in Internet of Things (IoT) using inverse weight clustering and decision tree," M.S. thesis, School Comput., Dublin Inst. Technol., Dublin, Republic of Ireland, 2017.

ARASP: An ASIP Processor for Automated Reversible Logic Synthesis

Zeinab Kalantari^{1*}, Marzieh Gerami², Mohammad Eshghi³

¹. Department of Computer Engineering, Rafsanjan Branch, Islamic Azad University, Rafsanjan, Iran

². Department of Computer Engineering, ShahreKord Branch, Islamic Azad University, ShahreKord, Iran

³. Faculty of Electrical and Computer Engineering, Shahid Beheshti University, Tehran, Iran

Received: 04 Jun 2021/ Revised: 04 Jan 2022/ Accepted: 02 Feb 2022

Abstract

Reversible logic has been emerged as a promising computing paradigm to design low power circuits in recent years. The synthesis of reversible circuits is very different from that of non-reversible circuits. Many researchers are studying methods for synthesizing reversible combinational logic. Some automated reversible logic synthesis methods use optimization algorithms. Optimization algorithms are used in some automated reversible logic synthesis techniques. In these methods, the process of finding a circuit for a given function is a very time-consuming task, so it's better to design a processor which speeds up the process of synthesis. Application specific instruction set processors (ASIP) can benefit the advantages of both custom ASIC chips and general DSP chips. In this paper, a new architecture for automatic reversible logic synthesis based on an Application Specific Instruction set Processors is presented. The essential purpose of the design was to provide the programmability with the specific necessary instructions for automated synthesis reversible. Our proposed processor that we referred to as ARASP is a 16-bit processor with a total of 47 instructions, which some specific instruction has been set for automated synthesis reversible circuits. ARASP is specialized for automated synthesis of reversible circuits using Genetic optimization algorithms. All major components of the design are comprehensively discussed within the processor core. The set of instructions is provided in the Register Transform Language completely. Afterward, the VHDL code is used to test the proposed architecture.

Keywords: Reversible logic; Optimization Algorithms; Application Specific Instruction Set Processors; ASIP; RTL.

1- Introduction

Application specific instruction set processors (ASIP) can compromise the advantages of custom ASIC chips and general DSP chips. In other words, ASIP chips utilize high performance and low power of ASIC chips and flexibility of DSP chips [1][2][3][4][5].

There is a tradeoff between cost and speed in ASIPs.

Programmability is the main advantage of ASIPs, which gives more flexibility to software developers. Other advantages are more convenient in the design and debugging process, predictability, and shorter time to market. Hardware and software are two aspects of ASIPs rather than one aspect of the task being dominant. Besides, compared to general-purpose processors, ASIP benefits from having specific instructions to perform a specific task faster and reduce programmer errors. Accordingly,

efficiency and programmability are both advantages of ASIPs compared to general-purpose processors.

In this paper, a novel ASIP-based processor for the synthesis of reversible circuits is proposed. This processor is used to synthesize reversible circuits using optimization algorithms. VHDL code is used to simulate and test the proposed architecture. The main objective of the proposed design is programmability as it is the major concept of ASIP. In addition, the suggested structure reduces hardware complexity.

The organization of the rest of the paper is as follows. The next section and subsequent sections present a background on the synthesis of reversible circuits. Section 3 details the proposed ASIP architecture model for the synthesis of reversible circuits. The testing process describes in section 4. Finally, section 5 concludes the paper.

2- Background

Reversible logic has applications in low power computing, quantum computing, nanotechnology, optical computing, and DNA computing. The design of the reversible circuits is quietly different from the design of conventional irreversible logic circuits [6] because of the different gates that are available in reversible logic.

The synthesis of reversible circuits differs significantly from synthesis using traditional irreversible gates. Many algorithms have been proposed for the synthesis of reversible circuits [7][8][9][10]. Dengli et al. proposed an improved KFDD based reversible circuit synthesis method [7]. Ahmed et al. suggested a synthesis approach using reorder algorithm [8]. Basak et al. presented an algorithm using the ESOP expressions [9]. Some automated reversible logic synthesis methods, such as genetic algorithms (GAs) are also presented [11][12][13][14][15][16]. These algorithms use optimization algorithms.

In optimization algorithms, the main process is generating some random circuits and then computing the output truth table of generated circuits. Then the hamming distance between the truth table of generated random circuits and the truth table of the given function is calculated. After that, according to the optimization algorithm the best circuit is selected. These operations are repeated while desired hamming distance is reached.

To implement the automated reversible synthesis algorithm, a background on reversible gates is needed. The next section is illustrated to introduce reversible logic gates. After that, a general algorithm for the synthesis of all reversible circuits is presented.

2-1- Reversible Gates

Since a serious problem in modern VLSI designs is power consumption, Low power circuit design is one of the most attractive subjects for hardware designers. Landauer has shown that for irreversible logic computations, each bit of information lost, generates $kT \ln 2$ joules of heat energy, where k is Boltzmann's constant and T is the absolute temperature at which computation is performed [17]. Bennett showed that $kT \ln 2$ energy dissipation would not occur if a computation is carried out reversibly [18]. This part of energy dissipation is independent of what the underlying technology is.

In reversible circuits, no bit of information is lost, and reversible computation in a system can be performed only when the system comprises reversible gates.

In a reversible gate, there is a one-to-one correspondence between its inputs and outputs. As a result, the number of outputs of a reversible gate is the same as the number of inputs, and for each input vector, there is a unique output vector and vice versa.

Some more common gates to design reversible logic circuits are Feynman Gate, FG [19], Toffoli Gate, TG [20], Fredkin Gate, FRG [21] are more common gates to design reversible circuits.

A 2*2 Feynman Gate, also known as controlled-NOT (CNOT), is depicted in Fig.1.a. It implements the logic functions: $P = A$ and $Q = A \otimes B$.

A 3*3 Toffoli Gate has 3 inputs: 2 control inputs, that are copied to the first 2 outputs and one other input that is complemented if all control inputs are 1s and are directly copied to the last output otherwise [20]. A 3- input, 3- output Toffoli Gate is shown in Fig.1.b. The inputs 'A' and 'B' are passed as first and second outputs, respectively. The third output is controlled by 'A' and 'B' to invert 'C'.

A 3*3 Fredkin Gate is depicted in Fig.1.c. Here the input 'A' is passed as the first output. Inputs 'B' and 'C' are swapped to get the second and third outputs, which are controlled by 'A'. If $A = 0$, then the outputs are simply duplicating of the inputs; otherwise, if $A = 1$, then the two input lines (B and C) are swapped.

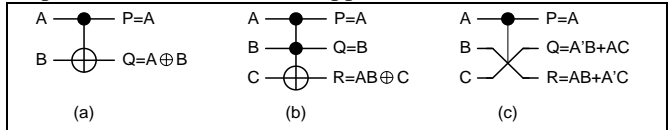


Fig. 1 (a) Feynman gate, (b) Toffoli gate, and (c) Fredkin gate

2-2- Synthesis of Reversible Circuits

Synthesizing a reversible circuit using a searching algorithm is a complex problem with a large amount of searching space. The number of combinations for placing one Toffoli $r \times r$ gate in an $n \times n$ circuit ($0 < r < n-1$) to synthesize a reversible $n \times n$ circuit is expressed by (1):

$$\rho = n \cdot \sum_{r=0}^{n-1} \binom{n-1}{r} = n \cdot 2^{n-1} \quad (1)$$

If the number of required gates to design a circuit is m , then the number of possible circuits is ρ^m . If Fredkin and Press are added to the set of Toffoli $r \times r$ gates, then the number of possible circuits is $(3\rho)^m$. So optimization algorithms, especially GA, are used to find the global minimum or maximum of a function, in an extensive searching space [11][12][13][14]. In the next subsection, a review of automated synthesis is presented.

2-3- Automated Reversible Logic Synthesis

The general algorithm for automated reversible logic synthesis is shown in Fig.2. This processor is used to synthesize reversible circuits.

The algorithm starts by generating a random configuration (a random state) of a circuit with one gate. We consider the hamming distance between the truth table of this circuit and the truth table of a given function as the cost function. Then for each new configuration, it's necessary

to compute the truth table. So the main operations in such algorithms are computing the truth table of each circuit and then comparing its truth table with the destination truth table to select the best circuit, based on the desired algorithm. So in a given iteration, the algorithm generates n new circuits at a time. Each new circuit is derived from the old configuration. The hamming distances of the two circuits are then compared.

After the process of finding a new circuit, comparing it to the current configuration, and either accepting or rejecting, it is done n times.

After all, if the stopping criteria of the algorithm, zero hamming distance, is not reached, the number of gates will increase and the algorithm is repeated for a new circuit with extra gates.

Set NoG=1 //the number of gates being used for the synthesis

Set S=S₀ // random initial state

Loop1:

Initiate a random circuit S using NoG

While (up to max-iteration)

```

{
  While(required number of circuits not generated)
  {
    Generate new circuit S' by perturbing S;
    For (all rows of truth-table)
      E=E+HD(des[i]&mask , Syn[i]&mask);
    ΔE=E(S')-E(S)
    If (ΔE<=0)
      S=S';
  }
  If (HD!=0)
  {
    NoG=NoG+1;
    Goto Loop1;
  }
}
else
  Print circuit; } }
    
```

Fig. 2 Algorithm description

This paper introduces an ASIP processor which is useful in the application domain of reversible circuits. This processor is called ARASP, an ASIP processor for automated reversible logic synthesis. The schematic of the ARASP processor is illustrated in Figure 2. The explanation of the proposed register configuration is in Figure 3.

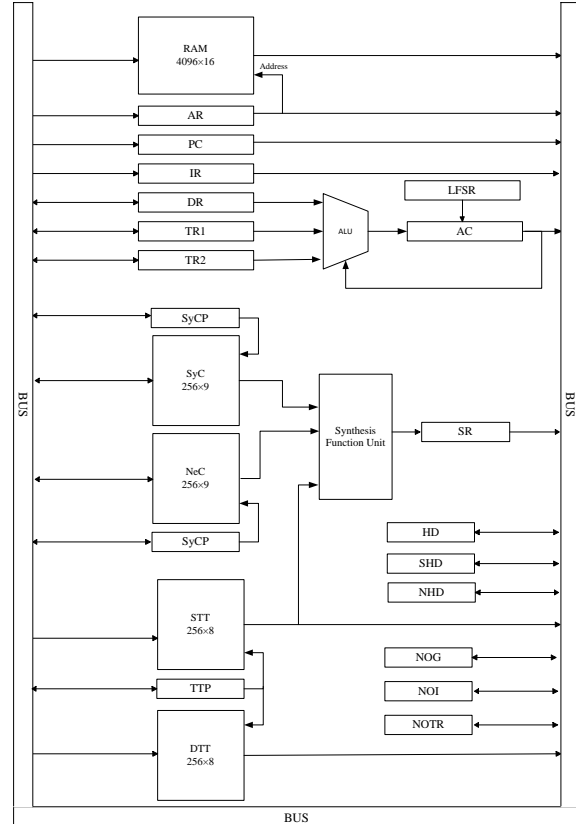


Fig. 4 ARASP register configuration

We proposed a 16-bit wide ASIP architecture with 16-bit integer operations for common arithmetic operations and a specific function unit for the automated synthesis of a reversible circuit.

3-1- Global View on ARASP

Registers have clock inputs that are all connected to the main system clock. Each AC, DR, TR₁, and TR₂ are 16-bit registers that provide operands of ALU. The output of ALU is connected only to the AC. The instruction register, IR, provides the instruction bits for the controller. The 12-bit program counter register is called PC, this register provides an address for current instruction through the memory address register, AR. This register also is a 12-bit binary up counter. The arithmetic logic unit, ALU, is a combinational logic unit with two 16-bit inputs, four flag inputs, and control inputs that specify the integer operations. The output of this unit is connected to the input of AC.

3- Architecture Overview

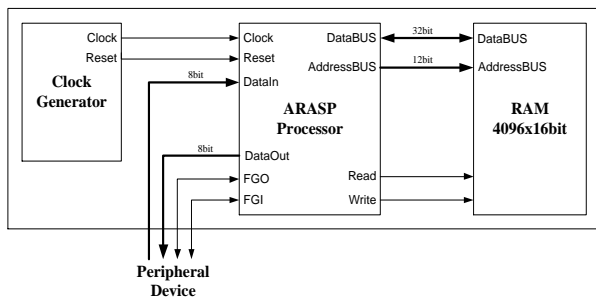


Fig. 3. ARASP processor

The next part is a function unit for the automated synthesis of a reversible circuit. This unit works as follows. At first, a random circuit with some gates that the NOG (Number OF Gates) register specifies is generated. The SyC (Synthesis Circuit) register bank maintains this random gate. The SyCP (Synthesis Circuit Pointer) register which is an 8-bit up-counter is used to specify each register of this register bank (each reversible gate) in each step. Same as SyC, NeC register bank is used to hold the neighboring circuit of synthesis circuit in each step, the size of this register is as SyC, also NeCP (Neighboring Circuit Pointer) is used to define a register of NeC register bank in each step of the synthesis.

DTT (Destination Truth Table) consists of 256 8-bit registers to maintain the truth table of a given function. As the algorithm generates a reversible circuit in each step, we need a register bank to maintain the truth table of the generated circuit. So STT (Synthesis Truth Table) is used for this aim.

The size of this register bank is the same as DTT. TTP is an 8-bit up-counter that refers to each row of DTT or STT in each step of the algorithm.

As said before the major time-consuming operation in the synthesis of a reversible circuit is calculating the output of the circuit for all combinations of inputs. So the synthesis function unit that is a combinational circuit calculates each row of synthesis truth table of desired circuit, SyC or NeC. After synthesis of the desired circuit, we need to compare this truth table with the destination truth table, DTT. So we have to calculate the hamming distance between synthesis truth table STT and destination truth table DTT, as the cost function. After that, the calculated hamming distance will set the SHD or NHD depending on the circuit that is synthesized. The final step in the algorithm is selecting the best circuit. So we need to compare SHD and NHD registers and set HD register with one of these registers.

4- Proposed Architecture

In this paper, the proposed CPU is referred to as ARASP. The proposed processor employs a reduced hardware requirement and application specific instruction set. Due to the size of its data register and buses, ARASP is considered to be a 16-bit processor. It has direct and indirect addressing modes. ARASP also has specific instructions and input-output interrupts.

4-1- Main Memory Organization

The ARASP is capable of addressing 4096 bytes of memory through its 12-bit address lines. This memory is addressed by a register called AR.

4-2- Register Configurations

The main data register of ARASP is AC, which is used in conjunction with most general instructions. This processor

has overflow, carry, zero, and sign flags (o, c, z, and s). These flags may be modified by arithmetic operations.

ARASP consists of two parts, global unit, and specific unit. The major components of the global unit are AR, PC, IR, DR, TR₁, TR₂, AC, LFSR, and ALU. Also, the components of the specific part are SyC and NeC register bank that consists of 256 9 bit registers. These register banks hold synthesis and neighboring circuits each consisting of at most 8 gates. DTT and STT register banks hold destination and synthesis truth tables respectively. According to the size of these register banks and hardware restrictions. The desired circuit can have at most 8 inputs (256 8-bit registers). HD, SHD, and HD registers that are 8-bit registers are used for holding hamming distance of the circuit throughout the running synthesis algorithm.

4-3- Instruction Types

The ARASP has a total of 47 instructions totally, and the specific instructions are summarized in Table 3. The Proposed processor has two different types of instruction sets (Table 1). The Memory reference instructions need the main memory address to do their operations and the Non-memory reference instruction set, which needs no memory for their operands. The ARASP'S memory instruction set can be used by direct and indirect addressing modes.

Table 1. Instruction Types and Addressing Modes of ARASP

Instruction Type	M	I	Address	Addressing Mode
Memory	1	0	No	Direct
Memory	1	1	No	Indirect
Others	0	×	Yes	-

As presented in Fig. 5, in memory reference instructions most significant bit of instruction (bit 15) is set, to specify the memory reference instruction type. Bit 14 called I, specifies direct or indirect addressing mode (0 for direct and 1 for indirect). The next 4 bits (bits 10-13) specify the operation of a memory reference instruction. As this type of instruction need a memory word for holding one of the operands, in these types of instructions we should refer to the main memory to read the operand. If it is set to 1 the operand's address is indirect and if it is set to 0 the operand's address is direct.

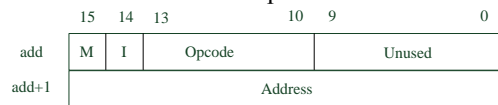


Fig. 5 Memory reference instruction format

These types of instructions occupy a byte whose most significant bit (bit 15) is 0. In this type of instruction, bit 14 specifies output and register instructions or specific instructions (0 for output and register instructions and 1 for specific instructions). The other 4 bits specify operations of instructions (Fig. 6).

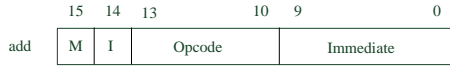


Fig. 6. Non Memory reference instruction format

The fetch, decode and calculation of effective address phases of the instruction cycle could be as follow:

Interrupt :

$$\text{IEN}(\text{FGI}+\text{FGO}) : \text{R} \leftarrow 1$$

$$\text{RT}_0 : \text{AR} \leftarrow 0$$

$$\text{RT}_1 : \text{M}[\text{AR}] \leftarrow \text{PC}, \text{PC} \leftarrow 0$$

$$\text{RT}_2 : \text{PC} \leftarrow \text{PC}+1, \text{IEN} \leftarrow 0, \text{R} \leftarrow 0, \text{SC} \leftarrow 0$$

Fetch :

$$\overline{\text{RT}}_0 : \text{AR} \leftarrow \text{PC}, \text{PC} \leftarrow \text{PC}+1$$

$$\overline{\text{RT}}_1 : \text{IR} \leftarrow \text{M}[\text{AR}]$$

Decode :

$$\overline{\text{RT}}_2 : \text{D}_{31} \dots \text{D}_0 \leftarrow \text{IR}[10-14], \text{AR} \leftarrow \text{IR}[0-8], \text{F} \leftarrow \text{IR}(9), \text{M} \leftarrow \text{IR}(15)$$

Address Fetch :

$$\text{MT}_3 : \text{AR} \leftarrow \text{PC}, \text{PC} \leftarrow \text{PC}+1$$

$$\text{MT}_4 : \text{AR} \leftarrow \text{M}[\text{AR}]$$

$$\text{M} \overline{\text{I}} \text{T}_5 : \text{nothing}$$

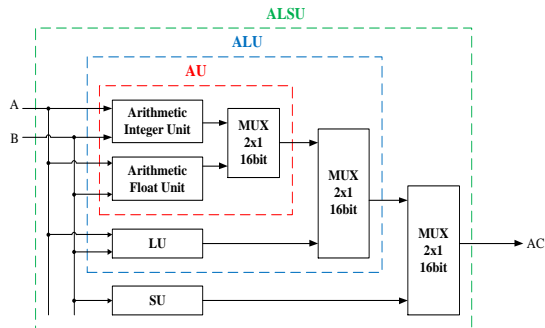
$$\text{M} \text{I} \text{T}_5 : \text{AR} \leftarrow \text{M}[\text{AR}]$$
4-4- Arithmetic Logic and Shift Unit

Fig. 7 Arithmetic Logic and Shift Unit

The presented processor supports basic arithmetic, logic, and shift units which are presented in

Table 2. ALU Operators

S4	S3	S2	S1	S0	Unit	Operation	Function
0	0	0	0	0	AU	ADD	A+B
0	0	0	0	1		SUB	A-B
0	0	0	1	0		DEC	A-1
0	0	0	1	1		INC	A+1
0	0	1	×	×		MUL	A×B
0	1	0	0	0	LU	AND	A∧B
0	1	0	0	1		OR	A∨B
0	1	1	0	0		XOR	A⊕B
0	1	1	1	0		NOT	¬B
0	1	1	1	1		PASS	A
1	×	×	0	0	SU	SHL	SHL(B)
1	×	×	0	1		SHR	SHR(B)
1	×	×	1	0		ROL	ROL(B)
1	×	×	1	1		ROR	ROR(B)

The 5-bit opcode (S_0 to S_4) hierarchically selects the proper operation. Besides the main results, five arithmetic flags (Carry, Overflow, Zero, and Sign) are set by the ALU. Each flag obtains the proper value by the Eq.s 2 to 5 considering that the input values are unsigned integer. AC, DR, TR1, and TR2 can be considered as both first and second operands.

C-Flag = Cout when ($\text{Op.}=\text{ADD}|\text{SUB}|\text{DEC}|\text{INC}|\text{SHL}|\text{SHR}$) (2)

O-Flag = '1' when ($\text{Op.}=\text{ADD}|\text{INC}|\text{SHL} \& C_{\text{out}}='1'$) | ($\text{Op.}=\text{SUB}|\text{DEC} \& C_{\text{out}}='0'$) | ($\text{Op.}=\text{MUL} \& 16\text{-bit MSB} \neq 0$) (3)

Z-Flag = '1' when (16-bit LSB=0) (4)

S-Flag = '1' when ($\text{Op.}=\text{SUB}|\text{DEC} \& C_{\text{out}}='0'$) (5)

4-5- Instruction Set

The instruction set for the ARASP is depicted in Table3.

The objective of synthesizing a reversible circuit is to compute a circuit for a given function. So we need to have a destination truth table of a given function. LDTT instruction reads the truth table of the desired function from the main memory to the DTT register bank. This instruction is a specific instruction that needs to refer to the main memory for its operation. Because of the hardware restrictions, it assumes that each function could have at most eight inputs. So the process of reading the rows of the destination truth table from memory reads some rows with the number that NoTR defines. The value of this register is set by SNOTTR instruction according to the value of the NOG register.

The STRC instruction stores the generated circuit of the synthesis process of a given function in the main memory.

The RAND instruction generates a 16-bit random number.

The CLRNOG instruction clears 8-bits of the NOG register, while INCNOG instruction increments the value of this register.

The SNOI instruction initializes the NOI register by an immediate number.

The GRNDC instruction generates a random circuit with some gates that are determined by the value of the NOG register. This instruction gets the value of the NOG register and the number of gates and sets a random value to some registers of SyC register bank according to the specified value of NOG.

The GNBRCROS instruction performs the crossover operation on a circuit to generate a new circuit called neighboring circuit from synthesis circuit. The instruction selects two random gates from the synthesis circuit that are in the SyC register bank. After that, it exchanges the position of these two gates to generate a new circuit called a neighboring circuit that is placed on the NeC register bank.

The other instruction that performs mutation operation is GNB RMUT. This instruction also generates a random new circuit from the existing circuit by mutation operation. The instruction selects a random gate from the synthesis circuit that is in the SyC register bank. After that, it exchanges the

position of control and main inputs to generate a new circuit called a neighboring circuit that is also placed on the NeC register bank.

To compute the truth table for each circuit it is necessary to set each row of the truth table with initial values 0 to $2^n - 1$. To achieve this goal, SIVDTT is used.

MASK instruction generates a mask pattern for output.

CALCSTT instruction computes the truth table of each generated circuit, SyC or NeC. In other words, this instruction, compute the value of the output for each combination of inputs in a circuit with some gate. By feeding the initial value of a given row of the truth table to the circuit as the first stage and synthesizing them the computation operation starts and then synthesis ALU calculates the output of this gate. The next gate gets the calculated output of the preceding gate and calculates the output. These operations continue while the output of the last gate is computed. This output will be replaced by the value of the given row. CALSTT repeats these operations for all combinations of inputs (all rows of truth table).

To calculate the hamming distance between syntheses or neighboring truth tables and destination truth tables, CALCHD instruction is used.

The *SBC instruction* selects the best circuits, a circuit with less hamming distance, between a circuit and a neighboring circuit.

Finally, SIZHD is used to determine the zero value of the HD register.

NOGTAC instruction is used to transfer the value of the NOG register to the AC.

Table 2. ALU Operators

S ₄	S ₃	S ₂	S ₁	S ₀	Unit	Operation	Function
0	0	0	0	0	AU	ADD	A+B
0	0	0	0	1		SUB	A-B
0	0	0	1	0		DEC	A-1
0	0	0	1	1		INC	A+1
0	0	1	×	×		MUL	A×B
0	1	0	0	0	LU	AND	A∧B
0	1	0	0	1		OR	A∨B
0	1	1	0	0		XOR	A⊕B
0	1	1	1	0		NOT	¬B
0	1	1	1	1		PASS	A
1	×	×	0	0	SU	SHL	SHL(B)
1	×	×	0	1		SHR	SHR(B)
1	×	×	1	0		ROL	ROL(B)
1	×	×	1	1		ROR	ROR(B)

5- Testing Process

In this paper, the proposed CPU is referred to as ARASP. The proposed

A structural VHDL code in Fig. 8 is used to test and verify the functionality of the given structure. Using the instruction set of the presented ARASP processor, the following code has to be programmed to generate a

random circuit for the desired N×N function which is stored from the memory address Addr1.

```

CLRNOG
INCNOG
NOI N
SNOTTR
LDDTT Addr1
GRNDC
SISTT
CLCSTT
CLCHD 0
GNBRMUT
SIVSTT
CLCSTT
CLCHD 1
SBC

```

Fig. 8 VHDL code

6- Conclusion

We showed that synthesizing a reversible circuit using a search algorithm is a complex task with a large number of search spaces. So optimization algorithms, especially Genetic Algorithm, GA, are used to find the global minimum or maximum of a function, in an extensive searching space. As in such algorithms, the process of calculating values of outputs is a time-consuming operation. So, we need a processor to speed up the process of synthesis. As a result, application specific flexibility is mandatory to meet the performance requirements of synthesis reversible circuits.

In this paper, we presented a novel design of the family of ASIP processors in the application domain of reversible circuits. The Providing programmability together with required specific instructions has been the main purpose of the automated synthesis of reversible circuits. The proposed processor that we referred to as ARASP is a 16-bit processor with a total of 47 instructions totally, which some specific instruction has set for automated synthesis reversible circuits. ARASP is specialized for automated synthesis of reversible circuits using optimization algorithms such as GA or simulated annealing.

The design steps of all the main components inside the processor core have been described in detail. Maximum specific instruction, GNBRMUT, needs 29 clock cycles for execution. Structural VHDL code has been used to test the proposed architecture. A pipeline technique could be used to enhance the speed and achieve a high throughput rate as future work.

As future work, the processor can be comprehensively implemented of this processor that will specialize in simulated annealing algorithm. It is suggested that the proposed work will provide a new focus in the reversible field making hardware more specific for such applications.

Appendix

Table 3. ARASP instruction set

I_i	Instruction	Name	Description	Ins. Reference	IR(9)	OpCode
I_0	INP	Input	$AC(L) \leftarrow INPR$	I/O		000000
I_1	OUT	Output	$OUTR \leftarrow AC(L)$	I/O		000001
I_2	SKI	Skip if FGI	$FGI: PC \leftarrow PC+1$	I/O		000010
I_3	SKO	Skip if FGO	$FGO: PC \leftarrow PC+1$	I/O		000011
I_4	ION	IEN On	$IEN \leftarrow 1$	I/O		000100
I_5	IOF	IEN Off	$IEN \leftarrow 0$	I/O		000101
I_6	CLA	Clear Accumulator	$AC \leftarrow 0$	Register		000110
I_7	CLE	Clear E	$E \leftarrow 0$	Register		000111
I_8	CMA	Complement Accumulator	$AC \leftarrow \overline{AC(L)}$	Register		001000
I_9	CME	Complement E	$E \leftarrow \overline{E}$	Register		001001
I_{10}	INC	Increment Accumulator	$AC \leftarrow AC+1$	Register		001010
I_{11}	ROL	Rotate Left Accumulator	$AC \leftarrow ROL AC$	Register		001011
I_{12}	ROR	Rotate Right Accumulator	$AC \leftarrow ROR AC$	Register		001100
I_{13}	SPA	Skip if Positive Accumulator	$\overline{S}: PC \leftarrow PC+1$	Register		001101
I_{14}	SZA	Skip if Zero Accumulator	$Z: PC \leftarrow PC+1$	Register		001110
I_{15}	SZE	Skip if Zero E	$\overline{E}: PC \leftarrow PC+1$	Register		001111
I_{16}	HLT	Halt	$SC \leftarrow \text{Disable}$	Register		010000
I_{17}	AND	AND	$AC \leftarrow M[AR] \wedge AC$	Memory		100000
I_{18}	OR	OR	$AC \leftarrow M[AR] \vee AC$	Memory		100001
I_{19}	XOR	XOR	$AC \leftarrow M[AR] \oplus AC$	Memory		100010
I_{20}	ADD	Addition	$AC(L) \leftarrow M[AR] + AC$	Memory	Int/Real	100011
I_{21}	SUB	Subtraction	$AC(L) \leftarrow M[AR] - AC$	Memory	Int/Real	100100
I_{22}	MUL	Multiplication	$AC(L) \leftarrow M[AR] \times AC$	Memory	Int/Real	100101
I_{23}	DIV	Division	$AC \leftarrow AC/DR$	Memory	Int/Real	100110
I_{24}	MOD	Power	$AC \leftarrow AC \% DR$	Memory	Int	100111
I_{25}	POW	Power	$AC \leftarrow AC^{DR}$	Memory	Int/Real	101000
I_{26}	EXP	ex	$AC \leftarrow e^{-DR}$	Memory	Real	101001
I_{27}	LDA	Load Accumulator	$AC(L) \leftarrow M[AR]$	Memory		101010
I_{28}	STA	Store Accumulator	$M[AR] \leftarrow AC$	Memory		101011
I_{29}	JMP	Jump	$PC \leftarrow AR$	Memory		101100
I_{30}	BSR	Branch and Save Return-address	$M[AR] \leftarrow PC, PC \leftarrow AR$	Memory		101101
I_{31}	DSZ	Decrement and Skip if Zero	$M[AR] \leftarrow M[AR]-1$ $Z: PC \leftarrow PC+1$	Memory		101110
I_{32}	RAND	Generate a random number	$AC \leftarrow LFSR$	Specific		010001
I_{33}	CLRNOG	Clear Number of Gate	$NOG \leftarrow 0$	Specific		010010
I_{34}	INCNOG	Increment Number of Gate	$NOG \leftarrow NOG+1$	Specific		010011
I_{35}	SNOI	Set Number of Inputs	$NOI \leftarrow \text{immediate}$	Specific		010100
I_{36}	SNOTTR	Set Number of Truth Table Rows	$NOTR \leftarrow 2^{NOI}$	Specific		010101
I_{37}	GRNDC	Generate Random Circuit	$SyC[0] \leftarrow \text{Random Number}$ $SyC[NOG-1] \leftarrow \text{Random Number}$	Specific		010110
I_{38}	GNBRC	Generate a Neighbor of Circuit (Select a random gate and exchange its main control and one of its input)	$NeC \leftarrow \text{Perturbing SyC}$	Specific		010111
I_{39}	SIVDTT	Set Initial Value for Destination Truth Table	$STT[0] \leftarrow 0$. . . $STT[2^{NOG}-1] \leftarrow 2^{NOG}-1$	Specific		011000
I_{40}	MASK	Generate a Mask For Output	$MASK \leftarrow$	Specific		011001
I_{41}	CALCSTT	Calculate Synthesis Truth Table		Specific		011010
I_{42}	CALCHD	Calculate Hamming Distance Between Synthesis Truth Table and Destination Truth Table	$SHD/NHD \leftarrow \text{Hamming Distance}$	Specific		011011
I_{43}	SBC	Select Best Circuit		Specific		011100
I_{44}	SETTEMP	Set Temperature	$TEMP \leftarrow \text{immediate}$	Specific		011101
I_{45}	DECTEMP	Decrement Temperature	$TEMP \leftarrow TEMP - 1$	Specific		011110

References

- [1] K. Kucukcakar, "An ASIP design methodology for embedded systems," in Proceedings of the Seventh International Workshop on Hardware/Software Codesign (CODES'99)(IEEE Cat. No. 99TH8450), 1999: IEEE, pp. 17-21.
- [2] M. Gries and K. Keutzer, Building ASIPs: The Mescal Methodology. Springer Science & Business Media, 2006.
- [3] R. F. Mirzaee and M. Eshghi, "Design of an ASIP IDEA crypto processor," in 2011 IEEE 2nd International Conference on Networked Embedded Systems for Enterprise Applications, 2011: IEEE, pp. 1-7.
- [4] K. Shahbazi, M. Eshghi, and R. F. Mirzaee, "Design and implementation of an ASIP-based cryptography processor for AES, IDEA, and MD5," Engineering science and technology, an international journal, vol. 20, no. 4, pp. 1308-1317, 2017.
- [5] M. Venkanna, R. Rao, and P. C. Sekhar, "An Efficient Design of ASIP Using Pipelining Architecture," in International Conference on Intelligent Computing and Applications, 2019: Springer, pp. 117-128.
- [6] D. Große, X. Chen, G. W. Dueck, and R. Drechsler, "Exact SAT-based Toffoli network synthesis," in Proceedings of the 17th ACM Great Lakes symposium on VLSI, 2007, pp. 96-101.
- [7] D. Bu and P. Wang, "An improved KFDD based reversible circuit synthesis method," Integration, vol. 69, pp. 251-265, 2019.
- [8] T. Ahmed, A. Younes, and A. Elsayed, "Improving the quantum cost of reversible Boolean functions using reorder algorithm," Quantum Information Processing, vol. 17, no. 5, pp. 1-16, 2018.
- [9] A. Basak, A. Sadhu, K. Das, and K. K. Sharma, "Cost Optimization Technique for Quantum Circuits," International Journal of Theoretical Physics, vol. 58, no. 9, pp. 3158-3179, 2019.
- [10] Z. Kalantari, M. Eshghi, M. Mohammadi, and S. Jassbi, "Low-cost and compact design method for reversible sequential circuits," The Journal of Supercomputing, vol. 75, no. 11, pp. 7497-7519, 2019.
- [11] M. Lukac, M. Perkowski, and M. Pivtoraiko, "Evolutionary approach to quantum and reversible circuits synthesis," Artificial Intelligence Review Journal, vol. 20, no. 3-4, pp. 361-417, 2003.
- [12] M. Lukac, M. Pivtoraiko, A. Mishchenko, and M. Perkowski, "Automated synthesis of generalized reversible cascades using genetic algorithms," 2002.
- [13] M. Haghparast, M. Mohammadi, K. Navi, and M. Eshghi, "Optimized reversible multiplier circuit," Journal of Circuits, Systems, and Computers, vol. 18, no. 02, pp. 311-323, 2009.
- [14] M. Mohammadi and M. Eshghi, "Heuristic methods to use don't care in automated design of reversible and quantum logic circuits," Quantum Information Processing, vol. 7, no. 4, pp. 175-192, 2008.
- [15] M. Y. Abubakar and L. T. Jung, "Synthesis of Reversible Logic Using Enhanced Genetic Programming Approach," in 2018 4th International Conference on Computer and Information Sciences (ICCOINS), 2018: IEEE, pp. 1-5.
- [16] T. Atkinson, A. Karsa, J. Drake, and J. Swan, "Quantum program synthesis: Swarm algorithms and benchmarks," in European Conference on Genetic Programming, 2019: Springer, pp. 19-34.
- [17] R. Landauer, "Irreversibility and heat generation in the computing process," IBM journal of research and development, vol. 5, no. 3, pp. 183-191, 1961.
- [18] C. H. Bennett, "Logical reversibility of computation," IBM Journal of Research and Development, vol. 17, no. 6, pp. 525-532, 1973.
- [19] R. P. Feynman, "Quantum mechanical computers," Foundations of Physics, pp. 507-531, 1986.
- [20] M. P. Frank, "Introduction to reversible computing: motivation, progress, and challenges," in Proceedings of the 2nd Conference on Computing Frontiers, 2005, pp. 385-390.
- [21] E. Fredkin and T. Toffoli, "Conservative logic," International Journal of theoretical physics, vol. 21, no. 3, pp. 219-253, 1982.

Propose an E-CRM Model based on Mobile Computing Technology in Pharma Distribution Industry

Alireza Kamanghad^{1*}, Gholamreza Hashemzadeh Khorasgani², Mohammadali Afshar Kazemi¹, Nosratollah Shadnoosh¹

¹. Department of Management, Central Tehran Branch, Islamic Azad University, Tehran, Iran

². Department of Management, Central (South) Tehran Branch, Islamic Azad University, Tehran, Iran

Received: 29 May 2021/ Revised: 04 Jan 2022/ Accepted: 02 Feb 2022

Abstract

In today's world, the competition between all business areas and companies including pharma distributor companies has increased dramatically, so it is very important for active companies in the pharma distribution industry which deal with a large number of customers in a B2B market to establish a deep and long-term relationship with their customers and manage that relationship effectively. Since the CRM which now is enriched by new emerging technologies in terms of e-CRM and m-CRM is under developing rapidly, it can play a critical role for empowering these companies to strengthen their relationship with their customers. In this research it has been tried to have a complete review of mobile computing technology concept and its effect on CRM. The research methodology is basically qualitative. After a literature review, using qualitative research methods and deep interviews with a group of 8 industry experts, the whole concept of the initial model was derived using Thematic Analysis method. The Grounded Theory approach was applied to extract the main factors and sub-factors of the final model. Additionally, some of research techniques such as Dematel, ANP and Super Decision software were used to investigate the interdependency, importance and priority of factors and sub-factors. At the last stage a new model for e-CRM in pharma distribution industry based on mobile computing technology has been proposed. The four key components of the model are Quality of Content and Services, Organizational Readiness, Quality of System and Communication, Customer Mobile App.

Keywords: e-CRM, m-CRM; Customer Relationship Management; Mobile Computing Technology; Pharma Distributors.

1- Introduction

In recent years, the competitive environment of various kinds of businesses has increased. This environment which includes the pharmaceutical distribution industry has affected most of pharma distributor companies. The diversity of goods and quality of their distribution services has increased significantly the selection power of customers in a B2B market. Due to the entrance of new distribution companies in the market and intensifying the competition, a national big challenge has been emerged recently [1]. So it is very important for these companies to establish a deep and long-term relationship with their customers and manage that relationship effectively which is the main mission of a CRM solution. In this situation, CRM should be considered as one of the most important key success factors for survival in competitive market. The CRM concept has evolved in such a way to maintain a

long-term relationship with the customers. The use of CRM systems is becoming increasingly important to improve customer life time value [2] and is more important in such companies that have a wide relationship with customers and deal with a large amount of customers especially in a B2B market such as pharma distributors. On the other hand, Internet has touched almost every sphere of our lives. The impact of the process of managing and interacting customers via the internet has affected CRM too. By leveraging the internet for customer management we get the new structure of CRM known as e-CRM [3]. E-CRM is a system which tries to provide better customer services with support [4]. In recent years, many organizations have identified the need to become more customer-facing with increased global competition. Hence, e-CRM has become an essential for many organizational strategies [5]. By emerging new technologies such as smartphones and mobile computing technology, all aspects of businesses including customer relationship management have been affected rapidly and

✉ Alireza Kamanghad
ali.kamanghad.mng@iauctb.ac.ir

mobile CRM (m-CRM) concept was introduced. Mobile Customer Relationship Management (m-CRM) system is one of the recent advancements in CRM solutions. Mobile CRM promotes satisfaction to customers through the mobile medium on communication [6]. mobile computing is a kind of combination which enables a real-time connection between a mobile device and other computing environments, such as the Internet or an Intranet. This innovation is creating a revolution in the manner in which people use computers. The new computing model is basically leading to ubiquity—meaning that computing is available anywhere, at any time [7]. As we know that the ubiquity is a main part of business environment of all pharma distribution companies and mobile computing technology is the core of ubiquity, we expect that those pharma distribution companies which has empowered their CRM strategy with new mobile computing technologies would be able to increase their customer’s satisfaction and loyalty more efficient and more effective.

So we’ve summarized the main research problem as follows: The Intensifying of competitive atmosphere in pharma distribution industry, has caused big challenges both in national level and enterprise level. Considering the critical role of effective customer relationship management in pharma distributors to survive and grow in such a competitive atmosphere, how we can apply mobile computing technology capabilities to help pharma distributors for establish an effective and closer relationship with their customers.

This research aims to extract and propose a suitable model to those pharmaceutical distribution companies willing to implement e-CRM system based on mobile computing technology. The main research questions are "which is the suitable model for pharma distribution companies to implement an e-CRM system based on mobile computing technology?" And "what are the main factors and sub-factors of the proposed model and “how is their interrelationship and their relative importance?”

2- Literature Review

2-1 Customer Relationship Management (CRM)

The customer relationship management concept has been born in ancient world and has been continued during centuries anonymously until the middle of the twentieth century [8]. Interest in Customer Relationship Management began to take its importance in 1990s [9]. Around the late 1990s, the first CRM systems were introduced [5] and most of businesses nevertheless of the size are still encouraged to adopt CRM to create and manage the relationship with customers well effectively and efficiently [9]. In recent years many organizations have identified the need to become more customer-facing with increased global competition. Hence, customer

relationship management (CRM) has become an essential for many organizational strategies [5]. Companies have, therefore, invested significantly in the implementation of CRM during the years [10].

CRM, is a business approach that seeks to create, develop, and enhance relationships with carefully targeted customers in order to improve customer value and corporate profitability and thereby maximize shareholder value [11]. Recent literature explained CRM conceptualizations according to specific implementation dimensions with each dimension representing a set of business activities [10]. CRM is a comprehensive approach for creating, maintaining and expanding customer relationships [12]. CRM aims at developing sustainable, long-lasting affiliations between companies and customers [3]. A list of desired CRM benefits is collected and summarized in the table 1.

Table1. Summary of CRM benefits [13],[14]

Authors	Core CRM Benefits
Chen and Popovich (2003)	<ul style="list-style-type: none"> •Increases data sharing across selling organization •Improves customer service •Improves cross-selling/up-selling •Improves customer targeting •Enables better personalization of marketing messages •Provides better self-service options for customers •Improves buyer–seller integration
Buttle (2004)	<ul style="list-style-type: none"> •Reduces cost to serve •Increases revenue •Increases customer satisfaction and loyalty
Jones, Brown, Zoltners and Weitz (2005)	<ul style="list-style-type: none"> •Improves customization of services and product offerings •Enhances ability to create long-term partnerships •Improves salesperson efficiency and effectiveness
Stan Maklam (2005)	<ul style="list-style-type: none"> •The ability to gather customer data •Identify the most valuable customers and increase customer retention is highly •Learning from customers (customer knowledge)
Eggert, Ulaga and Schultz (2006)	<ul style="list-style-type: none"> •Improves support for product development •Increases supply-chain efficiencies via personal contact •Enhances supplier know-how
Blery & Michalakopoulos (2006)	<ul style="list-style-type: none"> •Closer relationship to its customers and offer phone services •Servicing customers and receive information to develop the level of service offered to customers
Richards, Keith & Jones (2006)	<ul style="list-style-type: none"> •Target commercial customers

	<ul style="list-style-type: none"> • Offerings from different channels • Enhanced customer service • Customized products and services
	<ul style="list-style-type: none"> • Advance responsiveness • Accelerate delivery lead-time • Enable customer knowledge management • Develop customer segmentation • Targeting the most profitable customers • Improve product and business innovations
Wang, Sedera & Tan (2009)	<ul style="list-style-type: none"> • Enhance customization of marketing efforts and messages to individual customers • Permit multi-channel integration • Allow multi-channel communication • Enable personalized products and services • Improve product separation • Focus on customers and their needs • Provide customers a “one-to-one” skill
	<ul style="list-style-type: none"> • Personalized services • Customers knowledge and experience empowered • Deliver high quality service • Meet customer needs
Popli & Rao (2009)	<ul style="list-style-type: none"> • Employee empower more time to serve up customers • Advanced satisfaction ratings • Targeted product and service contributions can be timed to match with customer actions and requirements
Keramati, Mehrabi & Mojir (2010)	<ul style="list-style-type: none"> • Individualizations of market • Customization of product and services
Schubert & Williams (2010)	<ul style="list-style-type: none"> • Improved responsiveness • Valuable time savings during reduction of the search effort • Seamless communication
Kuo, Wu & Peng (2011)	<ul style="list-style-type: none"> • Enhancing customer’s attentiveness • Consolidating helpful services
	<ul style="list-style-type: none"> • Describe diverse customer group that will be served in different ways
Kiat Loh et al. (2011)	<ul style="list-style-type: none"> • Customer service and support service operations • Predict potential and personal customer’s behavior
Amoako (2011)	<ul style="list-style-type: none"> • Improved capability to target profitable customers • Integrated contributions across channels • Individualized marketing communication
	<ul style="list-style-type: none"> • Identify and target their best customers • Allowing the formation of individualized relationships with customers
Vazifehdust et al. (2012)	<ul style="list-style-type: none"> • Identifying the most profitable customers and providing them the highest level of service • Understand and identify customer needs

Information system, technology, e-business, management, knowledge management, human resources management and marketing are the key disciplines of CRM [15]. CRM is touted as an imperative strategy to enhance a firm’s competitive advantage [16].

2-2- E-CRM

The internet is creating tremendous impact on businesses also in interacting, nurturing, maintaining their customer bases. The impact of the process of managing and interacting customers via the internet has affected CRM too. Because of the growth of the Information Technology the usage of internet began to grow up and this in turn provided opportunities to marketing through transform the way of relationships between businesses and their customers [9]. By leveraging the internet for customer management we get the new structure of CRM known as e-CRM. E-CRM is all about managing customers online using internet as the primary channel of interaction [3]. E-CRM refers to CRM using internet technology plus a database, OLAP, data warehouse, data mining, etc. [17]. More and more businesses begin to attach great importance to electronic customer relationship management (e-CRM), which focuses on customers instead of products or services, that is, considering customer’s needs in all aspects of a business, ensuring customers’ satisfaction [2].

The distinction is made between CRM and e-CRM, on the basis of three parameters of Approach, Cost and Service [3]. The Fig 1 shows the difference between CRM and e-CRM.

Table 3 –the difference between CRM and e-CRM [3]

	Approach	Cost	Service
CRM	Fragmented	High	Efficient
e-CRM	Consolidated	Low	Effective

The purpose of e-CRM is not only to bring about changes in the area of marketing, but also to improve the company's efficiency in managing customers, then to increase customer service, safeguard precious customers, and to help provide organizations with analytic capabilities [18].

2-3- Mobile Computing Technology

Mobile computing, is a computing paradigm designed for workers who travel outside the boundaries of their organizations or for any other people traveling outside their homes. As an example, salespeople were able to make proposals at customers’ offices. This, enables a real-

time connection between a mobile device and other computing environments, such as the Internet or an intranet. This innovation is creating a revolution in the manner in which people use computers. Mobile computing and commerce are spreading rapidly, replacing or supplementing wired computing. IT involves mostly wireless infrastructure and may reshape the entire IT field [7].

Nowadays, Smartphones have numerous information and communication technology functions that are comparable to those of old computers. It is estimated that the global revenues from apps will make abundant business opportunities [19]. Current smartphones and tablets contain more computing power than many of the formerly known supercomputers, which used to fill an entire room. The shift in devices has already occurred in many countries across different continents as more people are using their smartphones rather than traditional PCs. As technology is progressing to miniaturize devices, increase computing power and, especially, decrease the price of electronics, smartphone adoption will only accelerate [20]. In 1985, the Cray-2 supercomputer was the fastest machine in the world. The iPhone 4, released in June 2010, had the power equivalent to the Cray-2; now, the Apple Watch has the equivalent speed of two iPhone 4s just five years later. Nearly everyone will soon have a literal supercomputer in their pocket [21].

2-4- Mobile-Customer Relationship Management

With the development of wireless technology, mobile devices, such as smartphones and smart watches, are becoming the most effective tools for communication in human's daily life. The popularity and availability of mobile devices can help mobile users enrich experience of various services provided by mobile applications without the constrain of time and place. Mobile applications are becoming increasingly ubiquitous and can provide better user experience on mobile devices [22]. The mobile application is in rapid growth and dissemination in business and enhancement of customer satisfaction has emerged as a core issue [23].

One primordial capacity needed for implementing a CRM strategy is the ability to communicate with customers on an individual basis. For that reason, mobile technologies represent an appealing additional channel which can complement the existing channels. Among the advantages of the mobile channel which are highly relevant to CRM are the personal character of mobile devices which allows an individual customer reach, the interactivity brought by its quick message delivery and response, its reachability and ubiquity. It is the only medium enabling a spontaneous, interactive, direct and targeted interaction with customers, anytime, anywhere. This makes it a valuable channel despite the drawbacks of mobile devices.

For that reason, the future CRM solutions is envisaged to combine traditional, Internet and mobile channels [24].

m-CRM has been defined as the communication, bilateral or unilateral, that is related to marketing activities via mobile phone in order to build and maintain relationships between the consumer and the company. For this purpose, a combination of strategy, technology, and human resources is required. There are two perspectives to understand CRM in the context of new technologies. On the one hand, from the perspective of technology, m-CRM is seen as a technological tool applied to marketing in order to reduce costs and increase the efficiency of the processing information between buyer and seller. On the other hand, from the strategic perspective and relationship marketing, m-CRM is seen as a long-term management approach that companies or organizations carry out via mobile channels in order to get very different benefits. In the first perspective, the benefits of the m-CRM are the result of the application of mobile technology to the management of relationships with customers. While in the second perspective, establishing and maintaining of mutually profitable and long-lasting relationships between the company and its customers through mobile channels are benefits of m-CRM [25]. The ubiquity of mobile computing devices, such as smartphones and tablets, and the proliferation of mobile customer relationship management (m-CRM) applications, may lead to increased CRM adoption and higher returns on CRM technology investments [26].

3- Research Methodology

This research is a fundamental because of its effort for combination of two areas of knowledge including customer relation management and mobile computing technology to develop a new model and also is an applied research, since its results would be useful in making decisions and formulating policies for pharma distribution industry and those companies active in this industry. It's more an exploratory research rather than descriptive. A qualitative research approach has been selected to conduct this research. Qualitative research is a type of research that explores and provides deeper insights into real-world problems. It gathers participants' experiences, perceptions and behavior [27]. The group of experts in this research includes 8 senior managers of 7 largest pharma distribution companies of Iran.

To contact this research towards answering the main research question, a list of steps followed as follows:

- i. A complete literature review through library studies and scanning the international scientific databases from 2005 for collecting secondary data.
- ii. marking, classifying and clustering of collected data for summarizing the findings of previous researches and models.

- iii.Using in-depth interview method for collecting primary data. The group of experts were interviewed separately to answer the main research questions.
- iv.Analyze data gathered from deep interviews by using Thematic Analysis method. Then the Grounded Theory method and its coding techniques were applied for conclusion the results towards a theoretical model.
- v.Extract of experts’ opinions about weigh or importance of model’s main factors and sub-factors also their interactions with the help of paired comparison questionnaire tool.
- vi.Determine the interdependencies and priorities of main factors and sub-factors using DEMATEL technique, ANP method and Super Decision software. Decision making trial and evaluation laboratory (DEMATEL) is considered as an effective method for the identification of cause-effect chain components of a complex system [28].
- vii.Summarize the research findings and propose the final model.

4- Data Analysis

4-1- Expressing Findings

After a vast review of literature and proposed models for recognizing the main factors of e-CRM in previous researches, a list of factors has been collected in a table which is displayed in table 2.

The most important finding of this step is to realize that the main goal of any CRM system is creating or enhancing customer loyalty to the organization. This loyalty itself is rooted to two main factors: customer satisfaction and customer trust.

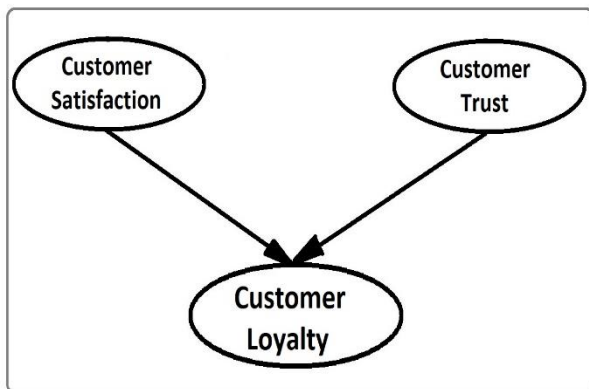


Fig. 1 Customer Loyalty is rooted in customer’s satisfaction and trust [29],[30],[31]

In the next step, several deep interviews with experts of pharma distribution industry were conducted to investigate the main factors and sub-factors of an effective e-CRM system based on mobile computing technology in this industry and extract the suitable model in this regard. To

achieve this, thematic analysis method was applied to analyze interviews and grounded theory approach was used. The group of experts were top managers of seven largest pharma distribution companies of the country. These managers were contacted through National Association of Distribution Industry and also direct relationship with their companies. An expertise criterion was applied for choosing experts. Using open coding, axial coding and selective coding techniques in grounded theory approach with a back and forth process, a complete list of categories was extracted from deep interviews. Table 3 and Table 4, show the results of open coding, Axial coding phases and categorizing the concepts.

Table 3. Open coding of research data

ID	Elements	Concepts
C1+ C2+ C3+ C4+ C7+ C12+ C15+ C23	Permanent access of customer to goods inventory of the company+ Instant access of customer to announced discounts+ Instant access of customer to announced promotions+ Permanent access of customer to his accounting information+ Financial transparency+ Access to required detailed data+ Providing useful information for customer	Required Content
C33+ C35	Customizing services for customers+ Possibility of real-time monitoring of orders	Customer’s Desired Services
C22	Possibility of systematic bilateral interaction	Systematic Bilateral Interaction
C5+ C32	Need for business process reengineering+ Rapid response of employees to the customer’s requirements	Business Process Reengineering
C18+ C20+ C21+ C28	Strong support by senior management+ Accompaniment of mangers of distribution centers+ Customer-oriented strategy+ Encourage innovation	Senior Management Support
C6+ C31	Direct communication channel of company for customers+ Stability of systematic communication	Direct Communication Channel
C26+ C27+ C34+ C39	System quality in terms of being error-free+ Simple and user-friendly system+ Flexibility of system for customers+ Regular updating of system	Qualitative Features of System
C30	Need for system security and privacy protection	System Security

C19+ C24+ C25	Staff training and culture development+ Motivating employees and incentive plan+ Recognition of system benefits by employees	Employees role
C10+ C11+ C13+ C17+ C36+ C37	Quickly obtain customer orders+ Online submitting complains and comments of customers+ Real-time scoring system for customers+ Ability to Perform system calculation+ Ability to geo-monitoring for customer shipment	Mobile Systematic features
C17	Ability to perform offline data processing by customer	Offline Processing
C8+ C14	Instant notification for customer Advanced notification for customer	Instant Notification
C9+ C38	Systematic customer satisfaction evaluation Instant satisfaction survey of customers	Real-time Satisfaction Survey
C29	Using mobile social networks	Mobile Social Networks
C40	Necessity of a customer-side mobile application software	Customer mobile app
C41+ C42+ C16	Equip the sales team with software+ Providing strong IT infrastructure+ Professional knowledge in IT department of company	Technology Infrastructure
C43+ C47	Customer loyalty as the main goal of CRM+ importance of customer satisfaction and customer trust in customer loyalty	Customer Loyalty
C44	Importance of investment by pharma distributors on new emerging technologies	Companies Investment on New Technologies
C45+ C58+ C59	Impact of macroeconomic policies on manager's decisions in pharma distributors+ Impact of business barriers on manager's decision making+ Impact of uncertainty in national economic outlook on manager's decision making	Macro economics Policies
C46	Importance of futurism in customer relationship management	Importance of Futurology
C48	Customer loyalty as the key factor for surviving in today's competitive environment	Survive in Competitive Environment

C49	Organizational culture as a key contextual factor for change	Contextual Impact of Organizational Culture
C50	Impact of macro organizational (holding) policies on manager's decisions	Holding's Policies
C51	Impact on healthy competition in pharma distribution industry	Healthy Competition in Health business
C52	Impact on general health	Improvement in General Health
C53	Collaborative activities of guild associations and NGOs on developing infrastructures	Collaborative Activities of Guild Associations
C54	The role of business alliance and coalition among pharma distribution on common investment	Business Alliances and Coalitions
C55	Impact of competitive intensity on manager's decision making	Intensifying Competitive Environment
C56	Impact of social culture on customer's welcome to new technologies	Social and Economical Context
C57	Impact of tendency to maintain current situation in pharma distributors	Tendency to Maintain Current Situation
C60	Importance of providing different resources in the company	Enterprise Resource Planning
C61	Impact of governmental health policies on manager's decision making	National Helath Macro Policies
C62+ C63	Impact of technical, communication infrastructure of the country for mobile-based services+ Impact of developing new generations of mobile communication system in the country	Technical, Network and Telecommunication Infrastructure
C64	Impact of monitoring and modeling of competitor's behavior by pharma distributors	Modeling from Competitors Behavior
C65	Impact of synergy of guild associations in pharma distribution industry	Collaborative Activities of Guild Associations
C66	Importance of research and development in distribution industry	Importance of R&D

C67+	Impact of customer loyalty on organization profitability+	Impact On Growth and Profitability
C68	Impact of customer loyalty on organizational growth	

Table 4 –Axial coding and categorizing the concepts

Categories	Concepts
Quality of Content and Services	Required Content
	Customer’s Desired Services
	Systematic Bilateral Interaction
Organizational Readiness	Business Process Reengineering
	Senior Management Support
	Technology Infrastructure
	Employees role
Quality of System and Communication	Direct Communication Channel
	Instant Notification
	Qualitative Features of System
	System Security
Mobile Social Networks	Mobile Social Networks
Customer mobile app	Mobile Systematic features
	Real-time Satisfaction Survey
	Offline Processing
Customer Loyalty	Customer Loyalty
Macroeconomics Policies	Macroeconomics Policies
Microeconomics Policies	Holding’s Policies
Health Macro Policies	National Health Macro Policies
Organizational Culture	Contextual Impact of Organizational Culture
National Technical Infrastructure	Technical, Network and Telecommunication Infrastructure
	Collaborative Activities of Guild Associations
Synergy of Associations and Forums	Synergy of Specialist Forums
	Business Alliances
Futurology	Importance of Futurology
	Importance of R&D
Growth and Profitability	Impact On Growth and Profitability
Survive in Competitive Environment	Survive in Competitive Environment
Healthy Competition	Healthy Competition in Health Business
General Health Improvement	Improvement in General Health
Intensifying Competitive Environment	Intensifying Competitive Environment
Social and Economical Context	Social and Economical Context
Conservatism	Tendency to Maintain Current Situation
Enterprise Resources Development	Enterprise Resource Planning
Investment on Emerging	Companies Investment on

Technologies	New Technologies
Modeling from Competitors	Modeling From Competitors Behavior

Following the grounded theory method, in addition to recognition of Customer Loyalty as the “core category”, other categories including “casual conditions”, “strategies”, “context conditions”, “intervening conditions” and “consequences” were recognized. These categories and their interrelationship is illustrated in a diagram named axial coding paradigm which is shown fig 2.

With the selective coding technique of grounded theory and reviewing findings by experts group as we can see on Fig 2, the final list of main factors for increasing customer satisfaction and customer trust and consequently customer loyalty, were recognized as follows:

- Quality of Content and Services
- Organizational Readiness
- Quality of System and Communication
- Customer mobile app

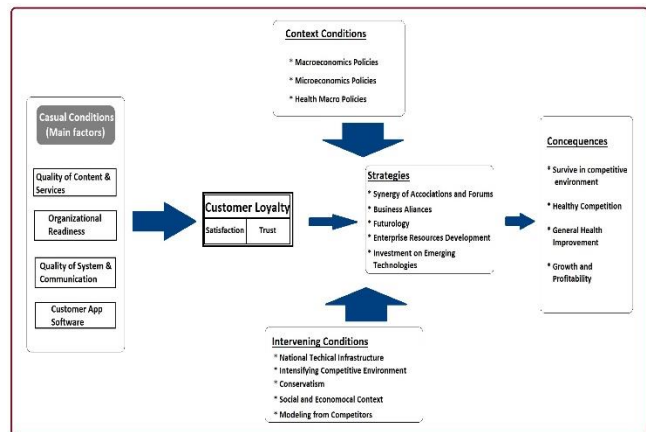


Fig. 2 Axial coding paradigm of research data

In addition, the final list of sub-factors categorized in Table 5.

Table 5 –Final list of main factors and sub-factors

Main Factors	Sub-Factors
Quality of Content and Services	Required Content
	Customer’s Desired Services
	Systematic Bilateral Interaction
Organizational Readiness	Business Process Reengineering
	Senior Management Support

	Technology Infrastructure
	Employees role
Quality of System and Communication	Direct Communication Channel
	Instant Notification
	Qualitative Features of System
	System Security
Customer Mobile App	Mobile Systematic features
	Real-time Satisfaction Survey
	Offline Processing

4-2- Interdependencies Between Factors

The content of interview process with experts and next talking to them, showed that there are interrelationships and interdependences between main factors and sub-factors. To evaluate these interdependencies, the Dematel method was applied using paired comparison questionnaire. The results of filled questionnaire by group of expert in Dematel calculations format for main factors are summarized in Table 6.

Table 6. R and J values of main factors

Main Factors	R	J	R+J	R-J
Quality of Content & Services	0.431	1.112	1.543	-0.681
Organizational Readiness	1.385	0.000	1.385	1.385
Quality of System & Communication	0.578	0.648	1.226	-0.07
Customer mobile app	0.537	1.172	1.709	-0.635

Table 6, shows that “organizational readiness” factor has the most impact on other factors and “customer mobile app” has the most impact from other factors. Furthermore, the “customer mobile app” factor has the most interaction with other factors whiles the intensity of “organizational readiness” on other factors is obviously more than others.

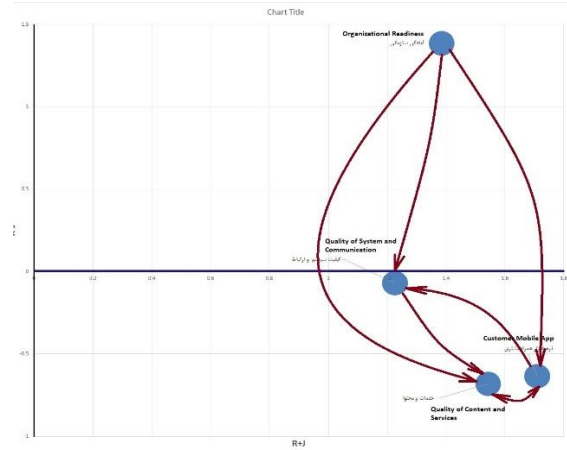


Fig. 3 Cartesian graph of Interdependencies between main factors

Fig 3 is the Cartesian Graph of interdependent relations between main factors based on Dematel method.

As the same way, the results of Dematel calculations for sub-factors are summarized in Table 7.

Table 7. R and J values of sub-factors

Sub-Factors	R	J	R+J	R-J	Affecting / Affected
Required Content	0	1.04	1.04	-1.04	Affected
Customer’s Desired Services	0.4	0.6	1	-0.2	Affected
Systematic Bilateral Interaction	1.24	0	1.24	1.24	Affecting
Business Process Reengineering	0.146	1.511	1.657	-1.365	Affected
Senior Management Support	1.364	0.134	1.498	1.23	Affecting
Technology Infrastructure	0.458	0.942	1.4	-0.484	Affected
Employees role	0.959	0.34	1.299	0.619	Affecting
Direct Communication Channel	0.4	0.571	0.971	-0.171	Affected
Instant Notification	0.92	0	0.92	-0.92	Affected
Qualitative Features of System	1.342	0.429	1.771	0.913	Affecting
System Security	0.607	0.429	1.036	0.178	Affecting

Mobile Systematic Features	1	0	1	1	Affecting
Real-time Satisfaction Survey	0	0.429	0.429	-0.429	Affected
Offline Processing	0	0.571	0.571	-0.571	Affected

Fig 4 illustrates the Cartesian Graph of interdependent relations between sub factors based on Dematel method.

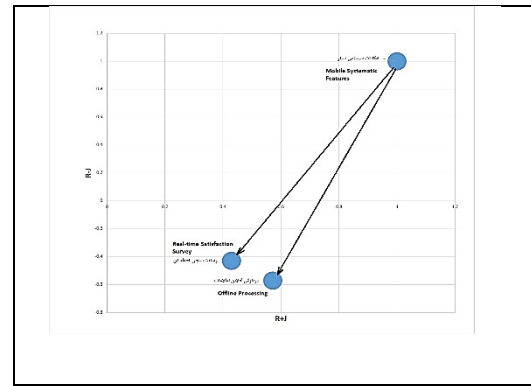
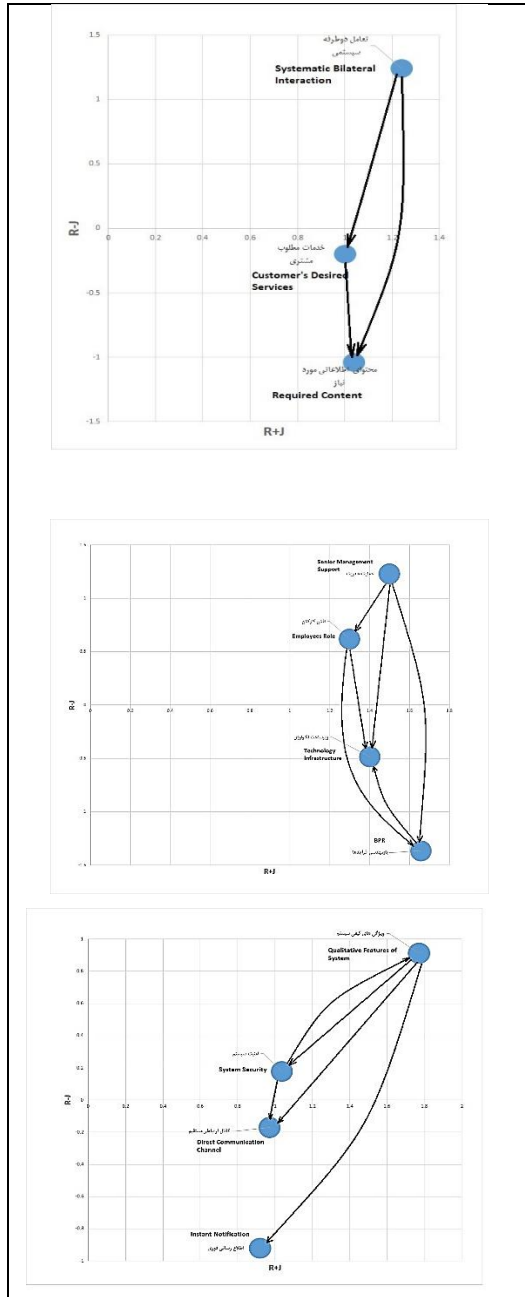


Fig. 4 C. graph of Interdependencies between sub-factors



4-3- Prioritization and Ranking of Factors

After clarifying interdependencies among main and sub-factors, the Analytic Network Process (ANP) was applied to evaluate the priorities of the all factors and to rank them based on their importance in final model. The Analytic Network Process (ANP) is a generalization of the Analytic Hierarchy Process (AHP). Priorities are established in the same way they are in the AHP using pairwise comparisons and judgment. Many problems cannot be structured hierarchically because they involve the interaction and dependence of higher-level elements in a hierarchy on lower-level elements [20]. So the ANP method was used considering the inner and outer interaction and dependence among main factors and sub-factors.

Fig 5 shows the feedback network map of relations between all factors.

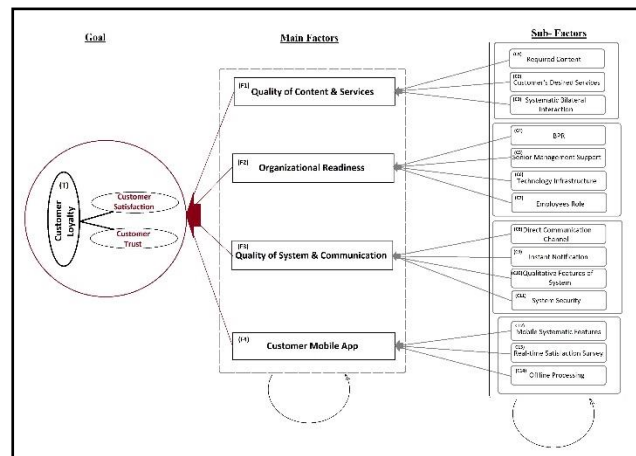


Fig. 5 Network map of relations among factors for ANP analysis

The preliminary data gathered from paired comparison questionnaire firstly was averaged in Excel software then using Super Decision software, the Weighted Super-matrix of all factors was obtained as shown on fig 6.

C43	C42	C41	C34	C33	C32	C31	C24	C23	C22	C21	C13	C12	C11	T-Loyalty	F4	F3	F2	F1		
0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.23524	0.161690	0.33330	0.16667	0.0000	F1
0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.11551	0.293810	0.23333	0.00000	0.11817	F2
0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.05647	0.04449	0.0000	0.166670	0.04097	F3
0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.59277	0.0000	0.233330	0.166670	0.34086	F4
0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	T-Loyalty
0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.5	0.666670	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.35443	C11
0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.5	0.0000	0.5	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.08931	C12
0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.33333	0.5	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.05626	C13
0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.66601	0.0549	0.33333	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.3252	0.0000	0.0000	C21
0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.26312	0.0000	C22
0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.333330	0.08234	0.0000	0.0000	0.0000	0.0000	0.0000	0.17181	0.0000	C23
0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.289740	0.333330	0.31503	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0252	0.0000	0.0000	C24
0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	C31
0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	C32
0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	C33
0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	C34
0.875	0.66667	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	C41
0.125	0.0000	0.14286	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.3333	0.0000	0.0000	0.0000	C42
0.0000	0.333330	0.25714	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.23333	0.0000	0.0000	C43

Fig. 6 Weighted Super-matrix of factors

The final results of Super Decision calculations based on ANP analysis are summarized in table 8.

Table 8- Ranking of Main Factors and Sub-Factors

Main Factors	Sub-Factors Code	Sub-Factors Weight	Sub-Factors Ranking	Factors Weight	Factors Ranking
Quality of Content & Services	C11	0.0855	4	0.2176	3
	C12	0.0769	5		
	C13	0.0684	6		
Organizational Readiness	C21	0.0369	11	0.3106	2
	C22	0.0908	3		
	C23	0.0514	9		
	C24	0.0568	7		
Quality of System & Communication	C31	0.0127	13	0.1500	4
	C32	0.0095	14		
	C33	0.0413	10		
	C34	0.0329	12		
Customer Mobile App	C41	0.1981	1	0.3219	1
	C42	0.0517	8		
	C11	0.01871	2		

The above table shows that on the basis of experts' opinion, the "Customer Mobile app" is the most important factor for an effective e-CRM system based on mobile computing technology in pharma distribution industry. Also the importance of "organizational readiness" factor is very close to first factor.

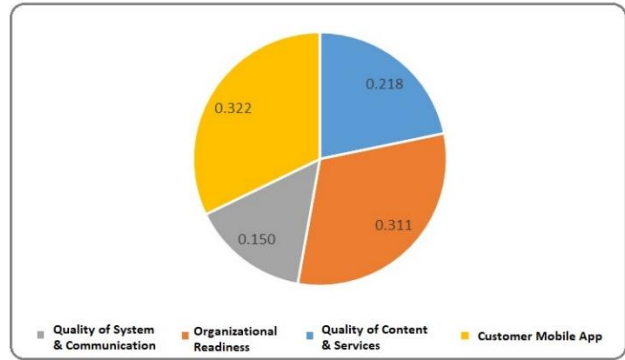


Fig. 7 Weight of main factors in a pie chart

The importance of main factors is illustrated in a pie chart (fig 7) and the importance of sub-factors with their ranking is illustrated in a bar diagram (fig 8).

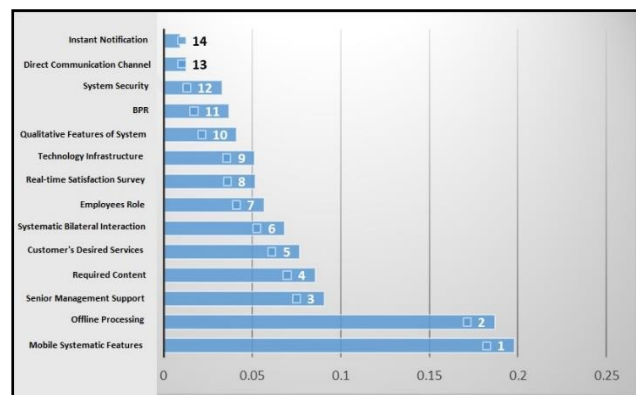


Fig. 8 Weight of sub-factors and their ranking

Fig 8 shows that "mobile systematic features" and "offline processing" are the most important sub-factors based on experts' opinion.

4-4- The Final Model

With extract of all main factors and their importance during research process, the final model for an effective e-CRM system based on mobile computing technology in pharma distribution industry has been designed and proposed as fig 9.

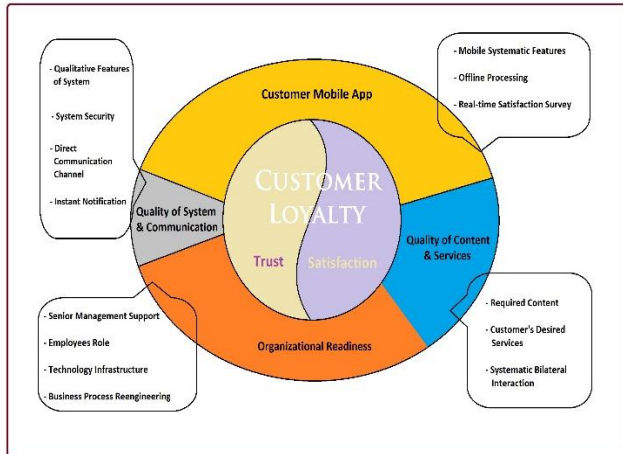


Fig. 9 The final model of research

In the process of proposing the final model, these fundamentals have been considered: Firstly, the positive impacts of implementing e-CRM system on different aspects of organization has been discussed in previous researches. In these researches the organization word refers to a general concept and includes all kind of economic firms in all industries. Whereas the pharma distribution industry in Iran has its own special characteristics such as direct B2B market, large number of customers, wide geographical dispersion, moving sales team, type of customers and special kind of products and also the competition environment has been intensified dramatically between companies in this industry, there is real need for a new technological e-CRM model dedicated to this industry. Secondly, almost in all previous researches we can find the customer loyalty is introduced as the main goal of all CRMs. This was confirmed too by the group of experts in this research which means that an effective new e-CRM model should increase customer loyalty towards the organization. Thirdly, the customer mobile app has discovered as a key factor in this research whilst has not been mentioned in none of previous researches. This is a main factor with highest importance in proposed model which seems is added to this area's body of knowledge and extracted from experts' opinion.

Effective e-CRM system based on mobile computing technology in pharma distribution industry was recognized as follows:

- Quality of Content and Services
- Organizational Readiness
- Quality of System and Communication
- Customer mobile app

5- Conclusion

Reviewing the research process and its findings, shows that the customer loyalty which is rooted to customer's satisfaction and trust is the main goal of customer

relationship management. To achieve this, a combination of two areas of knowledge including e-CRM and mobile computing technology can provide a significant capacity. In these days, for pharma distributors which have special organizational characteristics and are doing business under a high competitive pressure, it is very important to acquire this capacity. A suitable model for acquisition this capacity is necessary. These companies can apply the proposed model in this research which is suitable for implementing an e-CRM system combined with mobile computing technology. This model has four key components. Although they have different weights and importance in the model but all should be considered simultaneously by the companies. For any key component there are several sub-factors to clarify dimensions of related key component and act as guidelines. These four key components are briefly summarized as follows: **a.** A mobile application software that is compatible with mobile devices (Smartphones, Tablets, ...) and is affordable to customers with ability to run on customer's smartphone. This application is equipped with offline processing capabilities and provides accessibility to required information for customers wherever and whenever they like. **b.** The type and quality of content, information and services that propose and present to customers via the mentioned application software. **c.** The level of organizational readiness of the pharma distributor company in terms of management, IT infrastructure, processes, resources, etc. for support of services that company wants to give its customers via the mentioned application software. **d.** Qualitative dimensions of whole company's system including the mentioned software application and the quality of its connection to other core systems of the company. This should include the quality of permanent connectivity and high performance of the application software which will be used by customer everywhere and every time.

Research findings indicate that designing and developing a customer-side mobile application has the highest priority for an effective m-CRM system in pharma distribution Industry. Not only this factor but also organizational readiness of the company and other both factors should be considered and invested enough by companies. Accordingly, we suggest the proposed model to those pharma distribution companies willing to have an effective customer relationship management with their customers. We suggest these companies to revise their customer-oriented strategies and establish or modify their CRM approach with mobile computing technology with the help of proposed model. We suggest them to emphasize on developing a customer-side mobile application software coincide with investment on improving level of e-readiness. Continuous improvement of updated content, information, services and facilities that are delivered to

customers and likewise system quality modification are recommended strongly at the next steps.

Additionally, we suggest all companies and business enterprises which are faced with a rising up competitive environment to utilize the findings and results of this research and its proposed model for adoption a modern approach in their CRM structure. The model can be a general guideline for all companies willing to establish an effective e-CRM or m-CRM comprehensive system.

Implementing an e-CRM based on mobile computing technology is an enterprise multi-dimension project that its most important prerequisite is company's e-readiness. Considering the importance of this factor and its effect on the other factors, we suggest all companies interested in mobile-CRM implementation to perform an e-readiness assessment project and evaluate their level of readiness before starting the main project.

Mobile computing technology is expanding rapidly both technologically and applicably. We tried to investigate the role of this technology in customer relationship management specially in pharma distribution. There are many opportunities for further studies to investigate the role of mobile computing technology in other business areas and other industries. The research on applying the proposed model of this research on other type of service-oriented companies is another subject for further studies.

References

- [1] <https://www.tasnimnews.com/fa/news/1394/07/20/886841>
- [2] A. Mishra, D. Mishra, "Customer Relationship Management: Implementation Process Perspective", *Acta Polytechnica Hungarica*, 6 (4), 83- 99, 2009.
- [3] A. Jafarnejad, C. Loox and A. Monshi, "Towards Electronic Customer Relationship Management: An CRM Solutions Development Mythology", *Iranian Journal of Management Studies (IJMS)*, 1(1), 73-89, 2007.
- [4] S. Jirehbandei and A. Nemmaney Pour, "A New Model for e-CRM in e-Commerce using Live-Operator", *International Journal of Social, Behavioral, Educational, Economic, Business and Industrial Engineering*, Vol 5, 8, 1006-1008, 2011.
- [5] P. Fottouhiyehpour, "Assessing the Readiness for implementing CRM in B2B Markets Using AHP Method", (Master Thesis). Lulea University of Technology, Sweden, 2008.
- [6] D. Verma and D. Singh Verma, "Managing Customer Relationship through Mobile CRM in Organized retail outlets", *International Journal of Engineering Trends and Technology (IJETT)*, Vol 4(5), 1697-1701, 2013.
- [7] Turban et al, "Information Technology for Management: Transforming Organizations in the Digital Economy", 6th ed. John Wiley & Sons, 2008.
- [8] A. Motameni and E. Jafari, "The role of Human resource in CRM implementation", *The Managers Covenant*, 49, 51-56, 2010.
- [9] M. Karunanithy and K. Kajendra, "An Evolution of Customer Relationship Management: A Conceptual Approach", *Proceeding of Jaffna University International Research Conference (JUICE 2012)*, published March 2014, Sri Lanka, 49-54, 2014.
- [10] I. Dalla Pozza, O. Goetz, Oliver and J. Sahut, "Implementation effects in the relationship between CRM and its performance", *Journal of Business Research*, 89, 391-403, 2018.
- [11] I. Grazdane, "A Customer Relationship Management Approach for Optical Retail Business", (Master Thesis), Helsinki Metropolitan University of Applied Sciences, 2013.
- [12] K. Anderson, and C. Kerr, "Customer Relationship Management", McGraw Hill Companies, Inc., 2006
- [13] K.A. Richards and E. Jones, "Customer relationship management: Finding value drivers", *Industrial Marketing Management*, 37,120–130, 2008.
- [14] N. Mohammadhossein, and Z. Nor Hedayati, "CRM Benefits for Customers: Literature Review", *International Journal of Engineering Research and Applications (IJERA)*, Vol 2, 6, 1578-1586, 2012.
- [15] M. Aloka, M. Alkhateeb, M. Abbad, and F. Jaber, "Customer Relationship Management: A review and Classification", *Transactional Marketing Journal*, 7, 2, 187-210, 2019.
- [16] R. Lin, R. Chen, and K. Shun Chiu, "Customer relationship management and innovation capability: an empirical study", *Industrial Management & Data Systems*, 110(1), 111-133, 2018.
- [17] M. Wang, "Measuring CRM service quality in the library context: a preliminary study", *The Electronic Library*, 26(6), 896-911, 2007.
- [18] F. Nuradlin, R. Ferdiana, R. and S. Fauziati, "Current Trend and Literature on Electronic CRM Adoption", *International Conference on Electronic Representation and Algorithm, ICERA 2019*.
- [19] F. Tseng, T.T.L. Phan, T.C.E. Cheng and C. Teng, "Enhancing customer loyalty to mobile instant messaging: Perspectives of network effect and self-determination theories", *Telematics and Informatics*, 35, 1133-1143, 2018.
- [20] Word Economic Forum, "Deep Shift Technology Tipping Points and Societal Impact", Retrieved from http://www3.weforum.org/docs/WEF_GAC15_Technological_Tipping_Points_report_2015.pdf, 2015.
- [21] K. Schwab, "The Fourth Industrial Revolution", Penguin Random House, UK, 2017.
- [22] F. Gu, J. Niu, Z. Qi and M. Atiquzzaman, "Partitioning and Offloading in Smart Devices for Mobile Cloud Computing: State of the Art and Future Directions", *Journal of Network and Computer Applications*, Accepted 19 June 2018.
- [23] Z. Gani and W. Maung, "Mobile Applications in customer relationship management to enhance empowerment of knowledge to customers", *Journal of Physics: Conference Series*, 1502 0120235, 2020.
- [24] G. Componovo, Y. Pigneur, A. Rangone, and F. Renga, "Mobile Customer Relationship management: An Explorative Investigation of the Italian Consumer Market", 42- 48. 10.1109/ICMB, 2005.
- [25] S. San-Martin, N.H. Jimenez and B. Lopez-Catalan, "The firm's benefits of mobile CRM from relationship marketing approach and the TOE model", *Spanish Journal of Marketing- ESIC 20*, 18-29, 2016.

- [26] M. Rodriguez and K. Trainor, "A conceptual model of the drivers and outcomes of mobile CRM application adoption", *Journal of Research in Interactive Marketing*, Vol 10 (1), 67-84, 2016.
- [27] S. Tenny, G.D. brannan, J.M. Brannan & N.C. Sharts-Hopko, "Qualitative Study", In *Statpearls*, Statpearls Publishing, 2021.
- [28] S. Li, X. You, H. Liu and P. Zhang, "DEMATEL Technique: A Systematic Review of the state-of-the-Art Literature on Methodologies and Applications", *Mathematical Problem in Engineering*, 101-134, 2018.
- [29] W. Lee and LS. Wong, "Determinants of Mobile Commerce Customer Loyalty in Malaysia", *Procedia- Social and Behavioral Sciences*, 224: 60-67, 2016.
- [30] MH. Ronaghi, A. Dehdarizadeh, S. Safaee and A. Asadpour, "An eCRM model for banking industry in Iran", *New Marketing Research Journal*, Special Issue, 1-12, May 2012.
- [31] R. Thakur, "Understanding Customer Engagement and Loyalty: A Case of Mobile Devices for Shopping", *Journal of retailing and Consumer Services*, 32: 151-163, 2017.

A Novel Approach for Establishing Connectivity in Partitioned Mobile Sensor Networks using Beamforming Techniques

Abbas Mirzaei^{1*}, Shahram Zandiyan¹

¹. Department of Computer Engineering, Ardabil Branch, Islamic Azad University, Ardabil

Received: 12 Jun 2021/ Revised: 04 Jan 2022/ Accepted: 23 Feb 2022

Abstract

Network connectivity is one of the major design issues in the context of mobile sensor networks. Due to diverse communication patterns, some nodes lying in high-traffic zones may consume more energy and eventually die out resulting in network partitioning. This phenomenon may deprive a large number of alive nodes of sending their important time critical data to the sink. The application of data caching in mobile sensor networks is exponentially increasing as a high-speed data storage layer. This paper presents a deep learning-based beamforming approach to find the optimal transmission strategies for cache-enabled backhaul networks. In the proposed scheme, the sensor nodes in isolated partitions work together to form a directional beam which significantly increases their overall communication range to reach out a distant relay node connected to the main part of the network. The proposed methodology of cooperative beamforming-based partition connectivity works efficiently if an isolated cluster gets partitioned with a favorably large number of nodes. We also present a new cross-layer method for link cost that makes a balance between the energy used by the relay. By directly adding the accessible auxiliary nodes to the set of routing links, the algorithm chooses paths which provide maximum dynamic beamforming usage for the intermediate nodes. The proposed approach is then evaluated through simulation results. The simulation results show that the proposed mechanism achieves up to 30% energy consumption reduction through beamforming as partition healing in addition to guarantee user throughput.

Keywords: Mobile Sensor Networks (MSNs); Connectivity Restoration; Network Partitioning; Cooperative Beamforming; Fault Recovery.

1- Introduction

Mobile sensor networks (MSN) are effective platforms for the industrial and military communications which are applied for commercial affairs of the manufacturing sector and the identification of enemy frontiers in the military. In most applications, sensors have the role of a data source and send information from event triggers to each eNodeB or central receiver. The unique role of the base station makes it a natural target for the enemies who intend to carry out the deadliest attack with the least possible effort against the mobile sensor network. Even if the mobile sensor network uses common security mechanisms such as encryption and authentication, the enemy may use traffic analysis techniques to identify the base station. However, the attractiveness of mobile sensor networks and their advantages make them vulnerable to potential attack by

the evil enemy. A typical mobile sensor network consists of several relays that iteratively transfer new information to the existing BS. In this model, because the unique role of the base station makes it possible to carry out the most effective attack against the target mobile sensor network with the least possible effort, this station becomes the center of enemy attacks. That is, the enemy assumes that a Denial of Service (DoS) attack against the base station will actually cripple the larger mobile sensor network, because the base station not only acts as a data well.

One of the main effective ways to protect a base station from a vicious enemy attack is to keep its role, identity, and location unknown. However, conventional security mechanisms that provide confidentiality, integrity, and authentication are not capable of this type of protection [1] [2]. One of the major portions of the studies relevant to anonymous communications were so far related to analyzing routing algorithms with the aim of concealing actual paths from the transmitter well [3], [4]. It should be

✉ Abbas Mirzaei
a.mirzaei@iauardabil.ac.ir

noted that, in spite of the fact that secure routing algorithms can greatly reduce path discovery attack, the enemy can gain important data via monitoring the link layer and the relevance between pairs of nodes, based on which it can identify the location and role of the base station [5] [6].

According to [7] [8], the authors suggested an approach in the lower layer which uses dynamic beamforming to further identify the base station. Nowadays, distributed beamforming seems a very attractive way to improve the network performance, throughput and power utility, provide data link safety, in addition to increasing SINR in the multi-layer cooperative systems [9] [10]. Based on dynamic beam modulation, several mobile sensor network nodes work together to share the existing propagation capabilities in order to create a dynamic multiple transmission network. Various relays are able to concurrently transmit information, taking into account the conditions of the wireless channel and the precise control of the signal phase, in such a way that all the signals are combined at the destination. For example, ideally, N transmitters send the same messages with the same power, while tolerating a path loss during the transmission of the signal to a normal destination increases the power at the destination by N times. This feature has been shown to increase base station anonymity in mobile sensor networks. This protocol appropriately disrupts the evidence hypothesis (EH) and it doesn't consider the real base station of the mobile sensor network well.

This protocol is an effective technique for enhancing the probability of low-cost, multi-hop paths usage, and the power needed to transmit the signal to the destination is used as the cost of the L-link. Because, the average energy consumption cost of the protocol increases with increasing the number of auxiliary relays $|L|$, the use of L link selecting the paths to maximize $|L|$ can increase the base station anonymity with energy costs equal to the mobile sensor networks with anonymous protection [11].

As far as we know, participatory communication was first used to reinforce the base station anonymity in Ref. [12]. As a result, the former studies of distributed beam formation and base station anonymity will be discussed separately. Researchers on the subject of base station anonymity initially defined a quantitative way of measuring anonymity. Some researchers developed sub-optimal effective approaches for measuring anonymity in the connection entropy [13], GSAT test [14], and belief [15] [16]. Entropy and GSAT methods impose certain limitations on the enemy. They give the a priori possibility that the location of the base station is known to the enemy or that the enemy can estimate the location of the base station. The functionality of the belief index, according to the evidence theory, does not have any of these hypotheses and so it has attracted a lot of attention as a metric for recognizing anonymity. In Section 4, we discuss the

evidence theory and the metric of belief to evaluate base station security. Many published techniques for dealing with the traffic insecurity in mobile sensor networks applied various approaches to make the location of data sources hidden [17] [18]. Such as [19] in which the authors propose various approaches such as uniform packet speed and false paths to confuse the enemy. Similarly, the authors in [20] suggest that network paths be modified by considering virtual sinks. Two techniques have been proposed in [21]. In the first technique, the base station re-transmits a package of received packets at various degrees and the base station looks like an ordinary node for the enemy. The base station can also be considered and can move to a safer location.

The above techniques are used in the network layer of communication protocols. The protocol uses distributed beamforming in the physical layer to improve the base station anonymity [22]. This paper compensates for this shortcoming of the previous protocol by providing a multi-layered routing algorithm considering data link constraints and auxiliary intermediate nodes in order to decrease the total power utilization.

The authors in [23] proposed an efficient smart control plan for the dynamic transmission in the wireless sensor networks using cooperative protocols. This algorithm supports the dynamic operations of block data and third-party public validation to provide high security against data forgery and replacement. In [24], a QoS model for resource allocation algorithm was proposed for data replicas based on the servers existing in a network in order to improve the connectivity approach service and decrease total cost. In [25], a distributed algorithm was proposed to reduce the access delay and expand the network bandwidth. In this scheme, a new data analysis strategy was proposed to mitigate the costs of data storage and information transfer for applications. In [26], a close-loop content-oriented scheme was reviewed achieving higher performance for data-intensive applications. Also, some researchers addressed main critical challenges of this criterion, such as energy efficiency [27] availability [28], and security [29] of data access. However, heterogeneous MSN has security challenges, including vulnerability for sensors and association acknowledgment, that delay the rapid adoption of computing models.

Unfortunately, the abovementioned works cannot be considered as a proper approach for large-sized networks due to reliability conditions and high computational complexity at the central unit that significantly increases to the number of sensors in the network.

In this paper, we present the cost of the L link, which has been optimized for multi-hop paths that minimize the average power consumption of the mobile sensor network. Using simulations, we show that the cost of our link is such that it maintains the anonymity of the base station while reducing the communication energy consumption.

This article continues as follows: Section 2 examines the system model and the problem formulation. Section 3 provides a framework for Distributed beamforming and the energy efficiency of the protocol. Section 4 describes the proposed approach in Power Optimization in mobile sensor networks. The numerical results and discussion of the proposed approach are presented in Section 5. Finally, Section 6 draws the conclusion and highlights future challenges to motivate the effective integration of beamforming-based mobile sensor networks with the diagnosis.

2- System Model and Problem Formulation

A- Network Model

In this paper, we consider a homogeneous model for a mobile sensor network in which all sensor nodes have the same capabilities in terms of battery life, type of radio communication, and network protocols. In this paper the sensors were considered as mobile nodes. The base station acts as a well for all data traffic generated by the sensor nodes. There is only one base station on the network. Our hypothesis is that the sensors can be aware about the locations of the base station and the neighbor sensors as well [30]. In addition, the cells are well-informed about the level of transmission energy needed to get all subsequent hops. Multi hop paths are followed to deliver data frames to the base station. Also, the cross-wave propagation model has been considered in this paper.

We assume that precautions are used in the design and operation of the base station to prevent enemy infiltration. For example, the base station maintains the transmission power level equal to the other cells (for example, updated path exploration and authentication messages) so that it cannot be detected from other sensor nodes by radio frequency analysis. Messages are transmitted with the header and encrypted message body. We assume that the TDMA Media Access Control Protocol (MAC) operates by synchronizing sufficient time on all wireless network sensors at tolerable shielding intervals [31]. All nodes in mobile sensor networks are considered as auxiliary relay options.

B- Problem Formulation

Assume that the mobile sensor network transmits the target sensitive data, which is a desirable target for the enemy. After identifying and abusing the base station, the enemy aims to carry out a DoS attack against the base station at any cost, such as physically destroying the base station. Also, the enemy is actively engaged in eavesdropping by being present in all parts of the mobile sensor network [32] [33]. The enemy is able to identify the location of all radio communications at the location of the network [34]. While the enemy monitors the traffic, we assume that the cryptographic system is robust enough so that the enemy cannot use the cryptographic system

analysis to retrieve the contents of the body or header. The enemy uses the evidence theory traffic analysis to localize the base station, unaware that the mobile sensor network is using the distributed beamforming.

The enemy starts via monitoring the transmit links demonstrated by $E(U)$ in which U is a direct connection among the nodes (S_i & D_i). It also obtains the paths by correlating all the evidence for the node pair. The overall path containing two or more nodes is denoted by V , and the associated evidence $E(V)$ is calculated as follows:

$$E(V) = \min_{U \in V} \{E(U)\}, \quad |V| \geq 2, \quad (1)$$

Normalized evidence $m(V) = E(V)/\Sigma E(V)$ shows a proportion of all the evidence gathered by the enemy that supports the $B(u)$ which indicates the enemy's certainty that there is a path of length n in any given node and is expressed as follows:

$$B(u) = \sum_{U|u \in V} n m(U). \quad (2)$$

In this paper, the Belief index has been applied in order to evaluate base station anonymity. The small belief metric means less confidence of the enemy or more anonymity of the base station coordination. To reduce the computational complexity of the calculations needed, it's supposed that the enemy splits the mobile sensor network into an $M \times M$ network consisting of N_c square cells. This means that the enemy only needs to identify the target location within the cell. As a result, the belief metrics $B(u)$ generated by cell analysis cause u to indicate that the sensor is not a specific sensor, but one of the N_c cells in the enemy's target network. Section 5.b has presented an example of evidence theory analysis. Figure 1 illustrates the network configuration and the communication beamforming links between the sensors and the base station.

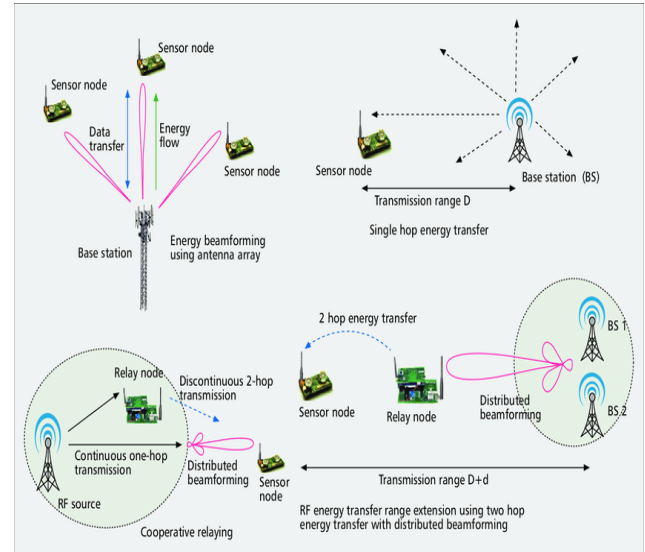


Fig. 1. The network configuration and the beamforming between sensors and the base station.

3- Distributed Beamforming Model

A. Distributed Beamforming

Based on the proposed methodology, a distributed beamforming approach has been applied to further identify the base station. Distributed beamforming uses the broadcast nature of wireless transmission. Adjacent nodes may also hear all the frames sent to a particular receiver. According to Figure 2, these adjacent nodes may act as auxiliary relays in cooperation with the transmission source S_i so that the transmitted signal travels to D_i through a diverse set of transmitters. Because each R_j relay sends the same message S_i with the exact time and synchronization of the carrier, the signals at the destination D_i are combined under the conditions of ideal scheduling and carrier synchronization [35].

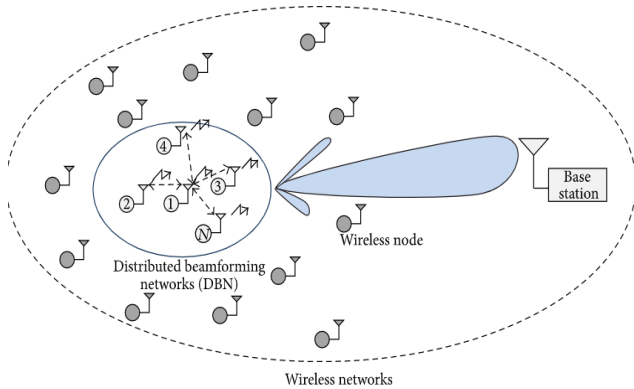


Fig. 1. Distributed beamforming model

Ideal received signal in D_i by source S_i and $|L|$ auxiliary relay R_j is sent as shown below:

$$r_{D_i}(t) \triangleq r_{S_i D_i}(t) + \sum_{j=0}^{|L|-1} r_{R_j D_i}(t)$$

$$= \Re(A_{S_i}(t) w_{S_i D_i} \beta h_{S_i D_i} e^{j(2\pi f_c t + \theta(t) + \varphi(t))}) + \Re\left(\sum_{j=0}^{|L|-1} A_{R_j}(t) w_{S_i D_i} \beta h_{S_i D_i} e^{j(2\pi f_c t + \theta(t) + \varphi(t))}\right) + n(t) \quad (3)$$

In this equation, h indicates the channel shock response, β indicates the sharing efficiency f_c is the carrier frequency, $\theta(t)$ illustrates the phase modulation expression, $\varphi(t)$ is total phase variation term and $n(t)$ represents the thermal noise contained in the D_i receiver. $r_{D_i}(t)$ consists of two expressions: Information received from source S_i (i.e., $r_{S_i D_i}(t)$) and total information received from the relay $|L|$ used $R_j \in L$ is equal to $\sum_{j=0}^{|L|-1} r_{R_j D_i}(t)$.

The meaning of Equation 3 in terms of achieving base station anonymity in the physical layer is that the

distribution of the distributed beamforming makes it possible for the component of the signal received $r_{S_i D_i}(t)$ from the transmitted information S_i decreases by $\sum_{j=0}^{|L|-1} r_{R_j D_i}(t)$ at the same time, the power level of the signal received by D_i remains constant during the phase offset $\varphi(t)$. Therefore, if S_i and $R_j \in L$ transmit only data at a specified SINR with the power required to reach D_i , each transmitter can reduce the resource via $10 \log(|L| + 1)$ dB and properly prevents the enemy from distinguishing $E(S_i, D_i)$. Accordingly, the elimination of $E(S_i, D_i)$ from the enemy evidence set increases the confidentiality of the base station, as it reduces its role in detecting $B(u = D_i)$. Figure 2 shows the functionality of the distributed beamforming used in any relay via the node S_i , which intends to apply dynamic beamforming to increase the base station anonymity while sending a packet to the next relay node D_i . S_i should first choose the appropriate subset of auxiliary nodes by handshaking (according to steps (a) & (b)). When a vector of auxiliary relays is used, S_i transfers the packet body in step d to $R_j \in L$. In step e , the distributed beamforming is submitted and then the authentication message is transmitted from node D_i to S_i in stage (f) to confirm the correct reception of the cooperative message.

B. Distributed Beamforming Protocol

Figure 3 is an example of an enemy theory (ET) analysis with and without using distributed beamforming for seven network sensors. In this example, the enemy divides the target area into $N_c = 9$ cells. In Figure 3, the transmission of information in the main method along with the cooperative transmission by the distributed beamforming. In continuation, the paper shows the evidence collected and the belief calculation in unit hops in the main method and the distributed beamforming protocol, in which the relay in cluster 28 sends the packet to the relay in cluster 35.

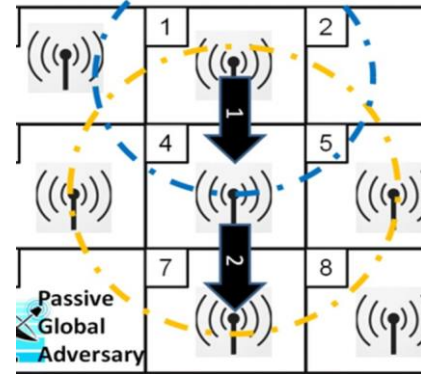


Fig. 3. Applying enemy-theory for a mobile sensor network

C. Distributed Beamforming Energy Analysis

In addition to the higher anonymity of the base station, the energy cost of communications resulting from the use of the distributed beamforming protocol should also be

considered. While it may save up to $10\text{Log}(|L| + 1)$ on the transmission power, additional signaling in the distributed beamforming protocol causes increase of the data overhead. Dedicated power for each sent bit ε_b can be computed by ratio of the mean transmit energy $\overline{P_T^{S_i, D_i}}$ needed to reach the S_i to D_i information by the signal to noise ratio (SNR) and the set speed r , where $\overline{\varepsilon_b^{S_i, D_i}} \triangleq \frac{\overline{P_T^{S_i, D_i}}}{r}$. Therefore, the mean communication energy consumed by the non-cooperative system is equal to $\overline{\varepsilon_{base}} \triangleq \overline{\varepsilon_b^{S_i, D_i} \bar{\beta}}$, where $\bar{\beta}$ is equal to the average body size. Accordingly, the mean distributed beamforming power is calculated based on a cooperative hop and is obtained based on the following formula:

$$\overline{\varepsilon_{DIBAN}} \triangleq \left(\varepsilon_b^{\overline{S_i R_j}} (\gamma_{RR} + \gamma_{data} + \bar{\beta} + K) + \frac{\varepsilon_b^{\overline{S_i D_i}}}{(|L| + 1)} \bar{\beta} + |L| \varepsilon_b^{\overline{R_j S_i}} \gamma_{Ack} + |L| \left(\frac{\varepsilon_b^{\overline{R_j D_i}}}{(|L| + 1)} \bar{\beta} + \varepsilon_b^{\overline{D_i S_i}} \gamma_{Ack} \right) \right), \quad (4)$$

Where Y_{RR} represents the Relay Request, Y_{Data} denotes the Data Multihop and Y_{Ack} is equal to the header size of acknowledgement/negative acknowledgement.

In this scenario, the channel was considered to be equal to K bits at time t microseconds (i.e., $K = t \times r$) [7].

In equation (4), we have two key observations of the average the distributed beamforming energy consumption per hop ($\overline{\varepsilon_{DIBAN}}$). First, S_i does not make saving on the distributed beamforming for power transmission ($\overline{P_T^{S_i, R_j}}$) which involves the use of relays. As a result, S_i uses relays that are as close as possible to S_i , thus minimizing ($\varepsilon_b^{\overline{S_i, R_j}}$). Second, S_i seeks to maximize $|L|$ by applying the maximum number of auxiliary relays possible and at the same time minimizing the use of transmission power ($\overline{P_T^{S_i, R_j}}$).

D- Selection of Relay in the Distributed Beamforming Protocol

The distributed beamforming protocol requires an approach to select a relay in order to use a set of $R_j \in L$ in each hop and has three objectives: higher anonymity of the base station, conservation or reduction of communication energy $\overline{\varepsilon_{DIBAN}}$ compared to the main system $\overline{\varepsilon_{base}}$ without distributed beamforming and achieving the best CSI measurement. Increased $|L|$ causes the reduction of the ability to communicate directly with S_i and D_i and thus improves the base station anonymity. But as we said before in the previous section, the increase of $|L|$ requires higher transmission power ($\overline{P_T^{S_i, R_j}}$) to be used during relay operation. By iteration, we obtain the expected number of potential relays that S_i made available by increasing the power level ($\overline{P_T^{S_i, R_j}}$). Given the constraint $\overline{P_T^{S_i, R_j}} < \overline{P_T^{S_i, D_i}}$, we maintain the anonymity of the base station to prevent

the enemy from collecting the evidence of $E(S_i, D_i)$ linking the transmitter and receiver.

First, the quantity of the potential intermediate nodes was considered equal to $|L_D| = \lambda \left(\frac{\pi d^2_{S_i, R_j}}{8} \right)$ where λ is the density of the node in the region, which is $\lambda = \frac{S_U}{M \times M}$ in a semicircle with radius $d_{S_i, R_j} = \frac{d_{S_i, D_i}}{\delta}$ is calculated based on the number of expected nodes in the receiving interval when S_i is transmitted with power $\overline{P_T^{S_i, D_i}}$. This algorithm adjusts $\overline{P_T^{S_i, R_j}}$ using δ through iteration where $\delta > 1$, because $\delta = 1$ means that $\overline{P_T^{S_i, R_j}} = \overline{P_T^{S_i, D_i}}$ and there is an undesirable link between S_i and D_i .

Our relay selection algorithm is briefly performed based on the below stages:

1- S_i chooses the primary amount of δ and mathematically calculate $|L_D|$. In the beginning, the set of the distributed beamforming relays is $L = \emptyset$.

2- S_i sends the relay request message and $\overline{P_T^{S_i, R_j}}$ is calculated based on the δ to reach $|L_D|$.

3- Nodes that respond to S_i with a confirmation message involves the number of available relays $|L_A|$ in a way that $R_j \in L_A$.

4. One of the following three results occurs for $R_j \in L_A$:

A- If $|L_A| < |L_D|$, then $\delta = \delta - \delta_{STEP}$ is used to increase $\overline{P_T^{S_i, R_j}}$ and reach more candidate relays in each iteration.

Return to step 2. If $\delta = \min(\delta)$, the algorithm ends with $L = \emptyset$ and the distributed beamforming is not used in this hop.

B- If $|L_A| = |L_D|$, then $L = L_A$ and the algorithm terminates.

C- If $|L_A| \geq |L_D|$, L becomes $|L_D|$ of the relay of the highest quality $R_j \in L_A$ which are prioritized based on the best condition of channel state information.

The proposed node selection approach is used at all relays next to the transmitter path to the base station. Decreasing δ is the only node solution to use more intermediate nodes

to enhance $\overline{P_T^{S_i, R_j}}$. Nevertheless, the multi-layer routing approach offers another option to increase $|L|$ for the mobile sensor network: In this case, the paths are selected as a functionality of the accessible auxiliary intermediate nodes $|L_A|$ in all hops.

The cost of a L_i link consists of two parameters, each with a unique objective for the distributed beamforming. First, energy saving is the most important factor in the initial design of mobile sensor networks and is one of the fundamental technical design constraints for further anonymity of the base station. Instantaneous energy required to use a set of the auxiliary relay $|L_A|$ is equal to:

$$\varepsilon_{RR} \triangleq \sum_n \left(\left(\varepsilon_b^{S_i R_j} \times \gamma_{RR} \right) + \gamma_{Ack} \sum_{j=0}^{|L_A|-1} \varepsilon_b^{R_j S_i} \right), \quad (5)$$

The reader is reminded that while selecting a relay, the energy to bit $\varepsilon_b^{S_i R_j}$ required to use $|L_A|$ the relay is dependent on δ . quantity of re-transfers needed for successful use of a suitable set of $|L_A|$ of the relay, a relay in which the relay selection condition is met is represented by n . Secondly, the average energy consumption of the distributed beamforming decreases with increasing the number of $|L_A|$ of the auxiliary relays available in each hop. As a result, the cost of the beam links is calculated based on the following formula:

$$\mathcal{L}_i = \begin{cases} \left(\frac{\varepsilon_{RR}}{|L_A|} \right) \text{ for } |L_A| > 0 \\ \infty \text{ for } |L_A| = 0 \end{cases} \quad (6)$$

When the auxiliary relays are available, the cost of link \mathcal{L}_i is decimal and when the relays are not available it becomes ∞ . The cross-layer scheme integrates with the distributed beamforming relay selection algorithm so that each node can calculate $|L_A|$ and make better use of the distributed beamforming by selecting routes with higher relay densities.

4- Power Optimization in Mobile Sensor Network

According to the main sources in this field, there is no standard model for energy consumption in beamforming-based mobile sensor backhaul networks. However, the application of nonlinear prediction energy consumption in such systems has attracted more satisfaction. Here, this paper uses adaptive resource allocation in which the backhaul connection has been modeled, in which C5 and C6 are the maximum transmission power constraints for sensors and macro base stations, respectively.

A. Content-Caching Model

In this network, we suppose that content can be modelled as a distinct set of packet data as $\mathcal{F} = \{\mathcal{F}_1, \mathcal{F}_2, \dots, \mathcal{F}_f, \dots, \mathcal{F}_F\}$ which \mathcal{F}_f represents the f -th data frame. The request probability for data frame f is expressed as

$$p_f (0 \leq p_f \leq 1), \quad \text{which, } \sum_{f=1}^F p_f \leq L_i, \quad \forall f \in \mathcal{F}, \quad (7)$$

It should be noted that the caching model presented for this paper is stochastic caching so that we can calculate the probability of caching data packet f via base station i $0 \leq q_{f_i} \leq 1$ where L_i illustrates the cache size. In addition, $\{q_{f_i}\}$ of base station i should satisfy below condition:

$$\sum_{f=1}^F q_{f_i} \leq L_i, \quad \forall i \in \mathcal{B}, f \in \mathcal{F}, \quad (8)$$

$L_i = L_M$ means the cell is macro cell, otherwise, $L_i = L_S$.

B. Resource Control Model

Based on the approach's principles, the resources of the base stations can be supplied by conventional smart grid and renewable energy harvesting. During this scenario, the transmission power relevant to base station i can be demonstrated by $P_i (i \in \mathcal{B})$, and the applied energy from the grid network is illustrated by G_i . The harvested renewable resource is shown as E_i . Based on the enabled power sharing capability, the shared power among cell i and cell i' is equal to $\varepsilon_{ii'}$, where $\beta \in [0,1]$ denotes the power-sharing index among base stations. So, we can conclude that $(1 - \beta)$ is equal to the loss percentage in the power sharing stage. The following condition should be satisfied during the power sharing process.

$$P_i < G_i + E_i + \beta \sum_{i' \in \mathcal{B}, i' \neq i} \varepsilon_{ii'} - \sum_{i' \in \mathcal{B}, i' \neq i} \varepsilon_{i'i}. \quad (9)$$

According to the defined conditions, the overall power efficiency can be affected by the transmission strategies, power sharing and the level of harvested energy from renewable resources.

C. Transmission Model

Taking fairness into account, we tried to provide data rate balancing throughout the network. In which, $x_{ij} (i \in \mathcal{B}, j \in \mathcal{U})$ denotes the association indicator, for example, $x_{ij} = 1$ represents that node j is associated with BS i and otherwise the node has not been associated with the base station. Subsequently, $k_i = \sum_{j \in \mathcal{U}} x_{ij}$ represents the number of sensors associated to cell i . $(\sum_{f=1}^F p_f q_{f_i})^{k_i}$ express the probability of serving k_i associated nodes by base station i . if $x_{ij} = 1$, we can calculate the efficiency of the j -th sensors as $\mu_{ij} = \log(R_{ij})$ which R_{ij} is the throughput so that the R_{ij} is obtainable as.

$$R_{ij} = \left(\sum_{f=1}^F p_f q_{f_i} \right)^{k_i} \frac{\mathcal{B}\beta}{\sum_{j \in \mathcal{U}} x_{ij}} \log(1 + \gamma_{ij}) \quad (8)$$

In this framework, the ratio of signal to interference-noise can be computed via (11)

$$\gamma_{ij} = \frac{P_i h_{ij}}{\sum_{i' \in \mathcal{B}, i' \neq i} P_{i'} h_{i'j} + \sigma^2} \quad (11)$$

In this formulation, h_{ij} and $h_{i'j}$ indicate the main channel gain and the interfering channel gain respectively, B denotes the frequency bandwidth. σ^2 is also a noise figure. We can consider the goal function equivalent to

minimization of the applied grid power. Consequently, we have the goal function as the following.

$$\mathbf{P1:} \max_{q,x,P,\varepsilon,G} \sum_{i \in \mathcal{B}} \sum_{j \in \mathcal{U}} x_{ij} \mu_{ij} - \eta \sum_{i \in \mathcal{B}} G_i \quad (12)$$

$$s. t. \quad C1: \sum_{i \in \mathcal{B}} x_{ij} \gamma_{ij} \geq \gamma_{min}, \forall j \in \mathcal{U},$$

$$C2: \sum_{i \in \mathcal{B}} x_{jm} = 1, \forall j \in \mathcal{U},$$

$$C3: P_i < G_i + \beta \sum_{i' \in \mathcal{B}, i' \neq i} \varepsilon_{i'i} - \sum_{i' \in \mathcal{B}, i' \neq i} \varepsilon_{ii'} + E_i, \forall i \in \mathcal{B},$$

$$C4: \sum_{f=1}^F q_f i \leq L_i, \forall i \in \mathcal{B}, f \in \mathcal{F},$$

$$C5: 0 \leq q_f \leq 1, \forall f \in \mathcal{F}, \forall i \in \mathcal{B},$$

$$C6: x_{ij} \in \{0,1\}, \forall i, j \in \mathcal{U},$$

$$C7: G_i \geq 0, \varepsilon_{i'i} \geq 0, \forall i \in \mathcal{B},$$

$$C8: 0 \leq P_i \leq P_{max}^i, \forall i \in \mathcal{B},$$

In which, $\mathbf{q}=[q_f]$, $\mathbf{X}=[x_{ij}]$, $\mathbf{P}=[P_i]$, $\boldsymbol{\varepsilon}=[\varepsilon_{i'i}]$, $\mathbf{G}=[G_i]$, γ_{min} illustrates the $SINR_{min}$ to guarantee the reliability of the connection between nodes and the base station. Also, η represents a weighting factor for evaluation of the power efficiency index. The multi hop strategy of backhauling in the presented mobile sensor network and the configuration of the network has been exhibited in Figure 4.

Figure 4 shows a simple network structure to illustrate the operation process of the proposed connectivity approach.

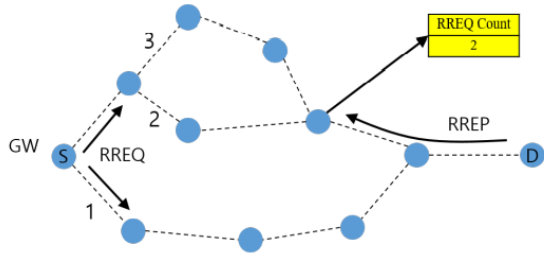


Fig. 4. Operation process of the proposed connectivity approach

In order to improve the reliability of the proposed model, we applied a path repair process with the branch node-based routing algorithm which is shown in Figure 8.

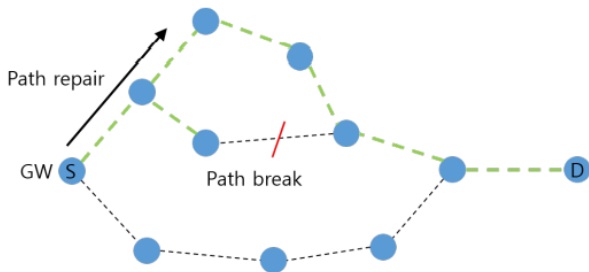


Fig. 5. Reliable distributed transmission model.

In this paper, we present the cost of the L link, which has been optimized for multi-hop paths that minimize the average power consumption of the mobile sensor network. Our distributed cross layer routing protocol uses a connection cost that can be added to the final goal function. Using simulations, the paper shows that the cost of our link is such that it maintains the anonymity of the base station while reducing the communication energy consumption.

The problem of association and power allocation in this approach may be modelled as problem P2 which itself can be considered as the optimal solution for the primary problem P1. Taking $k_i = \sum_{j \in \mathcal{U}} x_{ij}$, into account, this problem is expressed as the following.

$$\mathbf{P2:} \max_{q,x,P,\varepsilon,G} \sum_{i \in \mathcal{B}} \sum_{j \in \mathcal{U}} x_{ij} \log(c_{ij}) + \sum_{i \in \mathcal{B}} k_i^2 \log \left(\sum_{f=1}^F p_f q_f i \right) - \sum_{i \in \mathcal{B}} k_i \log(k_i) - \eta \sum_{i \in \mathcal{B}} G_i \quad (13)$$

$$s. t. \quad C1, C2, C3, C4, C5, C6, C7, C8,$$

$$C9: \sum_{j \in \mathcal{U}} x_{ij} = k_i, \forall i,$$

$$\text{where, } c_{ij} = B \log(1 + \gamma_{ij}).$$

D. Data Caching-based User Association Algorithm

In this framework, P2 as a NL mixed integer programming problem is not a convex problem and as Lemma 1 indicated, the sub gradient method will be the best approach to solve it. Taking $\{P, \varepsilon, G\}$ into account, the sensor association problem will be mathematically modeled as follows.

$$\mathbf{P2.1:} \max_{q,x} \sum_{i \in \mathcal{B}} \sum_{j \in \mathcal{U}} x_{ij} \log(c_{ij}) + \sum_{i \in \mathcal{B}} k_i^2 \log \left(\sum_{f=1}^F p_f q_f i \right) - \sum_{i \in \mathcal{B}} k_i \log(k_i) \quad (14)$$

$$s. t. \quad C1, C2, C4, C5, C6, C9$$

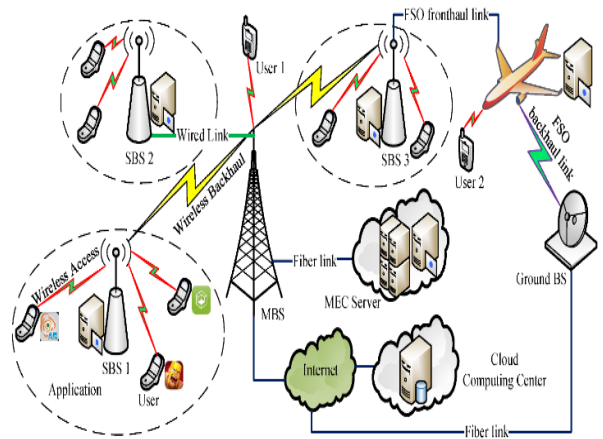


Fig. 4. Beamforming model in the mobile sensor architecture

Lemma 1: considering $p_{(1)} \geq \dots \geq p_{(f)} \geq \dots \geq p_{(F)}$ as the probability of demanded payload (f), the optimal solution for P2.1 is achievable as the following.

$$q_{f_i}^* = \begin{cases} 1, & f_i = (1), \dots, (L_i) \\ 0, & \text{otherwise} \end{cases}, \quad \forall i \in \mathcal{B}. \quad (15)$$

Proof 1: as mentioned before, P2.1 shows that obtaining the optimal value for $\sum_{f=1}^F p_f q_f$ is the target in which, the demanded payload is itself consists of L_i sections $\mathcal{F}_l (l = 1, \dots, L_i)$, and the probability \mathcal{F}_l is more than \mathcal{F}_{l+1} . Therefore, we have:

$$\sum_{(f) \in \mathcal{F}_l} q_{(f)_i}^l = 1, \sum_{l=1}^{L_i} q_{(f)_i}^l = q_{(f)_i} \text{ and } \bigcup_{l=1}^{L_i} \mathcal{F}_l = \mathcal{F}$$

Also,

$$\begin{aligned} \sum_{f=1}^F p_f q_{f_i} &= \sum_{l=1}^{L_i} \sum_{(f) \in \mathcal{F}_l} p_{(f)} q_{(f)_i}^l \leq \sum_{l=1}^{L_i} p_{(l)} \left(\sum_{(f) \in \mathcal{F}_l} q_{(f)_i}^l \right) \\ &\Rightarrow \sum_{f=1}^F p_f q_{f_i} \leq \sum_{l=1}^{L_i} p_{(l)}, \end{aligned}$$

Consequently, according to (15), this theory is confirmed. So, we can conclude that

$$\begin{aligned} \tilde{\text{P2.1:}} \quad & \max_x \sum_{i \in \mathcal{B}} \sum_{j \in \mathcal{U}} x_{ij} \log(c_{ij}) + \sum_{i \in \mathcal{B}} k_i^2 \log \left(\sum_{f=1}^{L_i} p_{(f)} \right) \\ & - \sum_{i \in \mathcal{B}} k_i \log(k_i) \end{aligned} \quad (16)$$

s. t. C1, C2, C6, C9.

For simplicity of process to obtain the best solution for $\tilde{\text{P2.1}}$ as a combination of several sub-problems, we work on its dual problem. Therefore, the target function should be reformulated as follows:

$$\begin{aligned} \mathcal{L}(x, k, \mu, \nu) &= \sum_{i \in \mathcal{B}} \sum_{j \in \mathcal{U}} x_{ij} \log(c_{ij}) + \sum_{i \in \mathcal{B}} k_i^2 \log \left(\sum_{f=1}^{L_i} p_{(f)} \right) \\ & - \sum_{i \in \mathcal{B}} k_i \log(k_i) - \sum_{j \in \mathcal{U}} \mu_j \left(\gamma_{\min} - \sum_{i \in \mathcal{B}} x_{ij} \gamma_{uj} \right) - \\ & \sum_{i \in \mathcal{B}} \nu_i \left(\sum_{j \in \mathcal{U}} x_{ij} - k_i \right), \end{aligned} \quad (17)$$

Where in this formulation, $\nu = [\nu_i]$, $k = [k_i]$ and $\mu = [\mu_j]$. It should be noted that ν_i and μ_j represent Lagrangian multipliers. In continue, we can define the problem's dual function $\mathcal{D}(\cdot)$ as the following

$$\mathcal{D}(\mu, \nu) = \begin{cases} \max_{x, k} \mathcal{L}(x, k, \mu, \nu) \\ \text{s. t. C2, C6.} \end{cases} \quad (18)$$

Subsequently, the dual problem of $\tilde{\text{P2.1}}$ (16) will be formulated as

$$\min_{\mu \geq 0, \nu \geq 0} \mathcal{D}(\mu, \nu). \quad (19)$$

μ_j and ν_i are coefficients of the dual problem and solution of the goal function can be obtained as the following steps

$$x_{ij}^* = \begin{cases} 1, & \text{if } i = i^* \\ 0, & \text{otherwise} \end{cases}, \quad (20)$$

In (20), $i^* = \arg \max_i (\log(c_{ij}) + \mu_j \gamma_{ij} - \nu_i)$.

Considering k_i , the second derivation of the goal function results in

$$\frac{\partial^2 \mathcal{L}}{\partial k_i^2} = 2 \log \left(\sum_{f=1}^{L_i} p_{(f)} \right) - \frac{1}{k_i}. \quad (21)$$

$$k_i^* = - \frac{W \left(-2 \log \left(\sum_{f=1}^{L_i} p_{(f)} \right) e^{\nu_i - 1} \right)}{2 \log \left(\sum_{f=1}^{L_i} p_{(f)} \right)}, \quad (22)$$

As it is obvious, $\sum_{f=1}^{L_i} p_{(f)} \leq 1$, so, $\frac{\partial^2 \mathcal{L}}{\partial k_i^2}$ cannot be a positive amount. With setting $\frac{\partial^2 \mathcal{L}}{\partial k_i^2}$ equal to zero, k_i^* is achieved as the optimum degree of k_i .

In (22), $W(z)$ shows the Lambert-W factor as a response for $z = we^w$. According to (20), the optimal solution (μ^*, ν^*) cannot be achieved by differentializing of $\mathcal{D}(\mu, \nu)$. Therefore, applying the iterative gradient approach will be useful.

$$\mu_j(t+1) = \left[\mu_j(t) - \delta(t) \left(\sum_{i \in \mathcal{B}} x_{ij}(t) \gamma_{ij} - \gamma_{\min} \right) \right]^+, \quad (23)$$

$$\nu_i(t+1) = \left[\nu_i(t) - \delta(t) \left(k_i(t) - \sum_{j \in \mathcal{U}} x_{ij}(t) \right) \right]^+, \quad (24)$$

In this formulation, $x_{ij}(t)$ and $k_i(t)$ can be renewed in an iteration via (20) and (22). The step size was shown by $\delta(t)$ and we have $[a]^+ = \max\{a, 0\}$, t also indicates the iterations quantity.

5- Numerical Results

In this section, we present the simulation results that demonstrate the effectiveness of the proposed Multi hop Cooperative Beamforming Mobile Sensor Network (MCB-MSN) approach. For simplicity, we assume that the harvested energy by base station during each time interval is constant. Following the former schemes, we modeled the energy harvest at each base station as the stationary stochastic process. In addition, we assume that the popularity of the content follows the introduced distribution model of [36] and that the contents of the F library have been sorted by popularity. Thus, the probability of the f^{th} demand for popular content is calculated based on [37], where α indicates the skewness of popularity. We compare the performance of the association design and our proposed power optimization based on the signal strength received from the Reference Signal Received Power (RSRP). In the simulation, the sensors randomly move in the macrocellular

geographical area and the main simulation parameters have been shown in Table 1.

We studied the power optimization in beamforming-based multi-layer heterogeneous networks. An association algorithm and cooperative power control were proposed to find the optimal data speed in addition to decreasing the network total power utilization. These schemes consider the system security, data speed and energy consumption to be relatively important. Also, in this paper, the effect of the number of sensors and the size of the cache was also investigated.

A- Simulation Environment

We evaluate the effectiveness of our approach using proprietary computer simulation with the Monte Carlo method. We analyze the results of implementing the proposed approach in MATLAB 2019 and CVX tool of Python programming language. The experiment environment was considered as a multi-layer heterogeneous system with a number of small cells within a micro layer.

The confidence interval of the results is 90%. The S_U sensor nodes are evenly distributed on the grid at $880 \times 880 m^2$. The base station is fixed without any mobility. But activated cells are able to submit packets to the base station in cell 35 via multihop paths. This procedure divides the target network into a grid of 36 cells measuring $167 \times 167 m^2$. We considered the sensitivity of each network node to receive a signal equal to -100 dBm, which is the common value in mobile sensor networks and the signal-to-noise (SNR) required a function which are able to estimate the base CSI and $P_T^{S_i, D_i}$ required to reach the next hop destination by observing the signal to noise (SNR). The maximum transmission power of each node is limited to 30 dBm. Distributed beamforming is used for each hop and occurs each time when $|L_A| > 1$.

We measure the performance of the protocol considering the context of base station anonymity (i.e., reducing the belief that the cell contains base station $B(u = 35)$). To assess MCB-MSN from the energy efficiency point of view, we compared it with three other schemes of mobile sensor networks: Fixed Power Allocation (FPA), Random Power Allocation (RPA) and Cooperative NOMA Simultaneous Wireless Information and Power Transfer (CN-SWIPT) [38].

B- Numerical Results

Figure 5 compares the average throughput of the proposed approach, MCB-MSN and CN-SWIPT scheme with equal maximum transmission power. Based on this figure, it is obvious that the average sum data rate increases with increasing the signal to noise ratio. It can also be seen that the MCB-MSN algorithm performs much better than other algorithms in terms of higher data rate. Because the MCB-MSN algorithm has the required flexibility to dedicate resources to the network entities.

Table 1. Main implementation factors

Parameter	Value
Configuration of the Network	Mobile Network, X-sectored BSs
sensor distribution model	uniform (U) and hotspot (Hs),
transmit backoff	1.5 dB
Base Station MTP	43 dBm
Codec strategy	Adaptive multi-rate
Rx loss & Tx loss	3 dB
Propagation model	Okumura-Hata
Fairness Index	Security/ Throughput
Upper bound of iteration	2000
L_{margin}	5 dBm
Learning factor $c_1 = c_2$	1.1
Weighting factor ω_{max} MAX	0.77
Weighting factor ω_{min} MIN	0.28
MAX Sensor Power P_m^o	90 w
MIN Sensor Power P_s^o	5 w
Sensor transmit range	30 m

This trend decreases slightly with an increase in N, because the algorithm reduces the throughput available to each of the nodes. In contrast, the demand for the throughput of each node is the same in all random power allocation, equal power allocation, CN-SWIPT algorithms, because they all provide the minimum throughput for each N. The throughput decreases exponentially with increasing number of N. Because with increasing N, the demand for data rate decreases. Because, the same MTP is shared equally between the nodes.

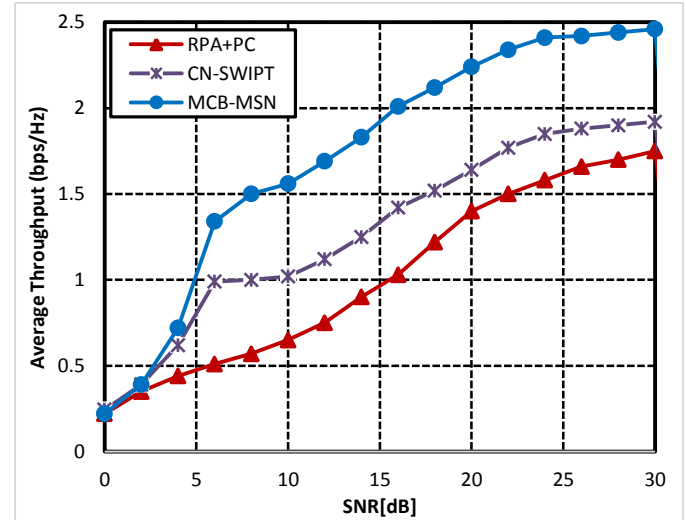


Fig. 5. Average throughput vs. signal to noise ratio

Based on the achieved results in Figure 6, the average sum-rate increases almost linearly with increasing N in the MCB-MSN algorithm. While all three other two algorithms, CN-SWIPT and RPA/PC show slight improvement

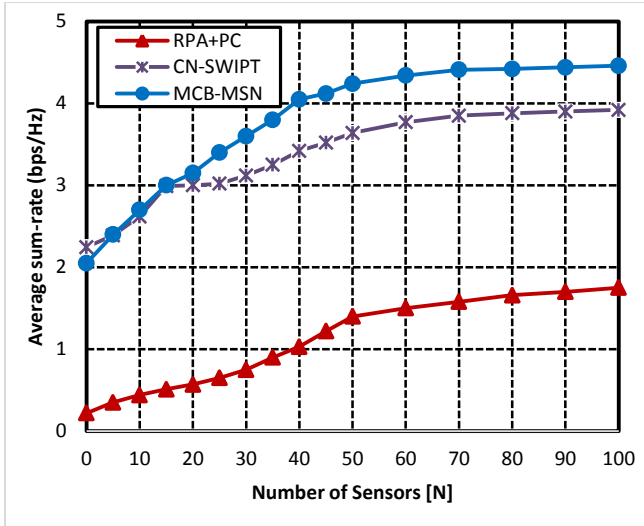


Fig. 6. Average sum rate vs. number of sensors

Figure 7 shows the capacity of backhaul links and their average traffic (link usage in percentage) for random power allocation, MSB-MSN and CN-SWIPT. Based on this figure, it can be seen that the MCB-MSN algorithm has the best performance in terms of load balancing and link usage and capacity. So, it has the highest possible efficiency in using backhaul links. Also, the high capacity of backhaul links reduces the potential for the backhaul link to be trapped in the bottleneck while sending the traffic flow to the central network.

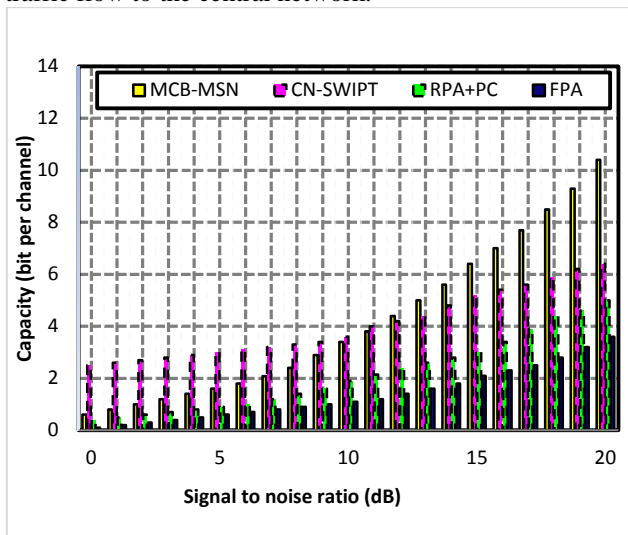


Fig. 7. Backhaul links capacity vs. signal to noise ratio

Figure 8 shows the performance of MCB-MSN and CN-SWIPT algorithms according to different numbers of nodes. As can be seen from the graph, energy efficiency is obtained when broader constraints on the total number of nodes are considered. Because, the feasible range of the problem

increases and the algorithm has more freedom to maximize the throughput and minimize the energy consumption. But as the upper bound of demand decreases, the feasible range becomes narrower and energy efficiency decreases. Further lowering the upper bound to $y_{min} = y_{max} = C_{ue}$ (where c_{ue} is the user equipment demand that must be met) results in identical efficiency of both the MCB-MSN and CN-SWIPT algorithms. Therefore, the MCB-MSN algorithm, which uses dynamic power optimization, performs better than the CN-SWIPT, which has strict constraints procedure.

Figure 9 shows that the MCB-MSN algorithm uses power sources better than other algorithms. Increasing the MTP to saturation increases the energy efficiency. After reaching saturation, increasing MTP does not affect the energy efficiency. Increasing MTP in Equivalent Power Allocation (EPA), Random Allocation (RA) and CN-SWIPT algorithms does not improve energy efficiency. In this case, the performance of these algorithms is slightly worse. In the form of 1000 independent simulations, we investigated how many times each algorithm successfully calculates the solution. Our criterion is actually the possible values used to summarize the result of each simulation result in CVX. CVX as a linear programming method is a powerful optimizer for solving iterative problems like the introduced main problem. Such convex-based tools can also be applied to analyze for rapid prototyping of models and algorithms incorporating convex optimization.

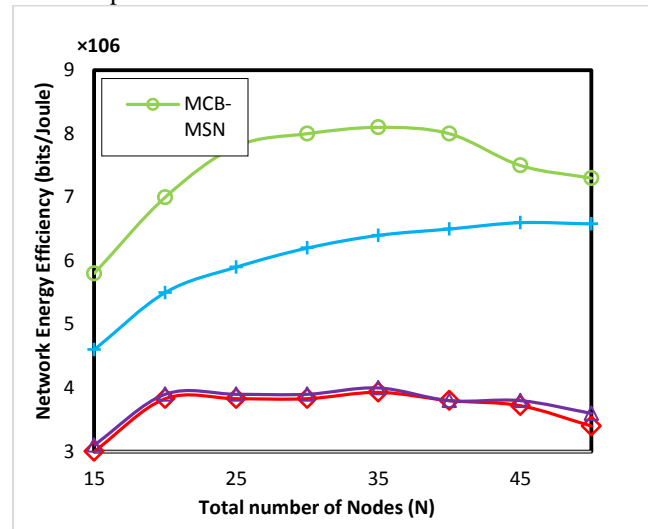


Fig. 8. Network energy efficiency vs. total number of nodes

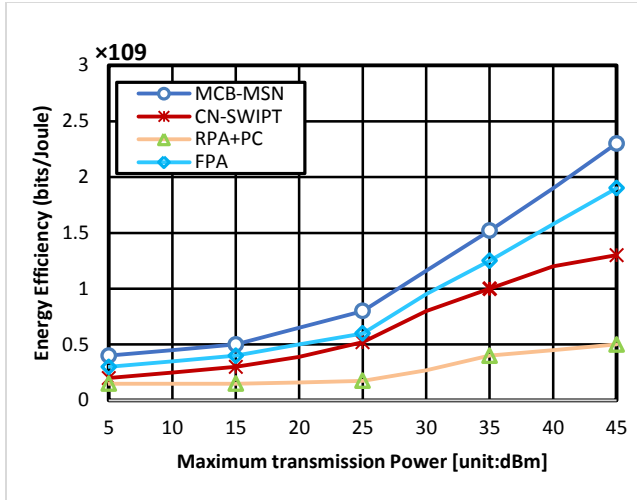


Fig. 9. Network energy efficiency vs. maximum transmission power

6- Conclusion

This paper presents a novel approach to Multi hop Cooperative Beamforming Mobile Sensor Network (MCB-MSN), which not only increases the anonymity of the base station but also maximizes the network energy efficiency in the distributed beamforming by choosing the routes with higher relay densities. In this paper, when the \mathcal{L}_i link cost of the MCB-MSN algorithm is used, the mobile sensor network maintains its level of anonymity significantly more than in a state where anonymity enhancement techniques are not used. In future studies, more MCB-MSN energy consumption should be evaluated in mobile sensor networks considering non-ideal cooperative beamforming conditions so that information needs to be transmitted frequently. In future, we plan to explore the potential of multi-agent smart queuing in various HetNet scenarios such as privacy-aware recommendation and store cell recommendation. We will also plan to examine how to exploit multi-modal data in the mobile sensor networks to further improve the proposed model.

References

[1] O. Cheikhrouhou, A. Koubaa, M. Boujelben, M. Abid. "A lightweight user authentication scheme for wireless sensor networks." In ACS/IEEE International Conference on Computer Systems and Applications-AICCSA 2010 May 16 (pp. 1-7). IEEE.

[2] Bhushan, Bharat, and Gadadhar Sahoo. "Recent advances in attacks, technical challenges, vulnerabilities and their countermeasures in wireless sensor networks." *Wireless Personal Communications* 98, no. 2 (2018): 2037-2077.

[3] AM Somarin, Y. Alaei, MR Tahernezhad, A. Mohajer, M. Barari. "An Efficient Routing Protocol for Discovering the Optimum Path in Mobile Ad Hoc Networks." *Indian Journal of Science and Technology*. 2015 Apr;8(S8):450-5.

[4] Radosavljević, Nemanja, and Đorđe Babić. "Overview of security threats, prevention and protection mechanisms in wireless sensor networks." *J. Mechatron. Autom. Identif. Technol* (2020): 1-6.

[5] L. Zhou, Y. Shan, X. Chen. "An Anonymous Routing Scheme for Preserving Location Privacy in Wireless Sensor Networks". In 2019 IEEE 3rd Information Technology, Networking, Electronic and Automation Control Conference (ITNEC) 2019 Mar 15 (pp. 262-265). IEEE.

[6] J. Yang, ME. Aydin, J. Zhang, C. Maple. "UMTS base station location planning: a mathematical model and heuristic optimisation algorithms". *IET communications*. 2007 Oct 1;1(5):1007-14.

[7] R. Mudumbai, G Barriac, U. Madhow. "On the feasibility of distributed beamforming in wireless networks." *IEEE Transactions on Wireless communications*. 2007 May 21;6(5):1754-63.

[8] R. Mudumbai, DR. Iii, U. Madhow, "Poor HV. Distributed transmit beamforming: challenges and recent progress." *IEEE Communications Magazine*. 2009 Feb 18;47(2):102-10.

[9] Mohajer, Amin, F. Sorouri, A. Mirzaei, A. Ziaeddini, K. Jalali Rad, and Maryam Bavaghar. "Energy-aware hierarchical resource management and Backhaul traffic optimization in heterogeneous cellular networks." *IEEE Systems Journal* (2022).

[10] H. Fang, L. Xu, KK. Choo. "Stackelberg game based relay selection for physical layer security and energy efficiency enhancement in cognitive radio networks." *Applied Mathematics and Computation*. 2017 Mar 1;296:153-67.

[11] AS. Abuzneid, T. Sobh, M. Faezipour, A. Mahmood, J. James. "Fortified anonymous communication protocol for location privacy in WSN: a modular approach." *Sensors*. 2015 Mar;15(3):5820-64.

[12] Mohajer, Amin, Mohammad Hasan Hajimobini, Abbas Mirzaei, and Ehsan Noori. "Trusted-CDS based intrusion detection system in wireless sensor network (TC-IDS)." *Open Access Library Journal* 1, no. 7 (2014): 1-10.

[13] V. Kumar, A. Kumar. "A novel approach for boosting base station anonymity in a wsn". *INTERNATIONAL JOURNAL OF ADVANCED COMPUTER SCIENCE AND APPLICATIONS*. 2017 Sep 1;8(9):114-20.

[14] U. Acharya, M. Younis. "Increasing base-station anonymity in wireless sensor networks." *Ad Hoc Networks*. 2010 Nov 1;8(8):791-809.

[15] YA. Bangash, LF. Zeng, D. Feng. "MimiBS: Mimicking base-station to provide location privacy protection in wireless sensor networks". *Journal of Computer Science and Technology*. 2017 Sep;32(5):991-1007.

[16] P. Gope, AK. Das, N. Kumar, Y. Cheng. "Lightweight and physically secure anonymous mutual authentication protocol for real-time data access in industrial wireless sensor networks". *IEEE transactions on industrial informatics*. 2019 Jan 24;15(9):4957-68.

[17] D. Sharma, A. Goap, AK. Shukla, AP. Bhondekar. "Traffic heterogeneity analysis in an energy heterogeneous WSN routing algorithm". In *Proceedings of 2nd International*

- Conference on Communication, Computing and Networking 2019 (pp. 335-343). Springer, Singapore.
- [18] AA. Mugheri, MA. Siddiqui, M. Khoso. "Analysis on Security Methods of Wireless Sensor Network (WSN)". Sukkur IBA Journal of Computing and Mathematical Sciences. 2018 Jun 26;2(1):52-60.
- [19] Z. Sun, M. Wei, Z. Zhang, G. Qu. "Secure Routing Protocol based on Multi-objective Ant-colony-optimization for wireless sensor networks". Applied Soft Computing. 2019 Apr 1;77:366-75.
- [20] H. Fakhrey, M. Johnston, F. Angelini, R. Tiwari. "The optimum design of location-dependent key management protocol for a multiple sink WSN using a random selected cell reporter. IEEE Sensors Journal". 2018 Sep 24;18(24):10163-73.
- [21] O. Oladayo, A. Ashraf. "A secure and energy-aware routing protocol for optimal routing in mobile wireless sensor networks (MWSNs)". International Journal of Sensors Wireless Communications and Control. 2019 Dec 1;9(4):507-20.
- [22] K. Indira, U. Sakthi. "An efficient anonymous authentication scheme to improve security and privacy in SDN based wireless sensor networks". Indian Journal of Computer Science and Engineering. 2020.
- [23] Mohajer, Amin, Mahya Sam Daliri, A. Mirzaei, A. Ziaeddini, M. Nabipour, and Maryam Bavaghar. "Heterogeneous Computational Resource Allocation for NOMA: Toward Green Mobile Edge-Computing Systems." IEEE Transactions on Services Computing (2022).
- [24] Lee, Jae Seang, Yoon-Sik Yoo, Hyungseok Choi, Taejoon Kim, and Jun Kyun Choi. "Group connectivity-based UAV positioning and data slot allocation for tactical MANET." IEEE Access 8 (2020): 220570-220584.
- [25] Zhang, De-gan, Hao Wu, Peng-zhen Zhao, Xiao-huan Liu, Yu-ya Cui, Lu Chen, and Ting Zhang. "New approach of multi-path reliable transmission for marginal wireless sensor network." Wireless Networks 26, no. 2 (2020): 1503-1517.
- [26] Jaber, Ghada, Rahim Kacimi, and Thierry Gayraud. "Efficient Interest Satisfaction in Content Centric Wireless Sensor Networks." In 2019 16th IEEE Annual Consumer Communications & Networking Conference (CCNC), pp. 1-5. IEEE, 2019.
- [27] M. Bavaghar, A. Mohajer, S. Taghavi Motlagh. "Energy Efficient Clustering Algorithm for Wireless Sensor Networks". Journal of Information Systems and Telecommunication (JIST). 2020 Jun;4(28):238.
- [28] Liu, Hao, Rongbo Zhu, Jun Wang, and Wengang Xu. "Blockchain-Based Key Management and Green Routing Scheme for Vehicular Named Data Networking." Security and Communication Networks 2021 (2021).
- [29] Gai, Keke, Kim-Kwang Raymond Choo, Meikang Qiu, and Liehuang Zhu. "Privacy-preserving content-oriented wireless communication in internet-of-things." IEEE Internet of Things Journal 5, no. 4 (2018): 3059-3067.
- [30] A. Mohajer, M. Barari, H. Zarrabi. "Big Data-based Self Optimization Networking in Multi Carrier Mobile Networks". Bulletin de la Société Royale des Sciences de Liège. 2016 Jan 1;85:392-408.
- [31] M. Dibaei, A. Ghaffari. "Full-duplex medium access control protocols in wireless networks: a survey". Wireless Networks. 2020 May;26(4):2825-43.
- [32] H. Jung, IH. Lee. "Secrecy performance analysis of analog cooperative beamforming in three-dimensional Gaussian distributed wireless sensor networks". IEEE Transactions on Wireless Communications. 2019 Feb 12;18(3):1860-73.
- [33] W. Ge, Z. Zhu, W. Hao, Y. Wang, Z. Wang, Q. Wu, Z. Chu. "AN-Aided Secure Beamforming in Power-Splitting-Enabled SWIPT MIMO Heterogeneous Wireless Sensor Networks". Electronics. 2019 Apr;8(4):459.
- [34] A. Angappan, TP. Saravanabava, P. Sakthivel, KS. Vishvakshan. "Novel Sybil attack detection using RSSI and neighbour information to ensure secure communication in WSN". Journal of Ambient Intelligence and Humanized Computing. 2020 Jul 7:1-2.
- [35] S. Kumar, H. Kim. "Energy efficient scheduling in wireless sensor networks for periodic data gathering". IEEE access. 2019 Jan 10;7:11410-26.
- [36] I. Tomić, JA. McCann. "A survey of potential security issues in existing wireless sensor network protocols". IEEE Internet of Things Journal. 2017 Sep 7;4(6):1910-23.
- [37] Y. Yuan, L. Huo, Z. Wang, D. Hogrefe. "Secure APIT localization scheme against sybil attacks in distributed wireless sensor networks". IEEE Access. 2018 May 15;6:27629-36.
- [38] Y. Liu, Z. Ding, M. Elkashlan, "Poor HV. Cooperative non-orthogonal multiple access with simultaneous wireless information and power transfer". IEEE Journal on Selected Areas in Communications. 2016 Mar 31;34(4):938-53.

Energy-Efficient User Pairing and Power Allocation for Granted Uplink-NOMA in UAV Communication Systems

Seyyed Hadi Mostafavi-Amjad¹, Vahid Solouk², Hashem Kalbkhani^{1*}

¹. Faculty of Electrical Engineering, Urmia University of Technology, Urmia, Iran

². Department of IT and Computer Engineering, Urmia University of Technology, Urmia, Iran

Received: 19 Sep 2021/ Revised: 04 Nov 2021/ Accepted: 07 Dec 2021

Abstract

With the rapid deployment of users and increasing demands for mobile data, communication networks with high capacity are needed more than ever. Furthermore, there are several challenges, such as providing efficient coverage and reducing power consumption. To tackle these challenges, using unmanned aerial vehicles (UAVs) would be a good choice. This paper proposes a scheme for uplink non-orthogonal multiple access (NOMA) in UAV communication systems in the presence of granted and grant-free users. At first, the service area users, including granted and grant-free users, are partitioned into some clusters. We propose that the hover location for each cluster is determined considering the weighted mean of users' locations. We aim to allocate transmission power and form NOMA pairs to maximize the energy efficiency in each cluster subject to the constraints on spectral efficiency and total transmission power. To this end, the transmission powers of each possible pair are obtained, and then Hungarian matching is used to select the best pairs. Finally, finding the flight path of the UAV is modeled by the traveling salesman problem (TSP), and the genetic algorithm method obtains its solution. The results show that the increasing height of the UAV and density of users increases the spectral and energy efficiencies and reduces the outage probability. Also, considering the quality of service (QoS) of granted users for determining the UAV's hover location enhances the transmission's performance.

Keywords: Energy Efficiency; NOMA; Power Allocation, Unmanned Aerial Vehicle (UAV), Uplink, Users Pairing.

1- Introduction

1-1- Motivation

Nowadays, power consumption is an important challenge in wireless networks, and saving users' energy to achieve high performance is essential. Since users' quality of service (QoS) should be satisfied, energy-saving is more challenging in wireless networks. In future wireless communication networks, particularly in the 5th generation (5G) of wireless communications and beyond, the application of unmanned aerial vehicles (UAVs) is proposed for operating as moving aerial relay nodes or moving aerial base stations (ABSs). UAVs, known as drones in a common tongue, has been the subject of a bunch of research over the past few years [1-5]. If they are fine established and well-operated, UAVs can provide reliable and cost-effective wireless communication solutions for many kinds of real-world scenarios [6].

UAVs can operate better and flexibly than the traditional relay nodes in dense areas. Considering them as ABS has several challenges: power consumption, handover management, channel modeling, low-latency control, 3D localization, and interference management [6-9]. Batteries provide UAVs' power in most cases; therefore, power saving is a severe problem in UAV-assisted communications.

1-2- Contributions

This paper investigates the uplink NOMA communication between the ground users, including granted and grant-free and ABS. Our goal is to partition the users into clusters, form NOMA pairs in each cluster, and allocate power to maximize energy efficiency. Each NOMA pair includes one granted and one grant-free user. To this end, at first, UAV flies right into the communication area and simultaneously starts to send a signal periodically to the ground users to get information about their communication conditions that it needs in the next step to clustering the

users. After that, power allocation for each possible pair in each cluster is performed to maximize energy efficiency, while the minimum QoS requirements of users should be satisfied. Finally, efficient pairs are the selection by the Hungarian algorithm. The hover locations for clusters are obtained considering the weighted mean of locations of granted users. Also, the flight path of UAV among different clusters is considered a traversal salesman problem (TSP), and its solution is obtained considering the genetic algorithm. In summary, the contributions of this work are as follows:

- 1) Introducing new ground-ABS uplink transmission scheme over the granted/grant-free ground users and UAV
- 2) Each NOMA group consists of one granted and one grant-free user to ensure the QoS of granted users
- 3) Problem formulation to maximize the energy-efficiency of the proposed transmission subject to the limitation on total transmit power and ensuring the QoS of the users in each NOMA group
- 4) Proposing weighted-mean based on the QoS of the users to determine the hover location of ABS
- 5) Proposing joint user grouping and power allocation that first obtains the transmit power for each possible NOMA pair and then finds the optimal pairs by Hungarian matching algorithm
- 6) Considering the flight path of UAV as TSP and utilizing the genetic algorithm to solve it.

The rest of this paper is organized as follows. The previous works are reviewed in Section 2. The system model is described in Section 3. The proposed user clustering, NOMA pair forming, and power allocation is presented in Section 4. Simulation results are given in Section 5. Finally, conclusions remarks are provided in Section 6.

2- Previous Works

Here, we review the related researches which consider UAV-assisted wireless communication topics. In [10], a rotary-wing UAV performs the send and collect data task to/from multiple ground users. This research aimed to optimize the total UAV power consumption by minimizing propulsion and data transmission powers while satisfying each ground node's minimum quality-of-service (QoS) requirement in the uplink direction. Energy-efficient UAV communication with a ground terminal in the downlink direction via optimizing the UAV's trajectory was studied in [11]. The authors aimed to design a new specimen that considers both the throughput and energy consumption of UAV together. Serving cell edge users by UAV and offloading the data from the base station in downlink direction by circle path for UAV was

studied in [12]. The goal was to optimize UAVs' resource allocation, user partitioning, and trajectory by maximizing energy efficiency. A new modularity-based dynamic clustering relying on UAVs' modified Louvain method was studied in [13]. The authors aimed to save the transmitted power of mobile devices in the uplink direction by locating the UAVs on the user clusters' centroids. Resource allocation and trajectory design in downlink direction were formulated as an energy-efficient problem in [14], which jointly optimizes the transmit power, user scheduling, and trajectory and velocity of UAV. A real-time resource allocation algorithm for maximizing the energy efficiency in downlink direction by jointly optimizing the energy-harvesting time and power control for the considered device-to-device (D2D) communication embedded with UAV was proposed in [15].

The optimum establishing of UAV as a relay for maximizing the reliability was studied in [16]. The total power loss, outage, and error rate were considered as the reliability parameters, and optimum height was investigated for static and mobile UAVs. However, it would be better when they consider more than one user on the cell edge. In [17], the effective use of flight-time constrained UAVs as aerial ABSs was investigated to provide coverage for ground users. Notably, a novel framework was proposed for optimizing the average number of bits transmitted to users and UAVs' hover duration. The authors in [18] proposed an optimum placement algorithm for UAVs that maximizes the number of covered users with minimum transmit power. They have detached the UAV located in the vertical dimension from the horizontal dimension, which simplifies the placement problem. The authors [19] characterized UAV-based communication's latency, reliability, and network availability of ultra-reliable and low-latency communications (URLLC). The height of UAVs and the bandwidth allocation were optimized to minimize the required total bandwidth of URLLC for a given density of UAVs. It was shown that the probability of the line-of-sight (LoS) path and the network availability is strictly concave for the distance between the ground user and UAV. The impact of the height of UAVs connected to the cellular network in uplink was studied in [19]. In [20, 21], a UAV was considered to collect data from a set of sensors with a fixed location. The goal was to minimize the UAV's total flight time while each sensor could successfully upload its data using a given amount of energy. The problem of trajectory design for UAVs to maximize satisfied users was studied in [22].

Using non-orthogonal multiple access (NOMA) has its benefits and challenges in comparison with other multiple access techniques such as orthogonal-frequency-division-

multiple-access (OFDMA) or orthogonal multiple access (OMA). Saving bandwidth by pairing strong and weak users in the same time slots is the most beneficial of NOMA. However, this pairing causes intra-time slot and intra-cell interference challenges. To tackle these challenges, the following works try to solve the problems. To facilitate the serving ground users in a cell, user clustering is a crucial element. Hence, dynamic user scheduling and power allocation problem was proposed in [23] to coordinate the intra-cell interference by minimizing the total power consumption. In [24], sum-rate maximization for uplink and downlink NOMA under the constraints of transmission power limitation, minimum rate requirements of users, and operation constraints were formulated. Machine learning-based user clustering and power allocation algorithms for mmWave-NOMA transmission were considered in [25]. Energy-efficient resource allocation for the uplink of hybrid NOMA and OMA transmission was considered in [26], obtained by jointly optimizing the user clustering, channel assignment, and power allocation. High-rate NOMA, where multiple users share a single zero-forcing beamforming vector, was proposed in [27]. The QoS of all clustered users was satisfied to maintain fairness among the users. In [28], ground-aerial uplink-NOMA of cellular networks was investigated, where ground base stations serve a UAV user and multiple ground users. They aimed to minimize the UAV mission completion time by jointly optimizing the UAV trajectory and association order while considering the UAV's interference to non-associated ground base stations. In [29], applying of NOMA technique to UAV to cellular BSs uplink communication, under the spectrum sharing with the existing ground users was investigated, and a new cooperative NOMA scheme was proposed to reduce the intense uplink interference due to the UAV's LoS channels with ground BSs in cellular-connected UAV communication. A combination of multi-UAV communication and NOMA was proposed in [30] to construct the high capacity uplink for the internet of things (IoT) which was achieved by jointly optimizing the sub-channel assignment, transmit power, and flying heights of UAVs. A novel framework for UAV networks with massive access capability supported by NOMA was proposed in [31].

User grouping into two sets to achieve low-latency access and reduce signaling overhead was investigated in [32], where the scheduled-access and random-access users are considered granted and grant-free users, respectively. In [33], NOMA-assisted semi-grant-free transmission was studied, which is investigated compromise between grant-free and grant-based users.

3- System Model

3-1- User Distribution and Transmission Model

As shown in Fig.1, it is assumed that N granted and $M > N$ grant-free terrestrial users are represented by $\mathbf{u} = \{u_1, u_2, \dots, u_N\}$ and $\mathbf{v} = \{v_1, v_2, \dots, v_M\}$, respectively. According to the spatial Poisson point process (SPPP) distribution with density λ_u , these users are distributed in the service area. The locations of the granted user u_i and grant-free user v_j in two-dimensional space are respectively denoted by (x_i, y_i) and $(\tilde{x}_j, \tilde{y}_j)$. These users are partitioned into n_c clusters $\{C_1, \dots, C_{n_c}\}$. The sets of granted and grant-free users belonging to the cluster C_k are given by \mathbf{u}_k and \mathbf{v}_k , respectively. The number of granted and grant-free users in the cluster C_k is given by $N_c^{(k)}$ and $M_c^{(k)}$, respectively, therefore we have $\sum_{k=1}^{n_c} N_c^{(k)} = N$ and $\sum_{k=1}^{n_c} M_c^{(k)} = M$. Suppose that \mathbf{S} is the binary matrices with the size of $(N + M) \times n_c$. If the granted user u_i belongs to the cluster C_k , we have $\mathbf{S}(i, k) = 1$, otherwise $\mathbf{S}(i, k) = 0$. On the other side, $\mathbf{S}(j + N, k) = 1$ denotes that grant-free user v_j is considered in the cluster C_k , else $\mathbf{S}(j + N, k) = 0$.

Considering Fig. 1, granted user u_i and grant-free user v_j belonging to cluster C_k form the two-user NOMA group (u_i, v_j) to transmit their data in uplink direction to ABS. The matrix \mathbf{G}_k , which presents the pairing of users in the cluster C_k , has the size of $N_c^{(k)} \times M_c^{(k)}$. If the granted user u_i and grant-free user v_j form the two-user NOMA group in the cluster C_k , then $\mathbf{G}_k(i, j) = 1$, otherwise $\mathbf{G}_k(i, j) = 0$.

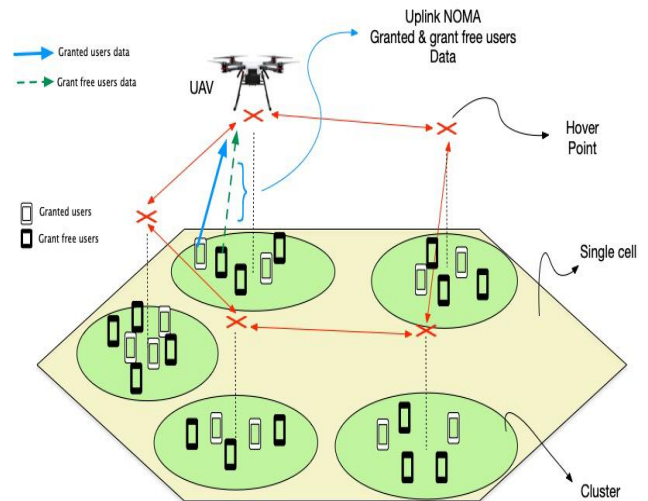


Fig. 1. The system and transmission model that considered in this paper.

The UAV has a hover for a specific time to collect data of the users of each cluster. Since the granted users have priority over the grant-free users, we propose that the hover location is determined based on the weighted mean of granted users in the cluster. Suppose that the hover location of UAV in three-dimensional space for the cluster C_k is denoted by $(\hat{x}_k, \hat{y}_k, \hat{\rho}_k)$, $k = 1, \dots, n_c$, where $\hat{\rho}_k$ is the height of UAV, then we have

$$\hat{x}_k = \sum_{i=1}^N \alpha_{i,k} x_i; k = 1, \dots, n_c \quad (1)$$

$$\hat{y}_k = \sum_{i=1}^N \alpha_{i,k} y_i; k = 1, \dots, n_c \quad (2)$$

where $\alpha_{i,k}$ denotes the weight of the granted user u_i for the cluster C_k , which is obtained based on the normalized QoS of the user. Since each user only belongs to one cluster, $\alpha_{i,k} \neq 0$ only for one cluster, and it is zero for other clusters. The minimum rate requirement of the users is considered as QoS in this paper. Let $\mathbf{q} = \{q_1, q_2, \dots, q_N\}$ and $\bar{\mathbf{q}} = \{\bar{q}_1, \bar{q}_2, \dots, \bar{q}_M\}$ respectively denote the QoS of granted and grant-free users in terms of spectral efficiency in bit/sec/Hz. Then, the weight $\alpha_{i,k}$ is computed as:

$$\alpha_{i,k} = \frac{q_i}{\sum_{i=1}^N S^{(i,k)} q_i} \quad (3)$$

3-2- Channel Model

Given the granted user u_i located at (x_i, y_i) and the ABS located at $(\hat{x}_k, \hat{y}_k, \hat{\rho}_k)$, the path loss between the ABS and user u_i will be [34]:

$$l_{i,k} = 20 \log_{10} \left(\frac{4\pi f_c d_{i,k}}{c} \right) + \vartheta_{i,k} \quad (4)$$

where f_c is the carrier frequency in Hz and c is the speed of light in m/s, and also $d_{i,k}$ is the Euclidean distance between UAV and user u_i in meter calculated as:

$$d_{i,k} = \sqrt{(x_i - \hat{x}_k)^2 + (y_i - \hat{y}_k)^2 + \hat{\rho}_k^2} \quad (5)$$

Also, $\vartheta_{i,k}$ is the normally distributed additional loss depending on environment conditions and distributed as $N(\mu_L, \sigma_L^2)$ and $N(\mu_{NL}, \sigma_{NL}^2)$ for the LoS and NLoS links, respectively. The probability of having an LoS link between the user u_i and ABS at hover location (x_k, y_k) is obtained as:

$$p_{i,k}^L = \frac{1}{1 + a \exp\left(-\frac{180}{\pi} b (\theta_{i,k} - a)\right)} \quad (6)$$

where a and b are constants, and they change depending on the environment and $\theta_{i,k} = \sin^{-1}(\hat{h}_k/d_{i,k})$ is the elevation angle of ABS for the desired user u_i . Then, the probability of having an NLoS link is $p_{i,k}^{NL} = 1 - p_{i,k}^L$ [35]. Eventually, the total path loss from granted user u_i ABS to $L_{i,k}$, is obtained as

$$L_{i,k} = p_{i,k}^L l_{i,k}^L + p_{i,k}^{NL} l_{i,k}^{NL} \quad (7)$$

A similar procedure presented in equations (4)-(7) can be used to obtain the path loss of grant-free user v_j , i.e., $\bar{L}_{j,k}$.

3-3- Spectral and Energy Efficiency

In two-user NOMA transmission, the user with the higher channel gain is called the strong user, and another one is the weak user. The transmitted signal of these users experiences distinct channel gains. In uplink two-user NOMA, the received signal at location k of ABS due to the transmission of pair (u_i, v_j) , i.e., $m_{k,i,j}^r$, can be obtained as [36]:

$$m_{i,j,k}^r = \sqrt{P_i} \square_{i,k} m_i^t + \sqrt{\bar{P}_j} \bar{\square}_{j,k} \bar{m}_j^t + n \quad (8)$$

where m_i^t and \bar{m}_j^t denote the transmitted signal from u_i and v_j , respectively. P_i and \bar{P}_j respectively signify the transmit power of granted user u_i and grant-free user v_j . $\square_{i,k}$ and $\bar{\square}_{j,k}$ represent the channel gain from u_i and v_j to ABS located at k th location, and n is the white noise with power spectral density P_{noise} . The results reported in [36] show that if there is enough separation between weak and strong users, the spectral efficiency of the strong user outperforms that of the weak user. Since satisfying the QoS of granted users has priority over grant-free users, granted users are considered strong users, and grant-free users transmit their data as weak users in each NOMA group.

Transmission of grant-free users interferes with the transmission of granted users and reduces the SNR of the granted user. On the other side, the transmission of the grant-free user receives zero interference from the transmission of the granted user during successive interference cancellation (SIC). Hence, the SINR of the transmission of pair (u_i, v_j) in ABS located in k th location is obtained as follows [36]:

$$\gamma_{i,k} = \frac{P_i L_{i,k}^{-1}}{P_j \bar{L}_{j,k}^{-1} + P_{noise}} \quad (9)$$

$$\bar{\gamma}_{j,k} = \frac{\bar{P}_j \bar{L}_{j,k}^{-1}}{P_{noise}} \quad (10)$$

Consequently, the SE of NOMA transmission of pair (u_i, v_j) at the k th location of ABS is obtained as:

$$\eta_{SE}^{i,j,k} = \log_2(1 + \gamma_{i,k}) + \log_2(1 + \bar{\gamma}_{j,k}) \text{ bits/sec/Hz} \quad (11)$$

With given the spectral efficiency of the pair (u_i, v_j) , the energy efficiency is calculated as follows:

$$\eta_{EE} = \frac{\eta_{SE}}{P_{tot}} (1 - P_{out}) \text{ bits/sec/Hz/J} \quad (12)$$

where η_{SE} is the total spectral efficiency, P_{tot} is the total power consumption, which is the sum of transmit power (P_{Tx}), circuit powers of terrestrial users (P_{cir}), and power consumed by UAV (P_{ABS}), P_{out} denotes the outage probability. Eventually, P_{tot} described as follow:

$$P_{tot} = P_{Tx} + P_{cir} + P_{ABS} \quad (13)$$

The circuit power of ABS consists of two parts, including the circuit power of ABS in hovering ($P_{\square over}$) and flying times (P_{flight}). The torque coefficient of UAV, q_c , is given as [10]:

$$q_c = \frac{\delta}{8} + (1+k) \frac{w^{1.5}}{\sqrt{2\rho^2 s A^2 \Omega^3 R^3}} \quad (14)$$

Therefore, the corresponding power consumption in hovering time of UAV can be described as [10]:

$$P_{\square over} = q_c \rho s A \Omega^3 R^3 \quad (15)$$

and by substitution of q_c in $P_{\square over}$, we have [10]:

$$P_{\square over} = \frac{\delta}{8} \rho s A \Omega^3 R^3 + (1+k) \frac{w^{1.5}}{\sqrt{2\rho A}} \quad (16)$$

The required power for the flight time of rotary-wing UAVs is more intricate than the fixed-wing peer. However, by some mild assumption, the pull coefficient of the blade area is constant, so the torque coefficient for the UAV in flight time with zero climbing angle and speed v_u is given as:

$$q_c = \frac{\delta}{8} (1 + 3\mu^2) + (1+k) \mu_i t_{CD} + \frac{1}{2} \hat{u}_u^3 d_0 \quad (17)$$

By substituting $\mu \approx \hat{v}_u = \frac{v_u}{\Omega R}$ and $t_{CD} = \frac{T}{\rho s A \Omega^2 R^2}$, q_c can be written as a function of forwarding speed v_u and rotor thrust T as follow:

$$q_c(v_u, T) = \frac{\delta}{8} \left(1 + \frac{3v_u^2}{\Omega^2 R^2}\right) + \frac{(1+k)T\lambda_i}{\rho s A \Omega^2 R^2} + \frac{1}{2} d_0 \frac{v_u^3}{\Omega^3 R^3} \quad (18)$$

Eventually, by definition of torque coefficient, the required power for flight time can be written as follow:

$$P_{flight} = q_c \rho s A \Omega^3 R^3 \quad (19)$$

3-4- Outage Probability

Outage probability is defined as the probability that the SNR or spectral efficiency at the receiver becomes lower than the predefined (or threshold) value. Considering the SIC process, for the strong (or granted) user u_i , the outage probability is obtained as [37]:

$$P_{out}^{u_i} = 1 - P_C^{u_i} \quad (20)$$

where $P_C^{u_i}$ is the probability that the transmit message of the strong user u_i is correctly detected at the receiver, which is calculated as:

$$P_C^{u_i} = \Pr\{\gamma_{i,k} \geq q_i\} = \Pr\left\{\frac{P_i L_{i,k}^{-1}}{P_j \bar{L}_{i,j}^{-1} + P_{noise}} \geq q_i\right\} \quad (21)$$

Similarly, for weak (or grant-free) user v_j , we have:

$$P_{out}^{v_j} = 1 - P_C^{v_j} \quad (22)$$

where the probability of correct detection of the message of the grant-free user is calculated as [37]:

$$\begin{aligned} P_C^{v_j} &= \Pr\{\gamma_{i,k} \geq q_i, \bar{\gamma}_{i,k} \geq \bar{q}_j\} \\ &= \Pr\left\{\frac{P_i L_{i,k}^{-1}}{P_j \bar{L}_{i,j}^{-1} + P_{noise}} \geq q_i, \frac{P_j \bar{L}_{i,k}^{-1}}{P_{noise}} \geq \bar{q}_j\right\} \end{aligned} \quad (23)$$

4- Proposed User Clustering, Power Allocation, and NOMA pair Forming

4-1- Problem Formulation

Over the past decades, energy efficiency has been studied from the information theory perspective. Due to the power limitation of ABS, it is worth considering the energy-efficient transmission scheme and maximizing energy efficiency. Therefore, the proposed scheme for users clustering and joint user pairing and power allocation can be formulated as follows:

$$(\mathbf{S}^*, \mathbf{G}^*, \mathbf{P}^*) = \operatorname{argmax}(\eta_{EE}) \quad (24)$$

subject to:

$$\begin{aligned} (S1) \quad & \gamma_{i,k} \geq \gamma_i^{th} \\ (S2) \quad & \bar{\gamma}_{j,k} \geq \bar{\gamma}_j^{th} \\ (S3) \quad & P_i + \bar{P}_j \leq P_{max} \\ (S4) \quad & \sum_{k=1}^{nc} \mathbf{C}(i, k) = 1, \forall i = 1, \dots, N \\ (S5) \quad & \sum_{k=1}^{nc} \mathbf{G}(i, j) = 1, \forall i = 1, \dots, N \\ (S6) \quad & \sum_{i=1}^{nc} \mathbf{G}(i, j) \leq 1, \forall j = 1, \dots, M \end{aligned} \quad (25)$$

The constraints S1 and S2 respectively demonstrate that the QoS of granted and grant-free users should be satisfied, where $q_i = \log_2(1 + \gamma_i^{th})$ and $\bar{q}_j = \log_2(1 + \bar{\gamma}_j^{th})$. Furthermore, S3 determines the upper limit of transmit power of terrestrial users. Constraint S4 specifies that each terrestrial user must be included only in one cluster. According to S5, each granted user can pair with a grant-free user and transmit its data. While, according to S6, there is no guarantee for grant-free users to transmit data.

Considering the energy efficiency optimization problem and constraints given in equations (24)-(25), this problem is non-convex, and obtaining a solution requires vast computational complexity. Hence, we propose to partition it into two sub-problems to obtain its solution. First sub-problem partitions the terrestrial users into several clusters. After clustering, the joint user pairing and power allocation problem was formulated for each cluster to maximize the energy efficiency of each cluster.

4-2- Proposed Solution

The proposed solution for the problem given in equations (24)-(25) is explained in Algorithm 1. The proposed solution generally consists of three steps, including cluster

forming, power allocation, and NOMA pair forming. In the following, each step is explained in detail.

Algorithm 1. The proposed solution for user clustering, power allocation, and NOMA pair forming

|| Cluster forming and obtaining flight path

1. UAV flights over the service area to partition the users to some clusters
2. Compute the hover location for each cluster considering equations (1)-(3)
3. Obtain the flight path by solving TSP via the genetic algorithm

|| Power allocation and NOMA pair forming

4. **for** each cluster, **do**
 5. Perform power allocation for all possible pairs of granted and grant-free users using equations (32)-(33)
 6. Compute energy efficiency of each pair using equation (34)
 7. Select the best pairs using the Hungarian algorithm (Algorithm 2)
 8. **end for**
-

4-2-1- User Clustering

In order to cluster the terrestrial users, the UAV starts to fly over the service area from a random location. UAV broadcasts the initialization signal and waits to receive the first acknowledgment signal from terrestrial users. This acknowledgment signal contains users' position, desired QoS, and type (granted or grant-free). When the UAV receives the first acknowledgment signal, it starts to create the first cluster. This process continues until the acknowledgment of the first signal is not received in the UAV. At this time, the first cluster is created, and all users whose acknowledgment signal was received are considered in the first cluster. Then, UAV continues to fly over the area to create the second cluster. This process will repeat until all users to be placed in one cluster.

After finishing the clustering process, the UAV computes the hover locations considering equations (1)-(2). In order to minimize the power consumption during fly and hover, the flight path should be minimized. We use TSP integer linear programming with the Dantzig-Fulkerson-Johnson formulation (DFJ) algorithm to obtain the shortest flight path. Suppose that the set $\{(\hat{x}_1, \hat{y}_1, \hat{z}_1), \dots, (\hat{x}_{n_c}, \hat{y}_{n_c}, \hat{z}_{n_c})\}$ contains the coordinates of hover locations. TSP solves the following problem:

$$\min \sum_{i=1}^{n_c} \sum_{j=1, j \neq i}^{n_c} \beta_{ij} D_{i,j} \quad (26)$$

subject to:

$$\begin{aligned} \beta_{ij} &\in \{0,1\}, i, j = 1, \dots, n_c \\ \sum_{i=1, i \neq j}^{n_c} \beta_{ij} &= 1, j = 1, \dots, n_c \\ \sum_{i=1, i \neq j}^{n_c} \beta_{ij} &= 1, i = 1, \dots, n_c \\ \sum_{i \in Q} \sum_{j \in Q} \beta_{ij} &\leq |Q| - 1, \\ &\forall Q \subset \{1, \dots, n_c\}, |Q| \geq 2 \end{aligned} \quad (27)$$

where D_{ij} is the Euclidean distance between hover locations i and j and β_{ij} is a binary variable defined as:

$$\beta_{i,j} = \begin{cases} 1 & \text{UAV goes from point } i \text{ to point } j \\ 0 & \text{Otherwise} \end{cases} \quad (28)$$

The solution given in [38] is utilized to obtain the optimum flight path of the UAV.

4-2-2- Joint Power Allocation and User Pairing

After clustering the users, the user pairing and power allocation should be performed for each cluster. Hence, the problem demonstrated in equations (24)-(25) is simplified for cluster C_k as:

$$(\mathbf{G}_k^*, \mathbf{P}_k^*) = \operatorname{argmax}(\eta_{EE}^{(k)}) \quad (29)$$

subject to:

$$\begin{aligned} (S1) \quad & \gamma_{i,k} \geq \gamma_i^{\text{th}} \\ (S2) \quad & \bar{\gamma}_{j,k} \geq \bar{\gamma}_j^{\text{th}} \\ (S3) \quad & P_i + \bar{P}_j \leq P_{\max} \\ (S4) \quad & \sum_{k=1}^{M_c^{(k)}} \mathbf{G}_k(i, k) = 1, \forall i = 1, \dots, N_c^{(k)} \\ (S5) \quad & \sum_{i=1}^{N_c^{(k)}} \mathbf{G}_k(i, j) \leq 1, \forall j = 1, \dots, M_c^{(k)} \end{aligned} \quad (30)$$

where $\eta_{EE}^{(k)}$ is the energy efficiency of cluster C_k which is computed as:

$$\eta_{EE}^{(k)} = \frac{\sum_{i=1}^{N_c^{(k)}} (\log_2(1+\gamma_{i,k}) + \log_2(1+\bar{\gamma}_{i,k}))}{\sum_{i=1}^{N_c^{(k)}} (P_i + \bar{P}_i) + 2N_c^{(k)} P_{\text{cir}, UE}} \quad (31)$$

where $P_{\text{cir}, UE}$ is the circuit power of each user. To solve this problem, at first, we perform power allocation for all possible pairs of granted and grant-free users, and then, the Hungarian algorithm is utilized to select the pairs that maximize the energy efficiency of the cluster.

Suppose that granted user u_i and grant-free user v_j form the two-user NOMA pair. We propose to minimize the power consumption to maximize energy efficiency. The transmit powers of them are calculated to satisfy the minimum QoS requirement of them as follows:

$$\bar{P}_j = \bar{\gamma}_{j,k} \bar{L}_{j,k} P_{\text{noise}} \quad (32)$$

$$P_i = \gamma_{i,k} (\bar{P}_j \bar{L}_{j,k}^{-1} + P_{\text{noise}}) L_{i,k} \quad (33)$$

After that, the energy efficiency of the pair (u_i, v_j) is computed as follows:

$$\eta_{EE}^{(k)}(i, j) = \frac{\log_2(1+\gamma_{i,k}) + \log_2(1+\bar{\gamma}_{j,k})}{P_i + \bar{P}_j + 2P_{\text{cir}, UE}} \quad (34)$$

Obtaining $\eta_{EE}^{(k)}(i, j)$ results in $N_c^{(k)} \times M_c^{(k)}$ the matrix for cluster C_k . The next step is to select the pairs from this matrix to maximize the energy efficiency of cluster C_k . After forming the Hungarian matrix, which contains the select energy efficiency of different pairs of users, the best

pairs should be selected from them. There are two ways to solve this problem; adjacency matrix and bipartite graph. The bipartite graph can easily represent by an adjacency matrix. As an example, suppose seven users where the users in rows of the matrix belong to granted users and the columns belong to grant-free users, which is shown as follow:

$$\eta_{EE}^{(k)} = \begin{matrix} & v_1 & v_2 & v_3 & v_4 \\ \begin{matrix} u_1 \\ u_2 \\ u_3 \end{matrix} & \begin{pmatrix} 2 & 1 & 3 & 4 \\ 8 & 4 & 2 & 6 \\ 3 & 12 & 6 & 9 \end{pmatrix} \end{matrix}$$

Note that the Hungarian method assigns a set of minimum optimal values of the matrix, and the energy efficiency problem must be maximized. Hence, we first convert all the arrays into the familiar form to get the maximized values from these minimum arrays as follow:

$$\begin{matrix} v_1 & v_2 & v_3 & v_4 \\ \begin{matrix} u_1 \\ u_2 \\ u_3 \end{matrix} & \begin{pmatrix} 2 & 1 & 3 & 4 \\ 8 & 4 & 2 & 6 \\ 3 & 12 & 6 & 9 \end{pmatrix} \end{matrix} \xrightarrow{\text{Inverse each elements}} \begin{matrix} v_1 & v_2 & v_3 & v_4 \\ \begin{matrix} u_1 \\ u_2 \\ u_3 \end{matrix} & \begin{pmatrix} 0.50 & 1.0 & 0.33 & 0.25 \\ 0.12 & 0.25 & 0.50 & 0.16 \\ 0.33 & 0.08 & 0.16 & 0.11 \end{pmatrix} \end{matrix}$$

It is essential to say that the Hungarian method is proper when the matrix is square. Therefore, if the assignment matrix is not square, we must turn it into square form by adding dummy rows or columns. The dummy arrays can be in two forms; they can be equal to the maximum matrix array or be a line with zero numbers; however, zero numbers are recommended. The solution for the Hungarian method is shown in Algorithm 2. The solution of the defined example with the Hungarian matrix method is shown step by step as follow:

(Step 1) Subtract the smallest value in each row from the other values in the row

(Step 2) Each column has zero, so no need to subtract the minimum value from each column.

(Step 1)

$$\begin{matrix} v_1 & v_2 & v_3 & v_4 \\ \begin{matrix} u_1 \\ u_2 \\ u_3 \\ Du \end{matrix} & \begin{pmatrix} 0.5 & 1 & 0.33 & 0.25 \\ 0.12 & 0.25 & 0.5 & 0.16 \\ 0.33 & 0.08 & 0.16 & 0.11 \\ 0 & 0 & 0 & 0 \end{pmatrix} \end{matrix} \begin{matrix} -0.25 \\ -0.12 \\ -0.08 \\ \end{matrix}$$

(Step 2)

$$\begin{matrix} v_1 & v_2 & v_3 & v_4 \\ \begin{matrix} u_1 \\ u_2 \\ u_3 \\ Du \end{matrix} & \begin{pmatrix} 0.25 & 0.75 & 0.08 & 0 \\ 0 & 0.13 & 0.38 & 0.04 \\ 0.25 & 0 & 0.08 & 0.03 \\ 0 & 0 & 0 & 0 \end{pmatrix} \end{matrix}$$

(Step 3) Draw lines through the row and columns that have the 0 entries such that the fewest possible lines are drawn. There are four lines drawn, which is equal to the matrix dimension, so there is the optimal number of zeroes.

Algorithm. 2. The Hungarian algorithm using an adjacency matrix

1. Convert all the arrays into the reciprocal form
2. **if** the number of rows and columns are not equal, **then**
3. Add dummy rows or columns to square the matrix
4. Subtract the smallest entry in each row from all the other entries in the row
5. **if** there is any column without zero, **then**
6. Subtract the smallest entry in each column from all the other entries in the column
7. Cover the rows and columns that have the 0 entries with the fewest lines possible are drawn
8. **if** there the number of lines drawn is equal to the number of rows, **then**
9. An optimal assignment of zeros is possible, and the algorithm is finished.
10. **else if** the number of lines is less than number of rows, **then**
11. The optimal number of zeroes is not yet reached.
12. **Go to** the next steps.
13. Find the smallest entry not covered by any line.
14. Subtract this entry from each row that is not crossed out,
15. Then add it to each column that is crossed out.
16. **end**

(Step 4) Highlight the selected zeros

(Step 3)

$$\begin{matrix} v_1 & v_2 & v_3 & v_4 \\ \begin{matrix} u_1 \\ u_2 \\ u_3 \\ Du \end{matrix} & \begin{pmatrix} 0.25 & 0.75 & 0.08 & 0 \\ 0 & 0.13 & 0.38 & 0.04 \\ 0.25 & 0 & 0.08 & 0.03 \\ 0 & 0 & 0 & 0 \end{pmatrix} \end{matrix}$$

(Step 4)

$$\begin{matrix} v_1 & v_2 & v_3 & v_4 \\ \begin{matrix} u_1 \\ u_2 \\ u_3 \\ Du \end{matrix} & \begin{pmatrix} 0.25 & 0.75 & 0.08 & 0 \\ 0 & 0.13 & 0.38 & 0.04 \\ 0.25 & 0 & 0.08 & 0.03 \\ 0 & 0 & 0 & 0 \end{pmatrix} \end{matrix}$$

(Step 5) Replace the original values. **(Step 6)** Replace the primary values of energy efficiency to get which users can be optimally pairs. As shown, there are four lines drawn, and it is equal to the dimension of the matrix, so the algorithm is finished optimally. However, if there are drawn lines less than the matrix dimension, it should follow the algorithm's rules. The Hungarian method, which is shown in Algorithm 2, forms pairs of users optimally. To the extent, one granted user should pair with a grant-free user, which are given as strong and weak users.

(Step 5)

$$\begin{matrix} v_1 & v_2 & v_3 & v_4 \\ \begin{matrix} u_1 \\ u_2 \\ u_3 \\ Du \end{matrix} & \begin{pmatrix} 0.5 & 1 & 0.33 & 0.25 \\ 0.12 & 0.25 & 0.5 & 0.16 \\ 0.33 & 0.08 & 0.16 & 0.11 \\ 0 & 0 & 0 & 0 \end{pmatrix} \end{matrix}$$

(Step 6)

$$\begin{matrix} v_1 & v_2 & v_3 & v_4 \\ \begin{matrix} u_1 \\ u_2 \\ u_3 \\ Du \end{matrix} & \begin{pmatrix} 2 & 1 & 3 & 4 \\ 8 & 4 & 2 & 6 \\ 3 & 12 & 6 & 9 \\ 0 & 0 & 0 & 0 \end{pmatrix} \end{matrix}$$

5- Simulation Results

5-1- Simulation Setup

This section provides numerical results to evaluate the performance of the proposed user clustering, power allocation, and NOMA pair forming to maximize energy efficiency. The UAV acts as a flying ABS in the simulation area and serves the users randomly distributed according to SPPP with density λ_u . The parameters used in simulations are given in Table 1.

Table 1. Parameters used in simulations

Parameter	Value
Simulation area	5×5 km
Density of users	(1~3)×10 ⁻⁴
Height of ABS	100 ~ 800 m
Maximum total transmission power	23 dBm
The minimum acceptable received power	-90 dBm
Noise power	-130 dBm
Carrier frequency	1.2 GHz
Minimum acceptable SNR of granted users	[2 8] dB
Minimum acceptable SNR of grant-free users	[1 3] dB

Results in terms of spectral efficiency, energy efficiency, and outage probability are obtained for each pair of UAV height and density of users. For each pair, we run *Monte Carlo* simulations for 10⁵ trials, and in each trial, the users' locations are generated using SPPP with specific density. Finally, results were averaged.

5-2- Spectral Efficiency

The spectral efficiency for different heights of UAV and density of users is given in Fig. 2 for the total transmission power of 23 dBm. It is observed that increasing the height of UAV and density of users increases the spectral efficiency. In traditional 2D wireless networks such as cellular networks, path loss increases as the distance of users from the base station increases. However, in UAV networks, increasing distance will not necessarily increase the path loss because the probability of the LoS link increases. Increasing the probability of LoS link reduces ν in path loss; hence, overall path loss reduces, SNR increases, and spectral efficiency increases. On the other side, increasing the density of users increases the number of users; therefore, there are more candidate users to form NOMA pairs with better channel conditions. As each pair transmits its data in a specific time slot, increasing the number of users does not increase the interference, and spectral efficiency increases. In summary, for the constant density of users, increasing the height of ABS enhances the spectral efficiency by reducing the path, resulting in higher SNR. For the constant height of ABS, increasing the density of users enhances the spectral efficiency by constructing the pairs with higher SNRs.

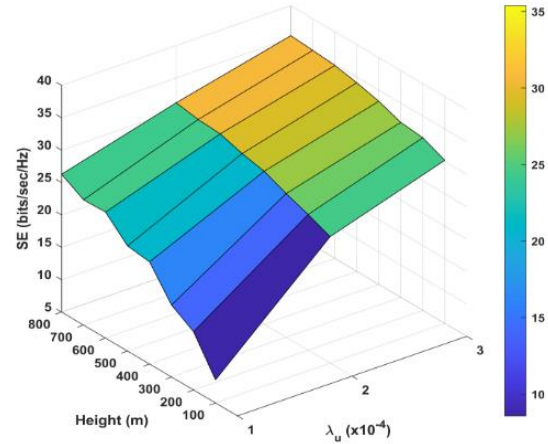


Fig. 2. Spectral efficiency for different densities of users and height of UAV.

Fig. 3 shows the impact of hover location on spectral efficiency. Three schemes are utilized to determine the hover location; 1) proposed weighted-mean (WM) of granted users, 2) equal-weight mean (EWM), where we simply consider the mean of the location of the granted user, and 3) random in which UAV randomly hover in the area of the cluster. It is observed that the proposed WM method achieves higher spectral efficiency than the other schemes. Proposed WM determines the hover location near granted users with higher QoS requirements, achieving higher spectral efficiency.

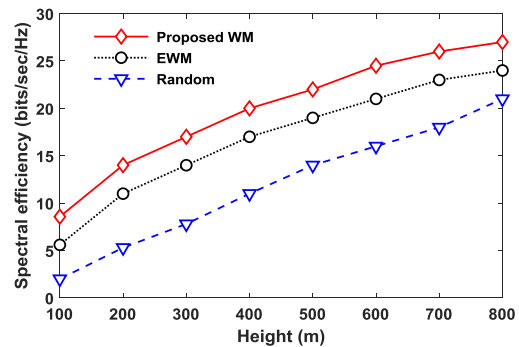


Fig. 3. The effect of hover location on the spectral efficiency

5-3- Energy Efficiency

In Fig. 4, the energy efficiency is given for several heights of UAV and the density of users. It is observed that similar to spectral efficiency, increasing the height of UAV and density of users increases the energy efficiency. Reducing path loss by increasing height reduces the transmit power to satisfy the spectral efficiency; hence energy efficiency increases. On the other side, increasing the density of users provides more candidates for NOMA pairs, which reduces the transmission power and increases energy efficiency. Also, Fig. 5 compares the energy efficiencies obtained for different hover locations. It is observed that the proposed WM scheme outperforms the other schemes considerably.

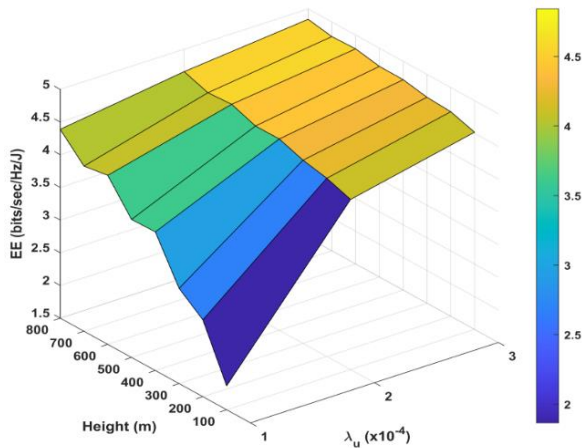


Fig. 4. Energy efficiency for different densities of users and heights of UAV.

5-4- Outage Probability

It was mentioned that increasing the height of the UAV enhances the uplink transmission by reducing the path loss; therefore, it is expected that outage probability reduces by increasing the height of the UAV, which is depicted in Fig. 6. As spectral and energy efficiencies, outage probability enhances by increasing the height of ABS location. Also, increasing the density of users reduces the outage probability. Fig. 7 compares the outage probability for different hover locations. As expected, the proposed WM scheme has the lowest outage probability since it determines the hover location of the UAV, considering the weighted mean of users' locations based on their QoS. This approach reduces UAV distance from the users with high QoS requirement and increases their SNR, which reduces the outage probability. Also, increasing the height of the UAV reduces the path loss by increasing the probability of LoS link resulting in lower outage probability.

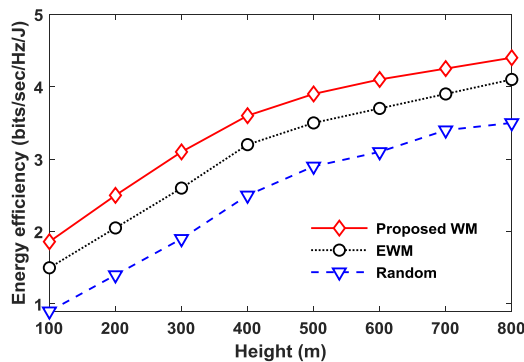


Fig. 5. The effect of hover location on the energy efficiency

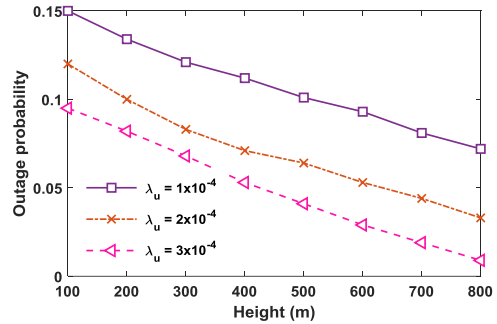


Fig. 6. Outage probability of network for UAV's height.

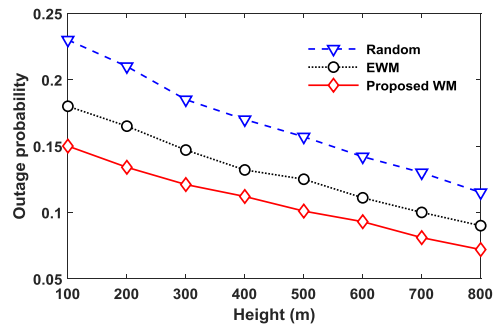


Fig. 7. The effect of hover location on the outage probability

5-5- Effect of Total Transmission Power

Fig. 8 demonstrates the effect of total transmission power on the energy and spectral efficiencies. As the maximum total power increases, the SNR of the links between ABS and ground users increases resulting in spectral efficiency enhancement. It is observed that increasing the total transmission power enhances spectral and energy efficiencies. As transmission power of granted and grant-free users increases, the SNR of links between them and ABS increases resulting in higher spectral efficiency values and lower outage probability values. Increasing the total transmission power increases the total power consumption. On the other side, the increase in spectral efficiency and $(1 - P_{out})$ is higher than the increase in total power consumption since most of the power consumption is related to flight and hover powers of ABS. Hence, increasing the total transmission power enhances energy efficiency.

5-6- Comparing Genetic Algorithm with PSO

Here we compare the performance of the genetic algorithm in finding the flight path of UAV with particle swarm optimization (PSO) and random approach, which selects the following hover location randomly among the possible locations. The flight path only affects the energy efficiency and does not affect spectral efficiency and outage probability since these metrics depend on the hover

location, user pairing, and transmission power and are independent of the flight path. The genetic algorithm, PSO, and random approach performances on the energy efficiency are demonstrated in Fig. 9. As observed genetic algorithm outperforms the PSO and random approaches and has higher energy efficiency.

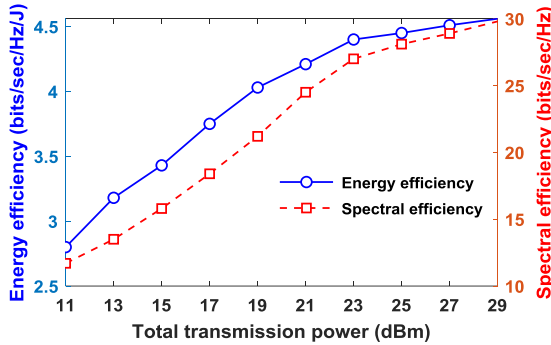


Fig. 8. Effect of total transmission power on energy and spectral efficiencies.

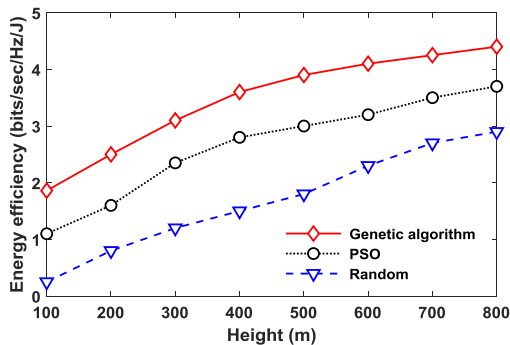


Fig. 9. The effect of the flight path on the energy efficiency

6- Conclusion

In this paper, the challenges of power allocation and NOMA form pairing in the uplink direction of UAV communication systems were investigated. Users in the UAV coverage area were divided into two classes: granted and grant-free, based on prioritizing the type of demands. Granted and grant-free users are respectively considered as strong and weak users in the NOMA pair. The main criterion for the stated challenges has been to maximize energy efficiency. The optimization problem was formulated to maximize the energy efficiency of transmission subject to the constraints on the minimum acceptable spectral efficiency and total transmission power. In order to solve the problem, at first, transmission powers were computed for each possible NOMA pair, and then, the Hungarian algorithm was employed to select the optimum pairs. The flight path of the UAV was modeled as TSP. The results demonstrated that increasing the

height of ABS enhances spectral efficiency, energy efficiency, and outage probability by reducing path loss. Also, increasing the density of users enhances the performance metrics. We also demonstrate that the hover location greatly impacts the performance metrics, and the proposed weighted-mean location outperforms the random and equal-weight men locations.

As future work, we can consider the methods based on machine learning, such as deep belief networks (DBN) for power allocation and pair forming. We can also consider the NOMA clusters with more than two users to support grant-free users in each time slot. Considering different heights for each cluster can be considered as another future work.

Conflicts of Interest

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

References

- [1] R. S. Stansbury, M. A. Vyas, and T. A. Wilson, "A survey of UAS technologies for command, control, and communication (C3)," in *Unmanned Aircraft Systems*: Springer, 2008, pp. 61-78.
- [2] K. P. Valavanis and G. J. Vachtsevanos, *Handbook of unmanned aerial vehicles*. Springer, 2015.
- [3] H. Luo, S.-C. Chu, X. Wu, Z. Wang, and F. Xu, "Traffic collisions early warning aided by small unmanned aerial vehicle companion," *Telecommunication systems*, vol. 75, pp. 169-180, 2020.
- [4] Y. Li *et al.*, "Air-to-ground 3D channel modeling for UAV based on Gauss-Markov mobile model," *AEU-International Journal of Electronics and Communications*, vol. 114, p. 152995, 2020.
- [5] S. Aggarwal and N. Kumar, "Path planning techniques for unmanned aerial vehicles: A review, solutions, and challenges," *Computer Communications*, vol. 149, pp. 270-299, 2020.
- [6] M. Mozaffari, W. Saad, M. Bennis, Y.-H. Nam, and M. Debbah, "A tutorial on UAVs for wireless networks: Applications, challenges, and open problems," *IEEE Communications Surveys & Tutorials*, vol. 21, no. 3, pp. 2334-2360, 2019.
- [7] J. H. Sarker and A. M. Nahhas, "A secure wireless mission critical networking system for unmanned aerial vehicle communications," *Telecommunication Systems*, vol. 69, no. 2, pp. 237-251, 2018.
- [8] S. Sudhakar, V. Vijayakumar, C. S. Kumar, V. Priya, L. Ravi, and V. Subramaniaswamy, "Unmanned Aerial Vehicle (UAV) based Forest Fire Detection and monitoring for reducing false alarms in forest-fires," *Computer Communications*, vol. 149, pp. 1-16, 2020.
- [9] Q. Liu *et al.*, "Joint power and time allocation in energy harvesting of UAV operating system," *Computer Communications*, vol. 150, pp. 811-817, 2020.
- [10] Y. Zeng, J. Xu, and R. Zhang, "Energy minimization for wireless communication with rotary-wing UAV," *IEEE Transactions on Wireless Communications*, vol. 18, no. 4, pp. 2329-2345, 2019.
- [11] Y. Zeng and R. Zhang, "Energy-efficient UAV communication with trajectory optimization," *IEEE Transactions on Wireless Communications*, vol. 16, no. 6, pp. 3747-3760, 2017.
- [12] M. Hua, Y. Wang, C. Li, Y. Huang, and L. Yang, "Energy-efficient optimization for UAV-aided cellular offloading," *IEEE Wireless Communications Letters*, vol. 8, no. 3, pp. 769-772, 2019.
- [13] J. Yu, R. Zhang, Y. Gao, and L.-L. Yang, "Modularity-based dynamic clustering for energy efficient UAVs-aided communications," *IEEE Wireless Communications Letters*, vol. 7, no. 5, pp. 728-731, 2018.
- [14] Y. Cai, Z. Wei, R. Li, D. W. K. Ng, and J. Yuan, "Energy-efficient resource allocation for secure UAV communication systems," in *2019 IEEE Wireless Communications and Networking Conference (WCNC)*, 2019: IEEE, pp. 1-8.
- [15] M.-N. Nguyen, L. D. Nguyen, T. Q. Duong, and H. D. Tuan, "Real-time optimal resource allocation for embedded UAV communication systems," *IEEE Wireless Communications Letters*, vol. 8, no. 1, pp. 225-228, 2018.
- [16] Y. Chen, W. Feng, and G. Zheng, "Optimum placement of UAV as relays," *IEEE Communications Letters*, vol. 22, no. 2, pp. 248-251, 2017.
- [17] M. Mozaffari, W. Saad, M. Bennis, and M. Debbah, "Wireless communication using unmanned aerial vehicles (UAVs): Optimal transport theory for hover time optimization," *IEEE Transactions on Wireless Communications*, vol. 16, no. 12, pp. 8052-8066, 2017.
- [18] M. Alzenad, A. El-Keyi, F. Lagum, and H. Yanikomeroglu, "3-D placement of an unmanned aerial vehicle base station (UAV-BS) for energy-efficient maximal coverage," *IEEE Wireless Communications Letters*, vol. 6, no. 4, pp. 434-437, 2017.
- [19] R. Amorim *et al.*, "Measured uplink interference caused by aerial vehicles in LTE cellular networks," *IEEE Wireless Communications Letters*, vol. 7, no. 6, pp. 958-961, 2018.
- [20] J. Gong, T.-H. Chang, C. Shen, and X. Chen, "Flight time minimization of UAV for data collection over wireless sensor networks," *IEEE Journal on Selected Areas in Communications*, vol. 36, no. 9, pp. 1942-1954, 2018.
- [21] D. Yang, Q. Wu, Y. Zeng, and R. Zhang, "Energy tradeoff in ground-to-UAV communication via trajectory design," *IEEE Transactions on Vehicular Technology*, vol. 67, no. 7, pp. 6721-6726, 2018.
- [22] X. Liu, M. Chen, and C. Yin, "Optimized trajectory design in UAV based cellular networks for 3D users: A double Q-learning approach," 2019.
- [23] D. Zhai, R. Zhang, L. Cai, B. Li, and Y. Jiang, "Energy-efficient user scheduling and power allocation for NOMA-based wireless networks with massive IoT devices," *IEEE Internet of Things Journal*, vol. 5, no. 3, pp. 1857-1868, 2018.
- [24] M. S. Ali, H. Tabassum, and E. Hossain, "Dynamic user clustering and power allocation for uplink and downlink non-orthogonal multiple access (NOMA) systems," *IEEE access*, vol. 4, pp. 6325-6343, 2016.
- [25] J. Cui, Z. Ding, P. Fan, and N. Al-Dhahir, "Unsupervised machine learning-based user clustering in millimeter-wave-NOMA systems," *IEEE Transactions on Wireless Communications*, vol. 17, no. 11, pp. 7425-7440, 2018.
- [26] M. Zeng, A. Yadav, O. A. Dobre, and H. V. Poor, "Energy-efficient joint user-RB association and power allocation for uplink hybrid NOMA-OMA," *IEEE Internet of Things Journal*, vol. 6, no. 3, pp. 5119-5131, 2019.
- [27] S. Dhakal, P. A. Martin, and P. J. Smith, "NOMA with guaranteed weak user QoS: design and analysis," *IEEE Access*, vol. 7, pp. 32884-32896, 2019.
- [28] X. Mu, Y. Liu, L. Guo, and J. Lin, "Uplink Non-Orthogonal Multiple Access for UAV Communications," *CoRR*, 2019.
- [29] W. Mei and R. Zhang, "Uplink cooperative NOMA for cellular-connected UAV," *IEEE Journal of Selected Topics in Signal Processing*, vol. 13, no. 3, pp. 644-656, 2019.
- [30] R. Duan, J. Wang, C. Jiang, H. Yao, Y. Ren, and Y. Qian, "Resource allocation for multi-UAV aided IoT NOMA uplink transmission systems," *IEEE Internet of Things Journal*, vol. 6, no. 4, pp. 7025-7037, 2019.
- [31] Y. Liu, Z. Qin, Y. Cai, Y. Gao, G. Y. Li, and A. Nallanathan, "UAV communications based on non-

- orthogonal multiple access," *IEEE Wireless Communications*, vol. 26, no. 1, pp. 52-57, 2019.
- [32] M. Yang, B. Li, Z. Bai, and Z. Yan, "SGMA: Semi-granted multiple access for non-orthogonal multiple access (NOMA) in 5G networking," *Journal of Network and Computer Applications*, vol. 112, pp. 115-125, 2018.
- [33] Z. Ding, R. Schober, P. Fan, and H. V. Poor, "Simple semi-grant-free transmission strategies assisted by non-orthogonal multiple access," *IEEE Transactions on Communications*, vol. 67, no. 6, pp. 4464-4478, 2019.
- [34] Q. Zhang, M. Mozaffari, W. Saad, M. Bennis, and M. Debbah, "Machine learning for predictive on-demand deployment of UAVs for wireless communications," in *2018 IEEE Global Communications Conference (GLOBECOM)*, 2018: IEEE, pp. 1-6.
- [35] M. Mozaffari, W. Saad, M. Bennis, and M. Debbah, "Optimal transport theory for power-efficient deployment of unmanned aerial vehicles," in *2016 IEEE international conference on communications (ICC)*, 2016: IEEE, pp. 1-6.
- [36] H. Tabassum, M. S. Ali, E. Hossain, M. J. Hossain, and D. I. Kim, "Uplink vs. downlink NOMA in cellular networks: Challenges and research directions," in *2017 IEEE 85th vehicular technology conference (VTC Spring)*, 2017: IEEE, pp. 1-7.
- [37] D.-T. Do and M.-S. Van Nguyen, "Outage probability and ergodic capacity analysis of uplink NOMA cellular network with and without interference from D2D pair," *Physical Communication*, vol. 37, p. 100898, 2019.
- [38] A. Hussain, Y. S. Muhammad, M. Nauman Sajid, I. Hussain, A. Mohamd Shoukry, and S. Gani, "Genetic algorithm for traveling salesman problem with modified cycle crossover operator," *Computational intelligence and neuroscience*, vol. 2017, 2017.