

# JOURNAL OF TELECOMMUNICATIONS AND INFORMATION TECHNOLOGY

2 / 2025

vol. 100

**Using a Half-mode SIW Loaded with Slots to Realize a Compact UWB Bandpass Filter**

*H. Al-Jeshami, H. Al-Saedi, M.F. Hasan, H.I. Khani, M.Y. Muhsin, and W.M. Abdel-Wahab* 1

**An Artificial Intelligence-based Handover Triggering and Management Mechanism for 5G Ultra-dense Networks to Improve Handover Authentication**

*P. Rajesh, A. Vijaya Lakshmi, and Ebenezer Abishek B.* 9

**Synthesizing Wide-beam Array Patterns Using Phase-only Control and Trapezoidal Amplitudes for Satellite-based Internet Access**

*Z. Turki Hassan and J.R. Mohammed* 21

**Analysis of Pyramidal Microwave Absorbers for Enhanced Performance in 1–10 GHz Frequency Range**

*A.R. Thanoon and K.H. Sayidmarie* 29

**UAV-BS-based Hybrid OMA-NOMA System with Multiple Antennas for Multi-user Communication**

*A.Y. Sadeeq and M.A. Ahmed* 38

**Using Modified Gorti-enhanced Homomorphic Cryptosystem to Improve Security of ECG Signal**

*F.Z. Besmi, S. Belkacem, and N. Messaoudi* 46

**ILP Optimized LSTM-based Autoscaling and Scheduling of Containers in Edge-cloud Environment**

*S. Singh, Narayan D.G., S. Mujawar, G.S. Hanchinamani, and P.S. Hiremath* 56

**TinyML-driven Sensor Nodes for Energy-efficient Acoustic Event Detection in Pervasive Acoustic WSNs**

*B.B. Roy, S. Das, and U. Kr. Mondal* 69

*(Contents continued on back cover)*

## *Editor-in-Chief*

**Adrian Kliks**, Poznan University of Technology, Poland

## *Editorial Advisory Board*

**Hovik Baghdasaryan**, National Polytechnic University of Armenia, Armenia

**Naveen Chilamkurti**, LaTrobe University, Australia

**Luis M. Correia**, Instituto Superior Técnico, Universidade de Lisboa, Portugal

**Pedro Crespo Bofill**, Universidad de Navarra, Spain

**Luca De Nardis**, DIET Department, University of Rome La Sapienza, Italy

**Nikolaos Dimitriou**, NCSR "Demokritos" Athens, Greece

**Ciprian Dobre**, Politechnic University of Bucharest, Romania

**Piotr Gawrysiak**, Warsaw University of Technology, Poland

**Filip Idzikowski**, Poznan University of Technology, Poland

**Andrzej Jajszczyk**, AGH University of Science and Technology, Poland

**Zbigniew Jaroszewicz**, National Institute of Telecommunications, Poland

**Albert Levi**, Sabanci University, Turkey

**Marian Marciniak**, National Institute of Telecommunications, Poland

**George Mastorakis**, Technological Educational Institute of Crete, Greece

**Constandinos Mavromoustakis**, University of Nicosia, Cyprus

**Takumi Miyoshi**, Shibaura Institute of Technology, Japan

**Klaus Mößner**, Technische Universität Chemnitz, Germany

**Imran Muhammad**, King Saud University, Saudi Arabia

**Mjumo Mzyece**, University of the Witwatersrand, South Africa

**Daniel Negru**, University of Bordeaux, France

**Jordi Perez-Romero**, UPC, Spain

**Michał Pióro**, Warsaw University of Technology, Poland

**Konstantinos Psannis**, University of Macedonia, Greece

**Salvatore Signorello**, University of Lisboa, Portugal

**Adam Wolisz**, Technische Universität Berlin, Germany

**Tadeusz A. Wysocki**, University of Nebraska, USA

## *Editorial Team*

Content Editor: **Robert Magdziak**

Managing Editor: **Ewa Kapuściarek**

eISSN 1899-8852

© Copyright by National Institute of Telecommunications, Poland 2025

# Using a Half-mode SIW Loaded with Slots to Realize a Compact UWB Bandpass Filter

Hussein Al-Jeshami<sup>1</sup>, Hussam Al-Saedi<sup>1</sup>, Mohammed F. Hasan<sup>1</sup>, Halah I. Khani<sup>1</sup>,  
Muhannad Y. Muhsin<sup>1</sup>, and Wael M. Abdel-Wahab<sup>2</sup>

<sup>1</sup>University of Technology – Iraq, Baghdad, Iraq,

<sup>2</sup>University of Waterloo, Waterloo, Canada

<https://doi.org/10.26636/jtit.2025.2.2085>

**Abstract** – This paper presents an ultra-wideband bandpass filter (UWB-BPF) based on a half-mode substrate-integrated waveguide (HMSIW) structure with a passband that spans from 4.12 to 11 GHz and has a center frequency of 7.56 GHz with a fractional bandwidth of 91%. Less than 0.3 dB of insertion loss was achieved through the passband, with the return loss being better than 19 dB. Two transmission zeros were generated, contributing to the 5.24 GHz out-of-band rejection. The proposed filter is compact, making it a good choice for space-constrained systems. In addition, a notch response has been utilized to comply with standards issued by applicable regulatory bodies, by blocking some bands located in the passband. Ansys EDT software was used to implement the design and determine the response of the proposed structure, while Keysight ADS was used to validate the results.

**Keywords** – compact filter, half-mode SIW, notch response, slots loaded, ultra-wideband

## 1. Introduction

Ultra-wideband technology (UWB) represents a remarkable step in the development of wireless communication systems. Due to its fractional bandwidth (FBW) of more than 20% [1], it has a wide range of applications and is used, inter alia, in medical imaging, ground penetration radars, surveillance, and portable electronic devices. The importance of UWB devices comes from their role in enhancing overall performance and, thus, increasing the reliability of wireless communication. The latter may be achieved mainly by minimizing interference to maintain integrity of signals.

Many papers focusing on UWB technology have been written. As far as filter design is concerned, the following works are worth mentioning as significant. Multiple UWB filters were discussed in [2], while a single notch UWB bandpass filter (UWB-BPF) was created in [3], [4] by using modified substrate-integrated waveguide (SIW) in combination with a dual-split square complementary split ring resonator (CSRR).

A miniaturized frequency-selective surface with intricate square-shaped rings and grids was presented in [5] to realize a dual-band BPF (DB-BPF). Although it has a compact

size, the high complexity of using cascaded layers is a significant drawback due to the manufacturing tolerances that are achievable at present.

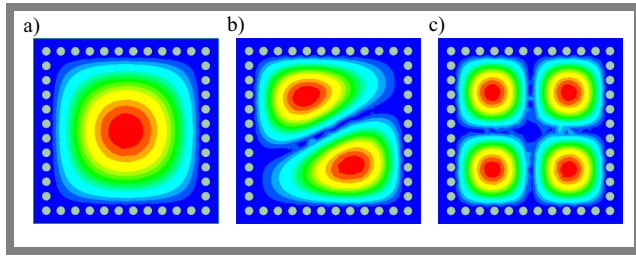
A half-mode SIW-based (HMSIW) resonator loaded with square-shaped slot was presented in [6]. It has a simple structure with a fairly small footprint. However, it provides narrow upper band rejection. In article [7], another SIW-based resonator loaded with U-shaped slots was proposed to illustrate the BPF response. The filter has two notches generated by further loading the resonator with lumped elements (capacitors). This filter is considered large and the design is characterized by poor impedance matching within the operating passband.

In [8], researchers presented an UWB-BPF component using a butterfly-shaped unit cell to load the coplanar waveguide technology. Additionally, T-resonators were added to the structure to obtain a notch response within the passband region. An UWB-BPF using open-stub resonators with stepped impedance is discussed in [9].

The resonators take advantage of using multimode quarter-wavelength series transformers. A defected ground structure (DGS) in a bowtie-shaped pattern and a grid of multi slots were engraved on the top layers of an SIW resonator to produce an UWB-BPF response in [10], while a combination of spoof surface plasmon polaritons (SSPP) and HM-SIW were used to form UWB-BPF response for Ka- and X-band applications in papers [11], [12].

A quasi-lumped SIW cavity was proposed in [13], despite its high-profile structure, because many SMD capacitors were integrated on the top layer to reduce the overall size. In [14], the researchers realized a hybrid fractal slot on the top of a HM-SIW resonator with the help of a square-ring CSRR unit cell, such as DGS, to create a BPF structure, while fractal-based miniaturization with meander line slots was applied on an SIW cavity in [15] to implement a wideband planar BPF.

Field perturbation achieved by applying upper surface slots in a triangular-shaped SIW cavity was carried out in [16] to obtain UWB-BPF response. Lastly, a single sector of a 2.2-degree SIW-based cavity was shown in [17] to deliver a dual-mode UWB response.



**Fig. 1.** Distribution of fields in an SIW cavity: a) fundamental mode, b) second mode, and c) fourth mode.

This paper presents a four-pole planar HMSIW filter with UWB-BPF response, where three slots were embedded on the top of the resonator to apply the miniaturization effect to the design. A notch response (NR) was implemented within the passband.

## 2. Filter Design

Field distribution is a crucial factor in the design of any SIW filter. The distribution of fields is affected by the mode in which the resonator is being excited at. Typically, an SIW cavity is constructed with the use of walls of metallic vias through which the traveling electromagnetic wave is guided and confined within the used substrate. The modes at which the SIW resonator works are transverse electric modes (TE). Ansys EDT simulator software was used to carry out the design and evaluation processes. Figure 1 shows the distribution of the electric field within an SIW cavity model.

Primarily, to design an SIW resonator, the following equations from [18], [19] can be used:

$$p \leq 2d, \tag{1}$$

$$0.05 < \frac{p}{\lambda_c} < 0.25, \tag{2}$$

$$d < \frac{\lambda_g}{5}, \tag{3}$$

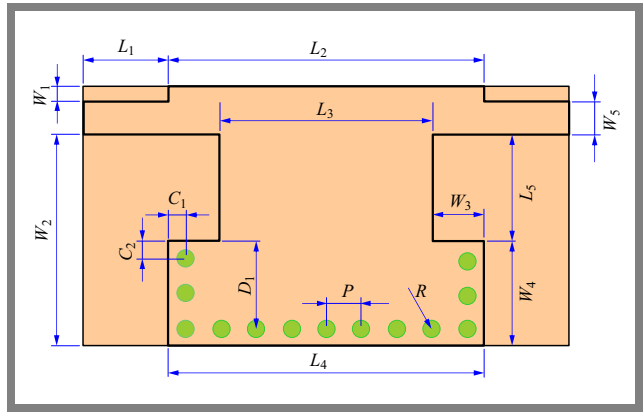
$$f_{c(TE_{mn})} = \frac{c}{2\pi} \sqrt{\left(\frac{\pi m}{W}\right)^2 + \left(\frac{\pi n}{L}\right)^2}, \tag{4}$$

$$W = W_{siw} - \frac{d^2}{0.95 p}, \tag{5}$$

$$L = L_{siw} - \frac{d^2}{0.95 p}. \tag{6}$$

In Eqs. (1)–(6),  $p$  is the pitch between two vias (measured from the center of the vias),  $d$  is the diameter of each via,  $\lambda_g$  is the guided wavelength,  $m$  and  $n$  are mode numbers with respect to the TE mode. Parameter  $W$  stands for the effective width,  $W_{SIW}$  is the width of the SIW,  $L$  is the effective length, and  $L_{SIW}$  is the length of the SIW.

To miniaturize the SIW structure, a half-mode technique has been used, where the width of the HMSIW cavity becomes  $\frac{W}{2}$  [20]. In this filter design, the Rogers RO4003 substrate with relative permittivity  $\epsilon_r = 3.55$ , loss tangent  $\tan \delta = 0.0027$ ,



**Fig. 2.** Initial filter design.

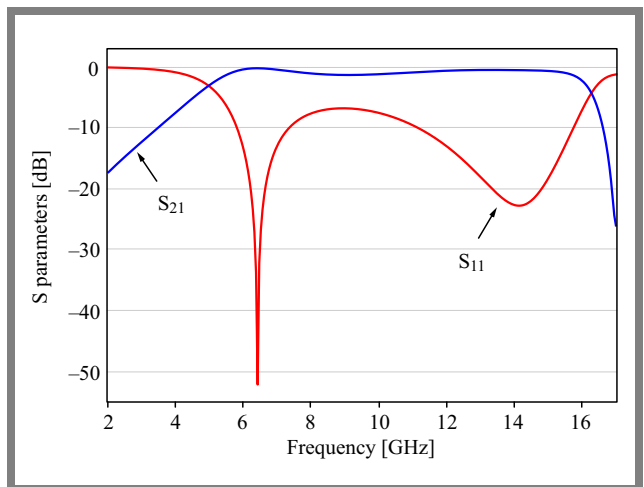
**Tab. 1.** Dimensions of the initial design (all in mm).

Desc.	Value	Desc.	Value	Desc.	Value
$L_1$	3.8	$W_1$	0.5	$R$	0.3
$L_2$	10.8	$W_2$	7.227	$P$	1.2
$L_3$	7.3	$W_3$	1.7	$D_1$	2.9
$L_4$	10.8	$W_4$	3.5	$C_1$	0.6
$L_5$	3.465	$W_5$	1.13	$C_2$	0.6

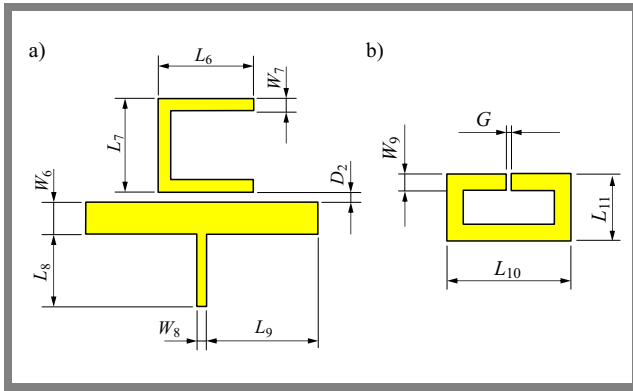
and height  $h = 0.508$  mm was used. The initial structure is based on Eqs. (1)–(6) and is shown in Fig. 2. The dimensions are listed in Tab. 1.

The response of the initial design is shown in Fig. 3. It can be noticed that the lower cut-off frequency  $f_1$  is 5.1 GHz, while the upper cut-off frequency  $f_2$  is located at 15.5 GHz. Furthermore, the passband response suffers from rippling, which degrades the performance of the filter. In other words, the filter shows a higher insertion loss (IL) in a specific part of the passband region. Additionally, the signal reflection at the ports is too high due to the band performance being driven by the  $S_{11}$  response.

To further miniaturize and improve the proposed filter, resonators were loaded onto the top layer of the SIW with complementary T-shaped and U-shaped slots, as illustrated in



**Fig. 3.**  $S_{11}$  and  $S_{21}$  parameters of the initial design.



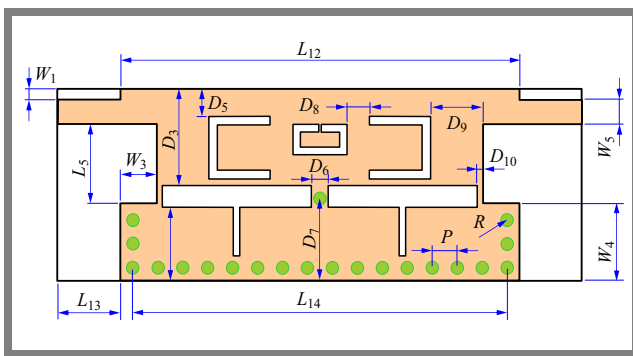
**Fig. 4.** Modification of the slots and dimensions of the design.

**Tab. 2.** Dimensions of the coupling slots (all in mm).

Desc.	Value	Desc.	Value	Desc.	Value
$W_6$	1	$L_6$	2.95	$L_{10}$	2.6
$W_7$	0.4	$L_7$	2.9	$L_{11}$	1.4
$W_8$	0.3	$L_8$	2.25	$D_2$	0.3
$W_9$	0.35	$L_9$	3.45	$G_1$	0.1

Fig. 4a. Furthermore, to increase the order of the filter and enhance its response, two HMSIW resonators were used, and a coupling complementary open-ring slot was added between them. The coupling slot with its parameter notations is shown in Fig. 4b, while the parameters of the slots are listed in Tab. 2.

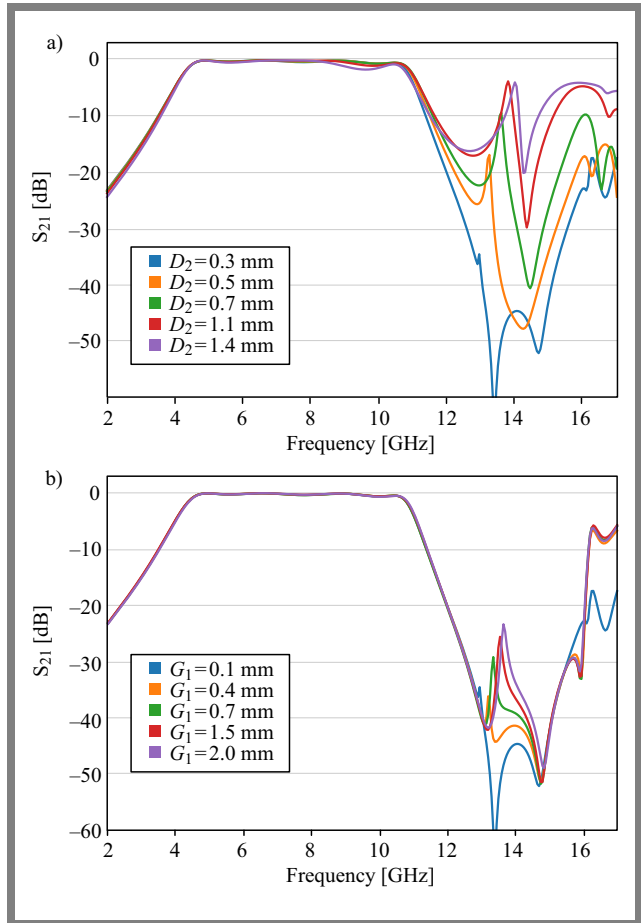
Such a modified UWB-BPF filter structure is depicted in Fig. 5, with its dimensions presented in Tab. 3.



**Fig. 5.** Improved UWB-BPF.

**Tab. 3.** Dimensions of the proposed UWB-BPF (all in mm).

Desc.	Value	Desc.	Value	Desc.	Value
$L_{12}$	19.2	$D_4$	2.95	$D_8$	2.6
$L_{13}$	3.8	$D_5$	1.26	$D_9$	2.45
$L_{14}$	18	$D_6$	1	$D_{10}$	0.2
$D_3$	4.465	$D_7$	3.8		



**Fig. 6.**  $S_{21}$  versus frequency for various  $D_2$  and  $G_1$  values.

### 3. Design Analysis

To achieve the best performance of the proposed filter, a set of parametric studies were performed. This research helps to find the design parameters that have the greatest effect on the filter's response. Mainly, parameters  $D_2$ ,  $G_1$ ,  $L_7$ ,  $L_9$ ,  $D_2$  and  $W_4$  are the ones with significant impact on the BPF response. It can be noticed that changes in  $D_2$  and  $G_1$  improve the out-of-band rejection (OBR) response, which is very important to suppress resonance harmonics. However, it is clear that  $D_2$  has the strongest impact on OBR, as illustrated in Fig. 6.

Then parameter  $L_7$  was varied, influencing the location of the upper band cutoff frequency (UBCF). As shown in Fig. 7, an increase in  $L_7$  results in shifting UBCF to lower values, without impact on the lower band cutoff frequency (LBCF). In other words, the bandwidth of this filter can be modified without changing LBCF.

The next analysis proves that changing  $L_9$  affects the location of LBCF, while UBCF remains the same. It can be concluded that an increase in  $L_9$  lowers LBCF, which allows to minimize the filter's size and makes it more suitable for use in space constrained devices that require the filter to work with lower frequencies. Figure 8 confirms that, opposite to the impact of  $L_7$ , an increase in  $L_9$  also leads to a higher bandwidth response.

In conjunction with the reflection coefficient at the designated ports of the proposed filter, it can be determined that parameter  $W_4$  has the biggest impact on  $S_{11}$  response, as presented in Fig. 9.

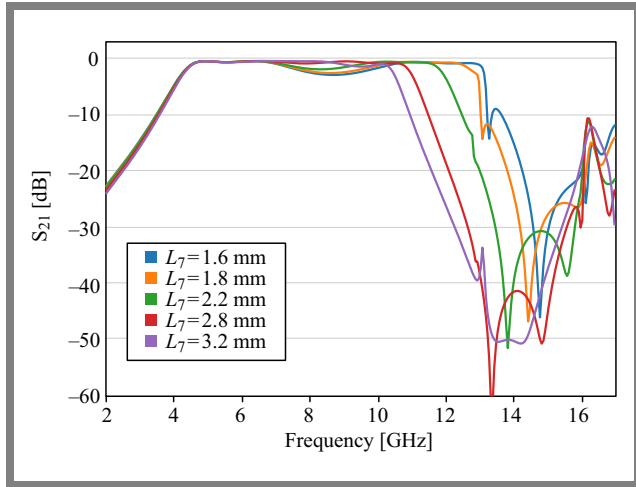


Fig. 7. Impact on  $S_{21}$  caused by variation of  $L_7$  parameter.

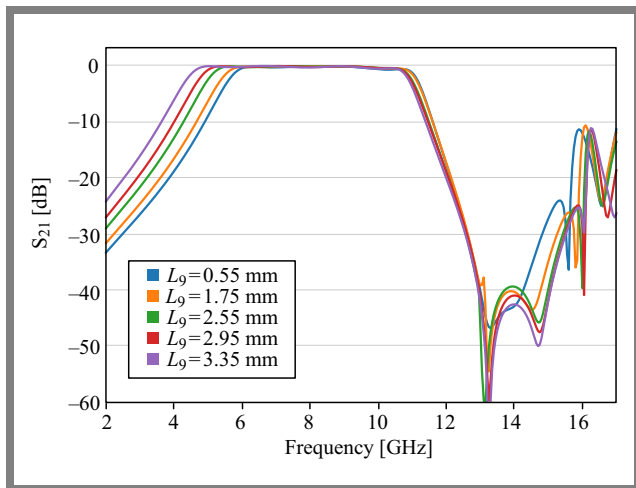


Fig. 8. Study of the  $L_9$  parameter.

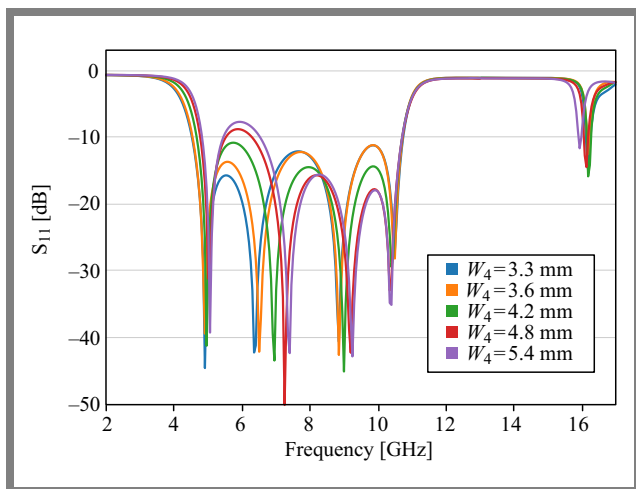


Fig. 9.  $S_{11}$  parameter versus frequency with variable  $W_4$ .

## 4. Design Optimization

To enhance filter response, a search-based multi-objective optimization technique available in the Ansys EDT simulator was utilized. The criteria were set to have the minimum possible insertion loss (IL) with the return loss (RL) greater than 20 dB across the passband.

Such a condition aims to minimize power losses suffered while passing through the filter and to diminish the reflection at the ports of the proposed filter for the entire required passband.  $L_1, L_5, W_3, W_4, D_1, D_6, D_8, D_9$  and  $D_{10}$  parameters were used in the optimization.

The simulation findings of the proposed BPF are shown in Fig. 10, where  $f_1 = 4.12$  GHz and  $f_2 = 11$  GHz. IL is found to be better than 0.3 dB with RL over 19 dB. Furthermore, two transmission zeros were realized at 13.66 GHz and 14.68 GHz. This helps to achieve an improved wideband rejection behavior, where it spans for 5.24 GHz after  $f_2$  with rejection of up to 50 dB.

The proposed UWB-BPF has a center frequency (CF) of 7.56 GHz and FBW of 91%. The filter size is calculated to be  $18 \times 8.265$  mm, or  $0.348 \times 0.758 \lambda_g$ , where  $\lambda_g$  is the guided wave length at the filter CF.

### 4.1. In-band Rejection Response

For any UWB-BF, it is important to have an NR in the passband (or band-stop behavior). This is a necessary step to suppress some signals that lie within bands that are required to be blocked due to regulatory compliance.

To produce NR, a modification must be applied to the proposed filter by loading the excitation ends of the filter with L-shaped shunt resonators, where they were integrated within the inset slots to maintain the same proposed filter size. The modified filter design structure is demonstrated in Fig. 11 and its characteristic dimensions are listed in Tab. 4.

A parametric study was conducted to examine the design parameter that has the most impact on NR. Figures 12-13 illustrate the variation in width  $W_S$  and length  $L_S$  of the L-shaped shunt resonators. One may notice that increasing

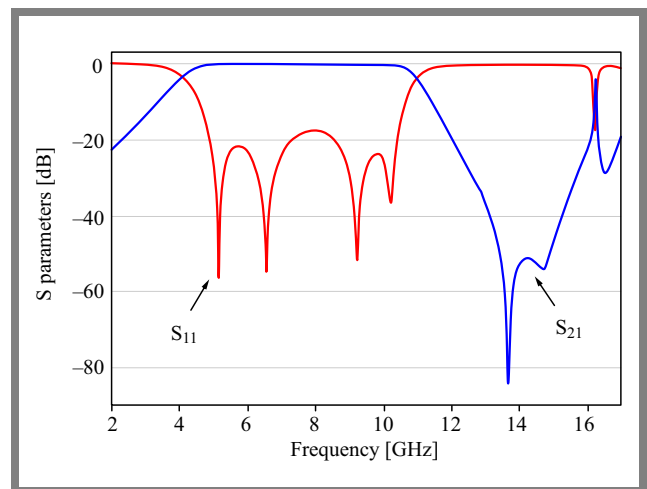


Fig. 10. Response of the S parameters of the proposed UWB-BPF.

$W_S$  and  $L_S$  leads to moving NR to lower frequencies, with almost no effect on any other regions in the filter's response. However,  $W_S$  variations tend to be more effective in shifting the location of NR to frequencies lower than that of  $L_S$ .

Frequency response of the proposed filter is presented in Fig. 14. By comparing the results shown in Figs. 10 and 14, one may conclude that the response of the modified UWB-BPF matches that of the originally proposed UWB-BPF, except at the location where the notch was applied.

In the next step, the Ansys EDT simulator was utilized to plot the electric field (EF) of selected frequencies. EFs at 2 GHz and 13.35 GHz, where no power was transferred through the filter, are illustrated in Figs. 15a and 15d, respectively. For the in-band region, EFs at 6.35 GHz and 10.5 GHz are shown in Figs. 15b-c.

Moreover, Fig. 15e illustrates the EF at which the notch was generated, where the notch is trapped in the notch structure and no power was transferred between the input and output ports.

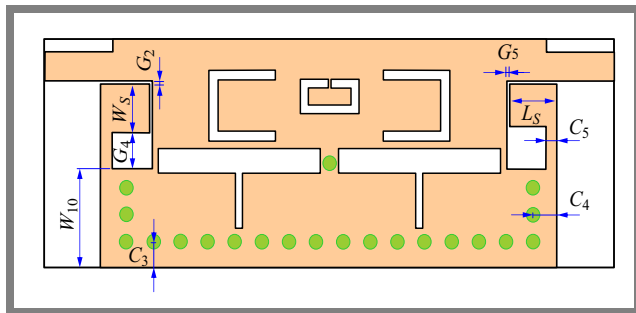


Fig. 11. Modified UWB-BPF structure.

Tab. 4. Characteristic dimensions of the modified UWB-BPF (all in mm).

Desc.	Value	Desc.	Value	Desc.	Value
$G_2$	0.1	$C_3$	0.6	$W_{10}$	3.5
$G_3$	0.1	$C_4$	0.8	$L_S$	1.8
$G_4$	3.265	$C_5$	0.2	$W_S$	0.2

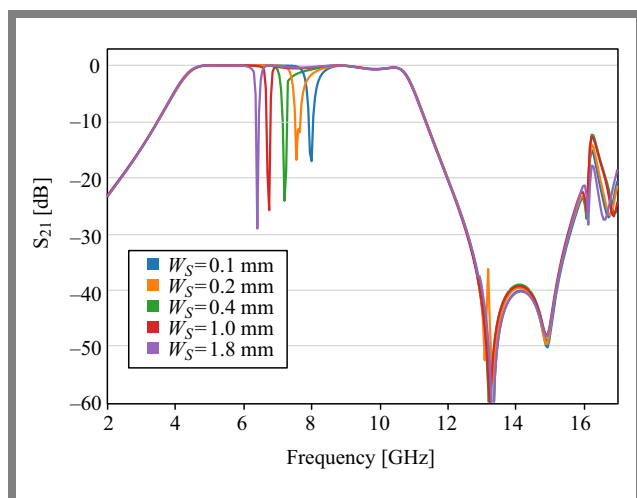


Fig. 12. Parametric study concerning  $W_S$ .

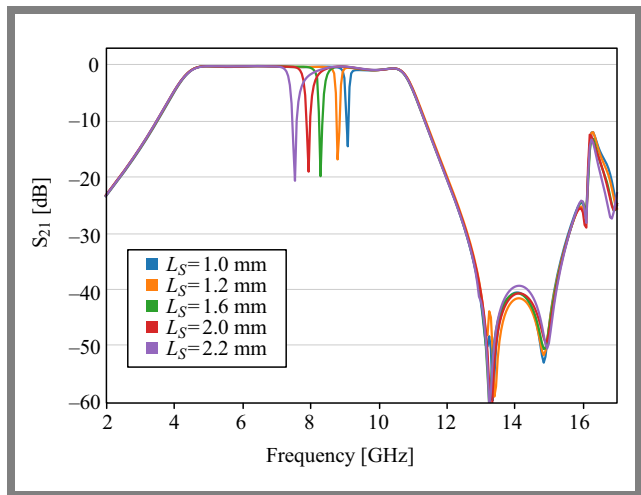


Fig. 13. Parametric study concerning the  $L_S$  parameter.

## 5. Experimental Evaluation and Verification

To validate the Ansys EDT simulated UWB-BPF responses, the proposed filter and its modified version were also simulated in Keysight ADS software (a 2.5D momentum EM simulator). The responses of the filters were in alignment with the findings obtained using Ansys EDT. Figure 16 illustrates the responses of the proposed UWB-BPF obtained using Ansys EDT and Keysight ADS.

It may be concluded that the very small differences between the two findings (Figs. 10, 14 and Figs. 16, 17) come from how the two EM simulators calculated the designs and the methods that were being used to solve the proposed structures. Being a full-wave simulator, Ansys EDT uses the finite element method (FEM), whereas Keysight ADS relies on the method of momentum (MoM) to solve the EM structures, hence the variations shown above.

Table 5 provides a detailed comparison between the most recent published work and the proposed filter. One may summarize that the designed filter structure achieved IL and RL parameters that were better than those obtained by

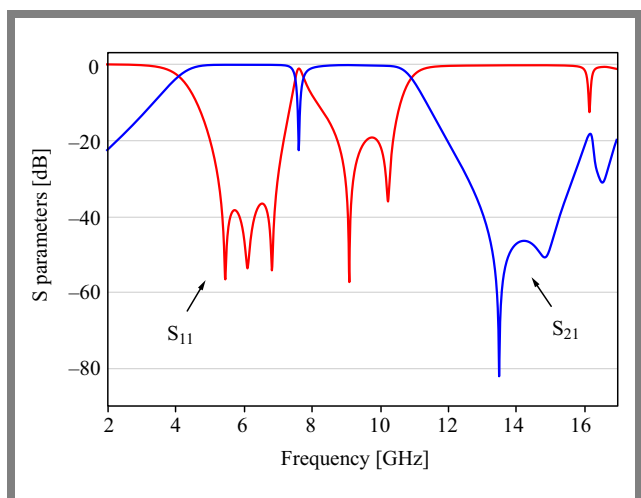
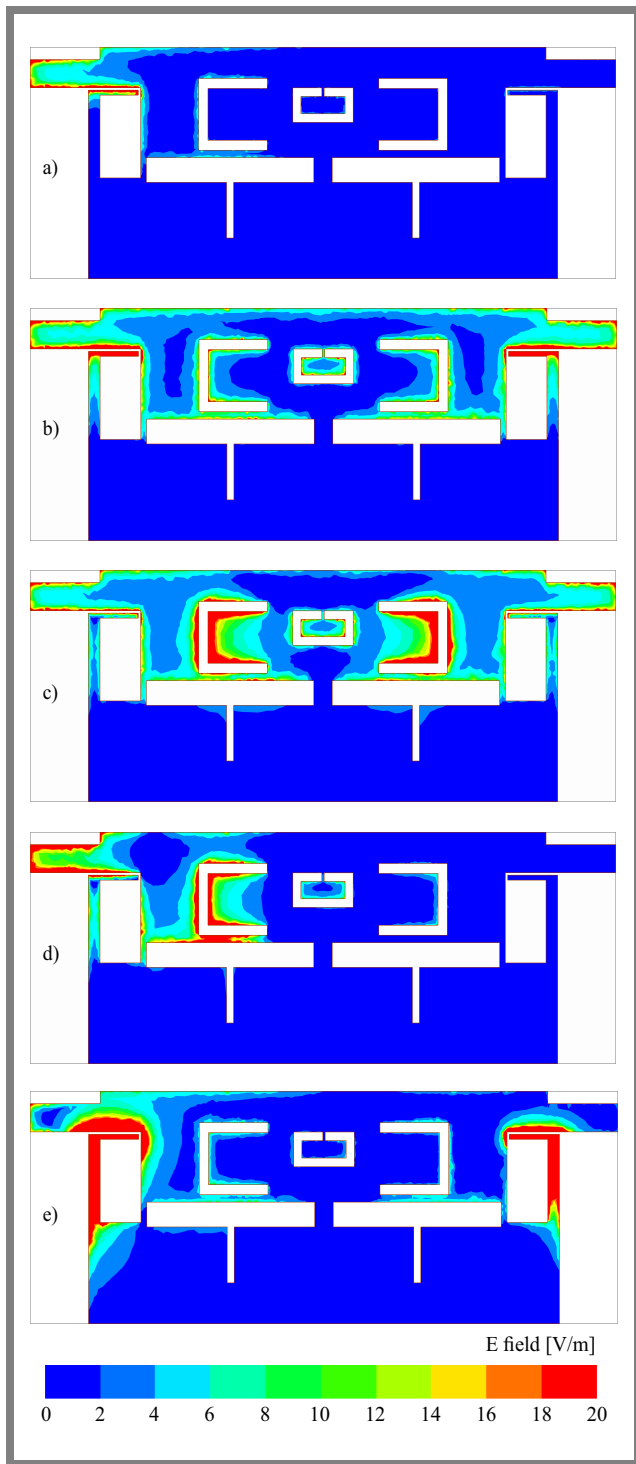


Fig. 14. Response of the S parameters of the modified UWB-BPF.

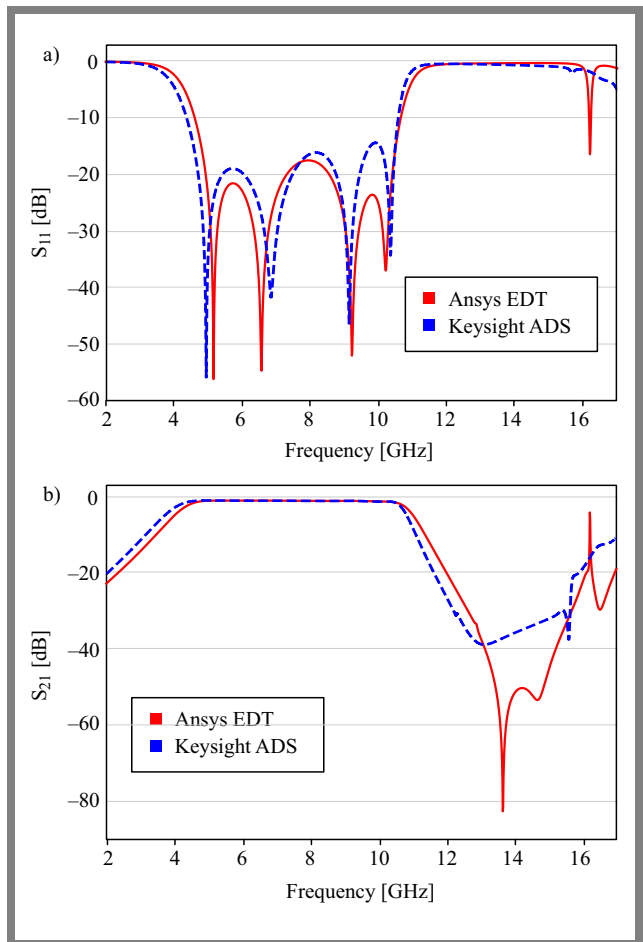


**Fig. 15.** Electric field at: a) 2 GHz, b) 6.35 GHz, c) 10.50 GHz d) 13.35 GHz, and e) 7.6 GHz.

most of its counterparts. Furthermore, the proposed filter is characterized by a smaller size in comparison with [3], [4], [6], but is larger when compared with [21].

### 6. Conclusions

In this paper, a UWB-BPF based on HMSIW and loaded with three different-shaped slots was proposed. The filter operates



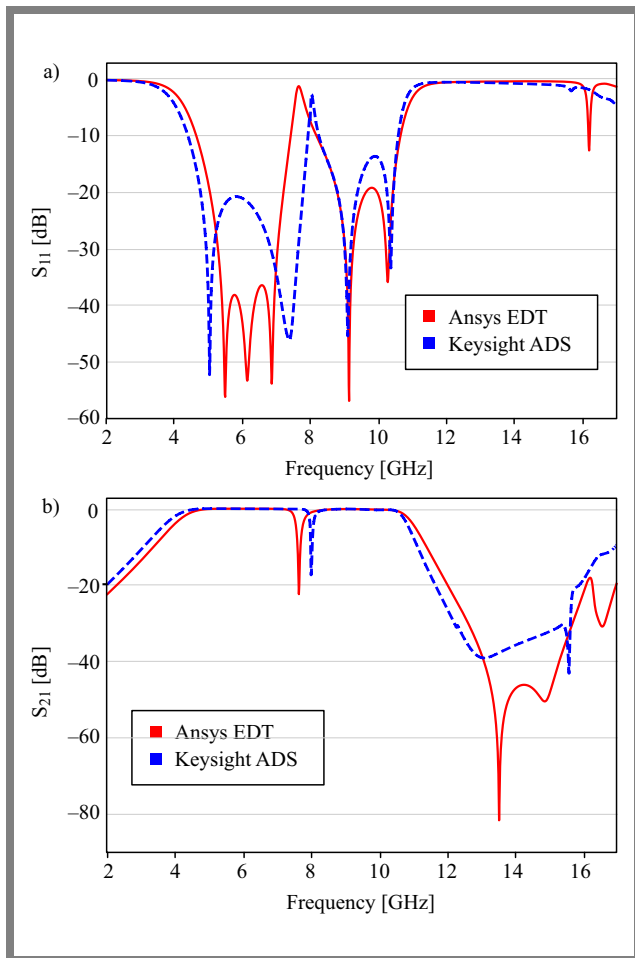
**Fig. 16.** S parameters of the proposed UWB-BPF obtained with Ansys EDT and Keysight ADS simulators: a)  $S_{11}$  and b)  $S_{21}$ .

in the frequency band between 4.12 GHz and a 11 GHz, with FBW of 91% and a low IL level of 0.3 dB at the CF. Its RL is better than 19 dB. The structure was designed with a small filter area of  $0.26 \lambda_g^2$  making it suitable for modern compact systems.

Moreover, integrated L-shaped slots were included within the same filter area to introduce an NR through the passband region. Compared with filters described in works published earlier, the proposed solution achieved a stronger overall balance of bandwidth, loss, and size.

**Tab. 5.** Comparison with other works.

Ref.	CF [GHz]	FBW [%]	IL [%]	RL [%]	Size [ $\lambda_g$ ]
[3]	6.6	106	< 1.8	> 10	0.37
[4]	6.7	107	< 1.6	> 15	0.37
[6]	11.65	119	< 0.3	> 20	0.29
[21]	6.5	92	< 1.74	> 15	0.24
This paper	7.56	91	< 0.3	> 19	0.26



**Fig. 17.** S parameters of the modified UWB-BPF obtained with Ansys EDT and Keysight ADS simulators: a)  $S_{11}$  and b)  $S_{21}$ .

## References

- [1] Code of Federal Regulations, *Part 15: Radio Frequency Devices*, 2023 [Online]. Available: <https://www.ecfr.gov/current/title-47/part-15>.
- [2] R.P. Verma and B. Sahu, "Structure and Performance Comparison of Ultra-wideband Bandpass Filter: Review Article", *2021 7th International Conference on Signal Processing and Communication (ICSC)*, Noida, India, 2021 (<https://doi.org/10.1109/ICSC53193.2021.9673323>).
- [3] S. Udhayanan and K. Shambavi, "Compact Single Notch UWB Bandpass Filter with Metamaterial and SIW Technique", *PIER Letters*, vol. 117, pp. 41–46, 2024 (<https://doi.org/10.2528/PIERL23113004>).
- [4] S. Udhayanan and K. Shambavi, "Metamaterial-based Compact UWB Bandpass Filter Using Substrate Integrated Waveguide", *PIER Letters*, vol. 120, pp. 1–6, 2024 (<https://doi.org/10.2528/PIERL24031107>).
- [5] Z. Li *et al.*, "A Miniaturized Ultra-wideband Dual Bandpass Frequency Selective Surface with High Selectivity", *IEEE Transactions on Antennas and Propagation*, vol. 72, pp. 6510–6519, 2024 (<https://doi.org/10.1109/TAP.2024.3423322>).
- [6] C. Li *et al.*, "An Ultra-wideband Bandpass Filter Based on Half-mode Substrate Integrated Waveguide Loaded with Rectangular-ring Slot", *2022 10th International Symposium on Next-Generation Electronics (ISNE)*, Wuxi, China, 2023 (<https://doi.org/10.1109/ISNE56211.2023.10221607>).
- [7] S. Keelailam *et al.*, "A Wideband Bandpass Filter Using U-shaped Slots on SIW with two Notches at 8 GHz and 10 GHz", *2022 IEEE*

*International Conference on Signal Processing and Communications (SPCOM)*, Bangalore, India, 2022 (<https://doi.org/10.1109/SPCOM55316.2022.9840774>).

- [8] Y.-F. Xie, Y.-H. Ma, and W.-S. Zhao, "A Single-notch Ultra-wideband Bandpass Filter Based on Bow-tie Cells", *2023 IEEE International Workshop on Electromagnetics: Applications and Student Innovation Competition (iWEM)*, Harbin, China, 2023 (<https://doi.org/10.1109/iWEM58222.2023.10234900>).
- [9] V.P.K. Kanaparthi, V.K. Velidi, R. Rajkumar, and R.R.T., "A Compact Ultra-wideband Multimode Bandpass Filter with Sharp-rejection Using Stepped Impedance Open Stub and Series Transformers", *IEEE Transactions on Circuits and Systems II*, vol. 69, pp. 4824–4828, 2022 (<https://doi.org/10.1109/TCSII.2022.3192512>).
- [10] T. Duraisamy *et al.*, "Compact Wideband SIW Based Bandpass Filter for X, Ku and K Band Applications", *Radioengineering*, vol. 30, pp. 288–295, 2021 (<https://doi.org/10.13164/re.2021.0288>).
- [11] K. Mahant *et al.*, "Spoof Surface Plasmon Polaritons and Half-mode Substrate Integrated Waveguide Based Compact Band-pass Filter for Radar Application", *PIER M*, vol. 101, pp. 25–35, 2021 (<https://doi.org/10.2528/PIERM20121803>).
- [12] R.S. Sangam and R.S. Kshetrimayum, "Hybrid Spoof Surface Plasmon Polariton and Substrate Integrated Waveguide Bandpass Filter with High Out-of-band Rejection for X-band Applications", *IET Microwaves Antennas and Propagation*, vol. 15, pp. 289–299, 2021 (<https://doi.org/10.1049/mia2.12049>).
- [13] E.M. Messaoudi, J.D. Martinez, and V.E. Boria, "Miniaturized Ultra-wideband Bandpass Filter Based on Substrate Integrated Quasi-lumped Resonators", *2021 IEEE MTT-S International Microwave Filter Workshop (IMFW)*, Perugia, Italy, 2021 (<https://doi.org/10.1109/IMFW49589.2021.9642382>).
- [14] N. Muchhal *et al.*, "Design of Hybrid Fractal Integrated Half Mode SIW Band Pass Filter with CSRR and Minkowski Defected Ground Structure for Sub-6 GHz 5G Applications", *Photonics*, vol. 9, art. no. 898, 2022 (<https://doi.org/10.3390/photonics9120898>).
- [15] A.J. Salim *et al.*, "Miniaturized SIW Wideband BPF Based on Folded Ring and Meander Line Slot for Wireless Applications", *2017 Second Al-Sadiq International Conference on Multidisciplinary in IT and Communication Science and Applications (AIC-MITCSA)*, Baghdad, Iraq, 2017 (<https://doi.org/10.1109/AIC-MITCSA.2017.8722992>).
- [16] W. Fu *et al.*, "Compact Bandpass Filter Based on Single Isosceles Right Triangular Substrate Integrated Waveguide Cavity and Modified Complementary Compact Microstrip Resonant Cell", *IEICE Electronics Express*, vol. 18, art. no. 20200397, 2021 (<https://doi.org/10.1587/elex.17.20200397>).
- [17] J. Liu, Y. Zhang, X. Wan, and H. Jing, "Miniaturized Dual-mode Ultra-wideband Filter Using Sector Substrate Integrated Waveguide", *Microwave and Optical Technology Letters*, vol. 63, pp. 2343–2347, 2021 (<https://doi.org/10.1002/mop.32912>).
- [18] D. Deslandes and K. Wu, "Accurate Modeling, Wave Mechanisms, and Design Considerations of a Substrate Integrated Waveguide", *IEEE Transactions on Microwave Theory and Techniques*, vol. 54, pp. 2516–2526, 2006 (<https://doi.org/10.1109/TMTT.2006.875807>).
- [19] Y. Cassivi *et al.*, "Dispersion Characteristics of Substrate Integrated Rectangular Waveguide", *IEEE Microwave and Wireless Components Letters*, vol. 12, pp. 333–335, 2002 (<https://doi.org/10.1109/LMWC.2002.803188>).
- [20] W. Hong *et al.*, "Half Mode Substrate Integrated Waveguide: A New Guided Wave Structure for Microwave and Millimeter Wave Application", *2006 Joint 31st International Conference on Infrared Millimeter Waves and 14th International Conference on Terahertz Electronics*, Shanghai, China, 2006 (<https://doi.org/10.1109/ICIMW.2006.368427>).
- [21] L. Huang and S. Zhang, "Ultra-wideband Ridged Half-mode Folded Substrate-integrated Waveguide Filters", *IEEE Microwave and Wireless Components Letters*, vol. 28, pp. 579–581, 2018 (<https://doi.org/10.1109/LMWC.2018.2835666>).

**Hussein Al-Jeshami, M.Sc.**

Electromechanical Engineering Department

 <https://orcid.org/0000-0001-7063-754X>

E-mail: hussein.m.abdul@uotechnology.edu.iq

University of Technology – Iraq, Baghdad, Iraq

<https://uotechnology.edu.iq>

**Hussam Al-Saedi, Ph.D.**

Communication Engineering Department

 <https://orcid.org/0000-0002-2029-7361>

E-mail: hussam.h.ali@uotechnology.edu.iq

University of Technology – Iraq, Baghdad, Iraq

<https://uotechnology.edu.iq>

**Mohammed F. Hasan, Ph.D.**

Communication Engineering Department

 <https://orcid.org/0000-0002-0244-9595>

E-mail: mohammed.f.hasan@uotechnology.edu.iq

University of Technology – Iraq, Baghdad, Iraq

<https://uotechnology.edu.iq>

**Halah I. Khani, M.Sc.**

Communication Engineering Department

 <https://orcid.org/0000-0001-5572-1073>

E-mail: halah.i.khani@uotechnology.edu.iq

University of Technology – Iraq, Baghdad, Iraq

<https://uotechnology.edu.iq>

**Muhannad Y. Muhsin, Ph.D.**

Electrical Engineering Department

 <https://orcid.org/0000-0003-3937-4467>

E-mail: muhannad.y.muhsin@uotechnology.edu.iq

University of Technology – Iraq, Baghdad, Iraq

<https://uotechnology.edu.iq>

**Wael M. Abdel-Wahab, Ph.D.**

Electrical and Computer Engineering

 <https://orcid.org/0000-0002-0994-3529>

E-mail: wmabdelw@uwaterloo.ca

University of Waterloo, Waterloo, Canada

<https://uwaterloo.ca>

# An Artificial Intelligence-based Handover Triggering and Management Mechanism for 5G Ultra-dense Networks to Improve Handover Authentication

P. Rajesh<sup>1</sup>, A. Vijayalakshmi<sup>1</sup>, and Ebenezer Abishek B.<sup>2</sup>

<sup>1</sup>Vels Institute of Science, Technology & Advanced Studies, India,  
<sup>2</sup>KCG College of Technology, India

<https://doi.org/10.26636/jtit.2025.2.2006>

**Abstract** — The emergence of 5G ultra-dense networks has gained considerable attention, as solutions of this kind enable rapid and intelligent device connectivity, thus ushering in a new era of high-speed communications. Nevertheless, the process of managing mobility across varying inter-frequency strategies increases interference and complexity. The development of a reliable handover algorithm is crucial for high-quality service, especially in ultra-dense networks with small cells. However, frequent handovers, ping-pong effects, and load-balancing issues arise due to the random and dense deployment of small cells. Additionally, ensuring secure and smooth handover authentication is critical, due to an increased risk of frequent transitions of users across different networks. In such a context, this research focuses on triggering handovers and managing 5G mobile networks, all while protecting sensitive data. We introduce an artificial intelligence-based approach aimed at improving handover initiation and management processes, leveraging Boruta random forest optimization (BRFO) to fine-tune handover margins and identify optimal trigger points for handovers. In addition, an impulsive graph neural network (IGNN) is utilized as a decision framework to predict and minimize unnecessary handovers, thus improving stability in small cell environments. Simulation results demonstrate that the proposed methodology significantly enhances handover management, strengthens authentication, and effectively mitigates a variety of attacks in 5G ultra-dense networks.

**Keywords** — authentication, communication security, handover, mobility management

## 1. Introduction

Ultra-dense networks (UDNs) have emerged as an innovation in 5G wireless communication [1]. By deploying a large number of small cells within a confined area, UDNs achieve significantly higher node density compared to traditional cellular networks [2]. They are engineered to meet the ever-growing demand for high data rates, ubiquitous connectivity, and ultra-low latency, challenges that conventional macrocell infrastructures struggle to address.

UDNs enhance network capabilities by improving coverage, boosting capacity, and ensuring better spectral efficiency

through small cell deployments. This strategy provides seamless connectivity and a superior user experience in dense urban environments, indoor settings, and high-traffic hotspots [3].

Moreover, UDNs integrate advanced technologies, such as mmWave communication, massive MIMO, and network densification to realize heterogeneous network architectures [4]. These advances are fundamental to supporting the newly emerging ultra-reliable low-latency communication (URLLC) services required by IoT, autonomous driving, AR, and VR applications [5]. UDNs offer improved flexibility, scalability and cost-efficiency compared to macro cell-based designs [5], optimizing spectrum usage and dynamic resource allocation, while simultaneously reducing construction and operational burdens [6].

Advances in wireless technologies and growing user demands have driven a significant evolution in mobile network handover (HO) mechanisms [7]. Initially, user handover decisions required manual intervention, which became impractical with the proliferation of mobile devices [8]. Automated HO systems were, therefore, introduced to facilitate seamless connectivity. In 2G (GSM) and 3G (UMTS) systems, network-controlled HO strategies became standard, leveraging signal strength, quality indicators, and mobility patterns to optimize HO decisions [9]. These automated techniques greatly improved service reliability and efficiency [10].

The arrival of 4G LTE further improved HO by introducing fast handover protocols that reduced latency and packet loss through proactive link establishment and network load balancing [11]. These improvements significantly improved mobile user experiences by enabling high-speed seamless connectivity [12].

In 5G, mobility management becomes even more critical, as URLLC, mMTC, and eMBB services demand different specialized handover mechanisms [13]. Technologies such as beamforming, flexible spectrum sharing, and network slicing help meet these demands by optimizing HO performance and reducing latencies [14]. As 5G deployments expand, handover protocols must ensure persistent connectivity and service continuity in diverse use cases.

Artificial intelligence (AI) has become a transformative force in improving 5G network performance [15]. Machine learning, deep learning, and natural language processing empower 5G networks with autonomous intelligence, enabling dynamic optimization and real-time decision-making. AI algorithms analyze massive datasets from network operations to detect trends, predict anomalies, and improve performance [16].

In dense 5G environments, AI facilitates dynamic network management, optimizing resource allocation, handover management, load balancing, and network slicing [17]. AI-driven automation reduces operating costs, improves service reliability, and accelerates service deployment [18].

Furthermore, AI enables predictive analytics and proactive maintenance by identifying potential failure points and mitigating risks before service interruptions occur [19]. It also plays a crucial role in fortifying the security of 5G networks. AI-based security solutions monitor traffic patterns, detect anomalies in real time, and respond proactively to cybersecurity threats [20]. These capabilities are essential for protecting sensitive information and critical infrastructures in next-generation networks against increasingly sophisticated cyberattacks.

In this work, we propose an AI-driven framework to enhance the efficiency, reliability, and security of handover triggering and management mechanisms, particularly emphasizing the improvement of handover authentication within 5G UDNs.

Optimizing handover triggering points is crucial in a mobile network to ensure efficient and seamless transitions between base stations (BSs), minimizing service interruptions, and maintaining optimal connectivity for users. The handover process in mobile networks is typically initiated based on certain trigger conditions, such as signal strength, speed of movement, and network load.

Below is a summary of the strategies we relied upon to optimize handover trigger points:

- 1) Handover optimization using Boruta random forest optimization (BRFO). We employ the BRFO algorithm to fine-tune handover parameters by dynamically adjusting the handover boundaries. This technique calculates the optimal conditions for initiating handovers, ensuring smooth transitions between base stations, while significantly minimizing connection interruptions. The integration of BRFO allows for adaptive optimization of the handover process, thus improving overall network continuity and user experience.
- 2) Reinforcement learning (RL) is used to continuously optimize handover triggers by rewarding the system for successful handovers (i.e. minimizing ping-pong effects and service interruptions) and penalizing failed handovers. Over time, the system learns the optimal handover trigger points for different scenarios.
- 3) Impulsive graph neural network (IGNN) for intelligent handover decision making. We utilize an IGNN model as a sophisticated decision-making tool to analyze network states and user mobility behaviors. This model predicts the need for handovers with a high degree of accuracy,

effectively reducing unnecessary handovers within small cell networks. By optimizing handover decisions, IGNN contributes to improved network efficiency, better resource allocation, and enhances the overall user experience.

The remainder of this paper is organized as follows. In Section 2, a review of the existing literature related to handover triggering and management in 5G UDNs is presented. Section 3 elaborates on the detailed operational workflow of the proposed system, focusing on BRFO-based handover optimization and IGNN-driven decision-making. Section 4 presents the simulation results and compares the performance of the proposed approach with existing methods, while Section 5 concludes the paper by providing final remarks and summarizing future research directions.

## 2. Literature Review

A review of the literature on handover triggering and management in 5G UDNs offers specific insights into the research conducted, methodologies adopted and challenges faced. It begins with an overview of the principles and architectures, highlighting such characteristics as low cell density and increased interference. Next, it examines previous studies on handover optimization, including signal strength-based and load balancing algorithms, as well as emerging AI and machine learning approaches.

Additionally, the review discusses security aspects, covering authentication protocols and encryption techniques. It also addresses challenges (as shown in Tab. 1), such as handover latency and interference mitigation, which are crucial for identifying research gaps and areas for improvement.

In [21], the authors discuss improved mobile broadband and extremely low dormancy communications that are supported by such technologies as 5G new radio and beyond. Due to the large number of mobile devices, it is important to manage high mobility in dense networks and constantly alter the time-to-trigger (TTT) and the hysteresis margin. The study suggests a mechanism for 5G and beyond that is based on online learning (learning-based intelligent mobility management – LIM2), to address these issues. For target cell selection, it uses SARSA-based reinforcement learning. For TTT and hysteresis adaptation, it uses the  $\epsilon$ -greedy strategy. This method shows promise as a means of improving mobility management and keeping advanced wireless networks connected without any interruptions.

The authors of [22] discuss the difficulties in managing handoffs in 5G mobile wireless networks that rely on UDN designs. Frequent turnover opportunities for user equipment in UDNs, which are defined by a large number of mmWave BSs, add complexity to the networks. Traditional handover plans simplify things too much, which results in more handovers than necessary and leads to poorer service quality. To address these problems, the authors of the study proposed a new transfer method known as FLDHDT. This technique relies on fuzzy logic to adapt the handover parameters, such as the handover margin (HOM) and TTT depending on the strength of the signal and the horizontal speed of the user equipment's move-

**Tab. 1.** Research gap summary.

Ref.	Methodology	Technique	Findings	Research gap
[21]	Mobility management in 5G	Adaptive time-to-trigger and hysteresis margin	Number of handovers and throughput	Misallocation of resources and overuse of electricity
[22]	Adaptive handover decision in UDNs	Fuzzy logic and time to trigger	Throughput and ping-pong ratio	Load balancing and inter-cell meddling have not been measured
[23]	MADM handover in 5G in UDNs	Fuzzy logic and MADM	Number of handovers and ping-pong handover	Lack of high-speed situations and ICI tests
[24]	ML protocol for secure 5G handovers	Burrows–Abadi–Needham (BAN) logic	Handover rate by 94.4%	The mobility of 5G UD HetNets needs to be clarified
[25]	Handover authentication mechanism in 5G HetNets	DHan_Auth and Conv_SLSTM	Attack detection accuracy 98.9832	This structure is vulnerable to DDoS outbreaks
[26]	Handover authentication in 5G HetNets	Fuzzy logic and key management	Latency and spatial complexity	Procedure flops to certify user discretion because of insecure channels
[27]	Hysteresis region authenticated handover for 5G HetNets	Artificial neural network and fuzzy logic (ANN-FL)	Handover success rate and communication overheads	A high number of needless handovers may occur in small cell networks at high speed
[28]	Secure handover protocol for 5G	ANN-FL	Handover success rate and handover failure rate	It is not appropriate for executing handovers in high-speed scenarios
[29]	Handover triggering estimation LTE-A/5G	Interval type II fuzzy logic system	Ping-pong handovers	Due to the restricted incidence choice, the recycling of incidence in 5G leads to co-channel interference
[30]	Proactive decision making for handover management 5G	Proactive decision making (PDM) and polynomial regression	Handover ping-pong, handover failure	The number of handovers would increase if MT travels at a high rate of speed

ment. By performing simulations and comparing them with traditional methods, the suggested plan is assessed. The results show that FLDHDT is effective in improving handover efficiency for 5G UDNs, compared to previous approaches. It reduces the number of handovers, lowers the ping-pong ratio, and overall system throughput.

Article [23] presents a new handover strategy to guarantee excellent service in UDNs and to reduce the impact of the aforementioned problems. The technique efficiently triggers handovers and transitions connections to nearby base stations by combining fuzzy logic with multiple-attribute decision algorithms (MADM). Fuzzy system membership functions are refined by subtraction grouping with past information available within the scheme, which improves performance. By reducing the frequency of handoffs, mitigating the impact of ping-pong, and maintaining high levels of service quality, the experimental results show that the suggested strategy outperforms traditional methods.

[24] introduces a machine learning-based handover authentication mechanism to tackle security, privacy, and efficiency issues. The protocol shows strong mutual authentication, protection of session keys, and resistance to several attacks in the course of a formal security analysis using the Burrows-Abadi-Needham (BAN) logic. It also guarantees user anonymity, mutual authentication, and complete confidentiality of key ciphers, according to informal security assessments. Compared to the enhanced 5G identification and key agreement (5G AKA) protocol, the simulation results show better performance metrics. It significantly improves the efficiency of handover signaling and achieves a staggering 94.4% drop-in turnover rate.

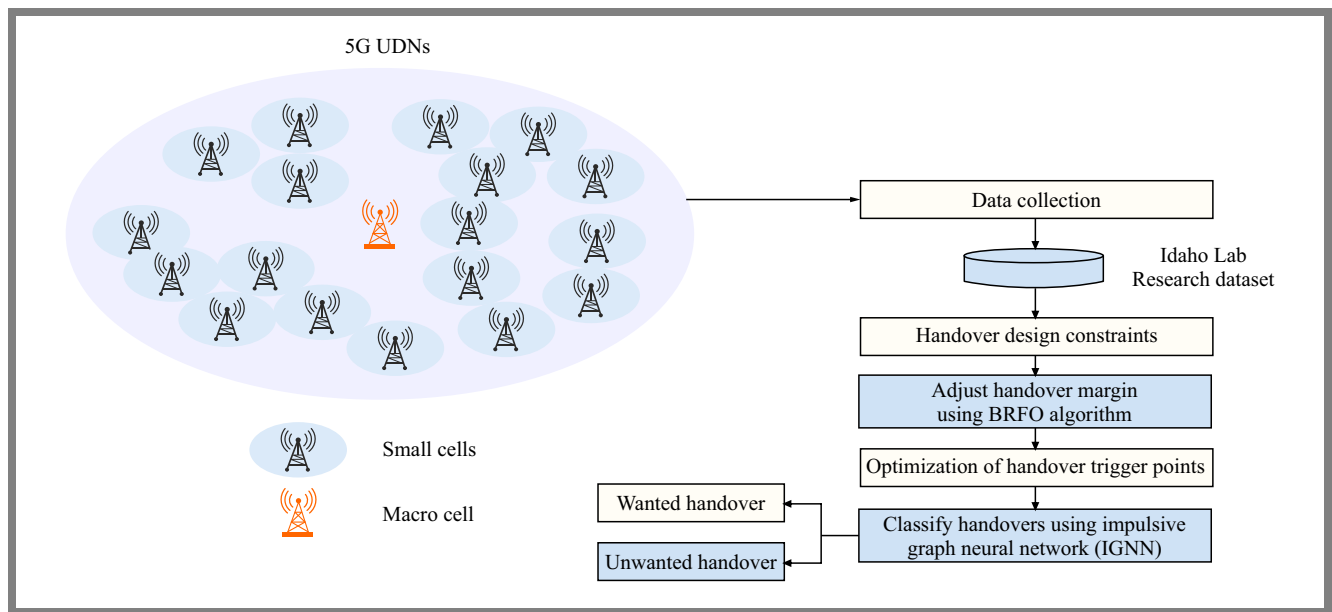
A deep learning-based handover authentication system is presented in [25] to solve these issues and enhance user experience. Using the 5G handover-authentication and key agreement (5G AKA) protocol, only data belonging to non-malicious users are authenticated once they have been classified using convolution-stacked long short-term memory

(Conv\_SLSTM) networks. Encryption and decryption are handled by the authentication process using extended elliptic curve cryptography (Ex\_ECC). An evaluation of the model’s performance on the Python platform shows that it improves handover processes and resists network attacks with a classification precision of 98.98% and a handover delay of 11.8 s for 200 nodes.

The authors of [26] suggested a solution that relies on fuzzy logic for key and handover management to improve the performance of cloud handover control and identification mechanisms in 5G networks. The goal of the fuzzy logic model is to reduce delays and maximize network efficiency by minimizing handovers and optimizing the selection of the target cell using several factors. The results showed that the model is capable of reducing latency, validating authentication threats, and handling geographic complexity, all of which are important concerns in the management and deployment of 5G networks.

To close the gap, 3GPP has included authenticating and key exchange protocols in its 16th release (3GPP R16). Privacy- and security-related standards applicable to 5G networks are quite high and, although there is a certain number of security protocols described in the literature, many of them are either ineffective or fail to meet these standards. A protocol that simultaneously prioritizes efficiency, security, and privacy is proposed in [27]. In order to ensure user devices, source gNB, and target gNB identification and session key establishment during handover, an intelligent model is created and deployed for target cell prediction using artificial neural networks and fuzzy logic.

Paper [28] presents an ANN-FL protocol that prioritizes security and service quality to solve the problems and meet the changing needs of 5G and beyond (B5G) networks. Simulation results show that thanks to reducing ping-pong handovers by 24.1%, increasing the success rate of handovers by 27.1%, and reducing the failure rate of handovers by 27.3%, the pro-



**Fig. 1.** Handover triggering and management mechanism for 5G UDNs using AI techniques.

protocol is robust against various attacks and effectively improves security and performance in 5G and B5G environments.

The authors of [29] present a new method to estimate the radio link quality (RLQ) of the serving and nearby cells, and then use fuzzy logic to trigger handover procedures. The system uses a simple fuzzy logic system for the prediction of RLQ and a second-order regressors for RLQ prediction. For handover trigger decisions, it uses a cascade fuzzy logic system and successfully mitigates ping-pong, premature, or delayed handovers. Simulation findings show that by using solely information on the quality of the radio connection, handover performance is significantly improved by 50% in high-speed situations, with the ns-3 LTE module. Significantly, the approach solves important 5G network management problems while remaining easy to implement and not being constrained by UE velocity, making it suitable for a wide range of applications, such as UAVs and IoT devices.

In [30], a suitable handover management strategy is proposed to solve mobility-related problems. Its primary objective is to investigate the impact of the handover control parameter on the operation of 5G networks through proactive decision-making in the cell selection process. The method was tested in 5G HetNet simulations to establish its impact on improving mobility management in these networks. The tests included measuring handover attempts, ping-pong transfers, handover mistakes, radio link mistakes, and transfer delays.

Several research challenges have been identified in the literature [21]–[31], including issues related to mobility management, handover optimization, and secure authentication mechanisms. Traditional handover algorithms may not be sufficient for the dynamic nature of ultra-dense networks, leading to problems such as handover latency, ping-pong effects, and load balancing concerns. Therefore, an efficient handover algorithm that optimizes handover triggers and minimizes the amount of unnecessary handover operations.

Additionally, as users move across multiple networks in 5G environments, robust handover authentication solutions are needed to protect against potential security threats.

Despite the progress made in network architecture design, significant challenges remain in ensuring both effective and secure handover authentications, especially in scenarios where pairing is unnecessary. Existing HO authentication protocols, such as identity-based cryptographic techniques, might not fully satisfy the rigorous security requirements of 5G networks, particularly during inter-operator handovers in smaller network regions. Furthermore, the decision-making processes during HO execution in 5G networks are complex, leading to architectural adjustments that may be economically inefficient. The computational burden of HO authentication protocols also poses challenges in meeting the stringent delay demands of 5G networks.

### 3. Methodology

The proposed mechanism for HO triggering and management in 5G UDNs, illustrated in Fig. 1, utilizes AI-driven techniques. Data from a 5G network are collected and stored within the Idaho Lab Research dataset. Handover constraints are addressed by dynamically adjusting the handover margin through the Boruta random forest optimization (BRFO) algorithm, which refines the handover trigger points. Finally, IGNN is used to classify and differentiate between legitimate and unnecessary handovers.

#### 3.1. Handover Optimization

To minimize HO-related issues in 5G networks, particularly in ultra-dense environments with small cells, optimization of the handover margin is essential. The handover margin sets the threshold signal strength difference required to trigger

a handover decision. An effective adjustment of this margin is capable of reducing unnecessary handovers and mitigating ping-pong effects, i.e. scenarios in which users frequently switch between adjacent base stations.

In the presented approach, we use BRFO [31] to fine-tune the handover margin. BRFO merges the Boruta feature selection algorithm and the random forest algorithm, leveraging their strengths to pinpoint the most influential features in order to optimize handover performance. By iteratively assessing significance of the HO margin in conjunction with other influencing factors, such as signal strength, interference, and user mobility, BRFO calculates the optimal HO trigger points.

The mean minimization accuracy or variant number of the randomization value for every input  $p_s$ , together with the matching shadier input  $p_s^b$  for the total number of trees is:

$$mda = \frac{1}{M_{tree}} \cdot \sum_{a=1}^{a_{tree}} \frac{\sum_{s \in oon} H(q_s = F(p_s)) - \sum_{s \in oon} H(q_s = F(p_s^b))}{|OON|}, \quad (1)$$

$H(\cdot)$  represents the indicator function, whereas OON refers to the predicted error of every sample used for training, calculated using bootstrap aggregation. Calculation of  $W$  scores is performed in the following manner:

$$W - score = \frac{mda}{sd}. \quad (2)$$

Let us compute the highest  $W$  score in the shadow characteristics by using the standard deviation  $sd$  of exact losses. The predictor applies a normalization technique to the data set, scaling the values from 0 to 1 to minimize the impact of extreme values.

$$\alpha_{norm} = \frac{\alpha - \varepsilon_{min}}{\alpha_{max} - \alpha_{min}}. \quad (3)$$

As a result, current input  $P_s$ , memory cell output  $i_{s-1}$  from the earlier time step  $s-1$ , and bias terms  $bf$  are used to calculate the activation values of forgetting gate  $ft$  at time step  $t$ . All activation values are divided by the sigmoid function between 0 (totally forget) and 1 (totally recall):

$$F_s = \text{sigmoid}(Z_{F,p} P_s + Z_{F,i} P_{s-1} + n_F). \quad (4)$$

Also, the second step defines the LSTM cover to be included in grid cell positions  $t_s$ . This job involves two actions [32]. First, we calculate applicant values that can be added to cell positions. Second, the input gate activation values are calculated as:

$$t_s = \tan i(Z_{t,p} p_s + Z_{t,i} i_{s-1} + n_t), \quad (5)$$

$$h_s = \text{sigmoid}(Z_{h,p} p_s + Z_{h,i} i_{s-1} + n_h). \quad (6)$$

In the third stage, the Hadamard product is defined by creating new cell locations  $t_s$  based on the outcomes of the preceding processes:

$$t_s = F_s \circ t_{s-1} + h_s \circ T_s. \quad (7)$$

Output  $i_s$  of the reminiscence cells is computed as the subsequent function, in the following manner:

$$o_s = \text{sigmoid}(Z_{o,p} p_s + Z_{o,i} i_{s-1} + n_o), \quad (8)$$

$$i_s = o_s \tan i(t_s). \quad (9)$$

At this stage, the system processes input  $s$  at each time point as defined by Eqs. (1)-(9). The output of each gate is obtained by a logic function and a non-linear alteration of the contribution. The following describes the link between input and outcome.

$$R(s) = \sigma_j(Z_R p(s) + u_R i(s-1) + n_R), \quad (10)$$

$$w(s) = \sigma_j(Z_w p(s) + u_w i(s-1) + n_w), \quad (11)$$

$$i(s) = (1 - w(s)) o(s-1) + w(s) o \hat{i}(s), \quad (12)$$

$$\hat{i}(s) = \sigma_i(Z_i p(s) + u_i R(s) o i(s-1)) + n_i, \quad (13)$$

where  $w(s)$  is the apprise gate trajectory,  $R(s)$  is the rearrange gate trajectory, with  $Z$  and  $u$  being stricture metrics and vector, respectively.  $\sigma_j$  is a sigmoid purpose and  $\sigma_i$  is referred to as a hyperbolic angle.

Algorithm 1 describes the process of optimizing HO using BRFO.

### 3.2. Handover Decision Model

The HO decision model plays a crucial role in managing handovers within small cell networks, where users' frequent mobility creates numerous handover opportunities. This model helps determine the optimal moments for handovers, ensuring

---

#### Algorithm 1 Handover optimization using BRFO

---

**Input:** HO design constraints, margin, trigger point

**Output:** HO optimization parameters

**Start**

- 1: Init. population  $P_s^b$  with candidate HO configurations
- 2: **for** each solution  $s \in P$  **do**
- 3:     **if** the input  $P_s^b$  for the total amount of trees **then**
- 4:         Randomly generate  $P_s^b$
- 5:     **end if**
- 6:     **if**  $P_s^b$  is defined **then**
- 7:         Compute  $W$  score using Eq. (2)
- 8:         Normalize the predictor of data set between 0 and 1 by Eq. (3)
- 9:     **end if**
- 10: **end for**
- 11: **for** each input  $P_s^b$  **do**
- 12:     Compute input gate activation values from Eq. (5)
- 13:     Formulate input-output relationship by Eq. (10)
- 14:     Find the fitness  $F_s$  value
- 15:     **if** better value  $F_s$  is found **then**
- 16:         Update final value  $F_s$
- 17:     **end if**
- 18: **end for**

**End**

---

that only necessary handovers occur. By effectively predicting unwanted handovers, the model reduces ping-pong effects and prevents excessive amounts of handover attempts, leading to a more efficient network.

An IGNN is used as the decision-making mechanism [35], [36]. IGNNs are specialized neural networks designed for processing graph-structured data, making them particularly suitable for network optimization and decision-making tasks. In HO management, IGNN analyzes key network parameters and mobility patterns to assess the probability of unwanted handovers. By learning from historical data on handovers and network performance, IGNN identifies patterns that suggest unwanted handovers, allowing it to make more accurate decisions. Let us consider a scenario where  $B$  represents a set of interconnected neural networks, each comprising identical types of networks, with both linear and quadratic components at each node of the  $B$ -dimensional system. The differential equation defining this network is described as follows:

$$\begin{aligned} \dot{p}_h = & -D_1 p_h(a) + D_2 p_h(ya) + N_1 F_1(p_h(a)) \\ & + N_2 F_2(p_h(ya)), \quad a \geq a_0. \end{aligned} \quad (14)$$

The state lattice of the  $h$ -th brain framework at a given time demonstrates the postpone importance and  $1 - y$  is ordinarily alluded to as the beginning significance and relative deferral:

$$\begin{aligned} F_R(p_h(a)) = & [F_{R1}(p_h(a)), F_{R2}(p_h(a)), \dots, F_{Rb}(p_h(a))]^S, \quad R = 1, 2. \end{aligned} \quad (15)$$

The ensuing straight-coupled differential capability depicts the fluctuating activities of interconnected brain organizations:

$$\begin{aligned} \dot{p}_h(a) = & -D_1 p_h(a) + D_2 p_h(ya) + N_1 F_1(p_h(a)) \\ & + N_2 F_2(p_h(ya)) + d \sum_{g=1, g \neq h}^R m_{hg} \Gamma(p_g(a) - p_h(a)), \end{aligned} \quad (16)$$

where  $d$  is the strength of connection  $m_{hg}$  and  $\Gamma$  is a remotely associated unequivocal positive network between two vertices  $h$  and  $g$ . It is defined as follows, when node  $g$  and node  $h$  are connected:

$$\text{if } g \neq h \text{ then } m_{hg} > 0. \quad (17)$$

otherwise

$$m_{hg} = 0, \quad m_{hg} = - \sum_{g=1, g \neq h}^B m_{hg}.$$

The state of the linking  $p_g(a) - p_h(a)$  is linked and nodes  $g$  and  $h$  vary due to excitement at a specific time  $a_K$ . Therefore, the neural networks associated with the stimuli can be obtained in the following form:

$$\begin{cases} \dot{p}_h(a) = -D_1 p_h(a) + D_2 p_h(ya) + N_1 F_1(p_h(a)) \\ + N_2 F_2(p_h(ya)) + d \sum_{g=1, g \neq h}^B m_{hg} \Gamma(p_g(a)), \quad a \neq a_K \\ p_g(a_K^+) - p_h(a_K^+) = j_K (p_g(a_K^-) - p_h(a_K^-)), \quad m_{hg} > 0 \end{cases} \quad (18)$$

where  $\varsigma = \{a_1, a_2, a_3 \dots\}$  is a rash series nutritious,  $a_{K-1} < a_K$  represents the number of careless occurrences of the impulsive sequence  $\zeta$  during the interlude  $(t, a)$  and  $j_K$  indicates the impulsive signal's gain. This is the Laplacian matrix of the compliance system topology. The impulsive sequence  $\zeta$  a  $V$ -asymptotic regular  $S_{asy}^V$  impetuous period is:

$$\lim_{K \rightarrow \infty} (V(a_{K+1}) - V(a_K)) = a_{asy}^V. \quad (19)$$

Let  $\bar{m}, \bar{n}$ , and  $q$  be real numbers, and let  $b$  be greater than 0. Let  $y$  be a real number between 0 and 1. Recognize that the given explanation serves as an explanation.

$$\begin{cases} \dot{p}(a) = \bar{m} p(a) + \bar{n} p(ya), \quad a \geq a_0, \quad a \neq a_k \\ p(a_K^+) = j_K p(a_K^-) \end{cases}. \quad (20)$$

Assuming that  $p(a)$  is greater than zero, let  $x(a)$  be a non-negative functional defined on interval  $[, +\infty)$  that satisfies:

$$\begin{cases} \dot{x}(a) \leq \bar{m} p(a) + \bar{n} p(ya), \quad a \geq a_0, \quad a \neq a_k \\ x(a_K^+) \leq j_K p(a_K^-) \end{cases}. \quad (21)$$

Given that 0 is less than  $x(a)$  and  $x(a)$  is less than  $p(a)$  for any  $s$  in interval  $[, ]$ , for all values of  $s$  greater than or equal to a certain value,  $x(a)$  is less than or equal to  $p(a)$ . Therefore, for  $p(a) x(a)$  with  $0 < x(a) < p(a)$  for  $s \in [, ]$ :

$$x(a) < p(a), \quad \text{for all } a \geq a_0. \quad (22)$$

where  $S > 0$ , such that set:

$$W = \{a \in (a_0, a) : x(a) \geq p(a)\}, \quad x(a^*) = p(a^*).$$

and

$$x(a) < p(a), \quad x(a^*) \geq p(a^*). \quad (23)$$

We compute the optimal threshold condition as follows:

$$\dot{x}(a) = \bar{m} x(a^*) + \bar{n} x(ya^*). \quad (24)$$

We compute the maximum and minimum range of threshold condition  $x(a^*) = p(a^*)$  and  $x(ya^*) = p(ya^*)$ , which generates the following set of conditions.

$$\begin{aligned} 0 & \leq \dot{x}(a^*) - \dot{p}(a^*) \\ & \leq (\bar{m} x(a^*) + \bar{n} x(ya^*)) - (\bar{m} p(a^*) + \bar{n} p(ya^*)) \\ & = \bar{n} x(ya^*) - p(ya^*) \\ & < 0 \end{aligned} \quad (25)$$

Worldwide  $\mu$ -dependability model follows the Dasey  $< \infty$  condition and  $S$  condition with drive-related brain networks when coordinated upgrades or non-synchronized improvements occur during the motivation span. In this way, drive-associated brain organizations can be reworked in the Kronecker item structure:

$$\begin{cases} \dot{p}(a) = -(H_B \otimes D_1)p(a) + (H_B \otimes D_2)p(ya) \\ + (H_B \otimes N_1)f_1(p(a)) + (H_B \otimes N_2)f_2(p(ya)) \\ + d(M \otimes \Gamma)p(a), \quad a \neq a_K, \quad K \in B \\ p_g(a_K^+) - p_h(a_K^+) = j_K (p_g(a_K^-) - p_h(a_K^-)), \\ \text{for } (h, g) \text{ satisfying } m_{hg} > 0. \end{cases} \quad (26)$$

In this case, the network topology exhibits robust connectivity, indicating that the Laplacian connected matrix  $A$  remains unchanged. Algorithm 2 outlines the operational procedure of the HO decision model employing IGNN.

---

**Algorithm 2** Handover decision model using IGNN

---

**Input:** Number of small cells and macro cells, threshold condition

**Output:** Handover decision wanted and unwanted

**Start**

```

1: Initialize the random population
2: if the network is initialized then
3:   Describe it using the Eq. (14)
4: end if
5: if  $i = 0$  then set  $j = 1$ 
6: end if
7: while condition is true do
8:   if the system study state then
9:     unwary arrangement  $\zeta$  is as given in Eq. (19)
10:  else if  $p(s)$  is valid then
11:    Recognize it as solution to Eq. (20)
12:  end if
13:  if a non-negative function  $\chi(a)$  exists on  $[y_{s0}, +\infty)$ 
14:    then
15:      Ensure it satisfies Eq. (21)
16:       $x(a) < p(a)$  for all  $s \in [a_0, a_1]$ ,
17:    else Revise
18:  end if
19: end while

```

**End**

---

## 4. Results and Discussion

In the next step, a comparative analysis between the proposed HO triggering and management mechanism and existing approaches is conducted. Performance is validated using the Idaho Laboratory Research dataset. The proposed handover trigger and management mechanism is implemented on the Google Colab platform using Python.

We compare the results of the BRFO+IGNN mechanism with those obtained using existing solutions, including conventional Event A3, FLDH [37] and FLDHDT [22]. Furthermore, the results of the handover authentication of the BRFO+IGNN mechanism are compared with existing mechanisms, such as transport layer security (TLS), fuzzy systems, fuzzy transport layer security (F-TLS) and convolutional SLSTM (CLSTM) [25].

For the handover decision-making process, we compare the performance of the BRFO+IGNN mechanism with several benchmark models, including random forest (RF), decision tree (DT), naive Bayes (NB), linear regression (LR), support vector machine (SVM) and XGBoost.

### 4.1. Simulation Setup

The data set utilized in this study includes both normal and attack data generated within a simulated setting. Data was

**Tab. 2.** Simulation setup.

Parameter	Value
Network size	1000 × 1000 m
Number of evolved nodes	3
Number of users	100–500
Amount of pieces of user equipment	5
Mobility model	2D random walk
Speed of user equipment	2–20 m/s
Power of evolved nodes	43 dBm
Power of next generation nodes	23 dBm
Frequency of evolved nodes	2.4 GHz
Frequency of next generation nodes	28 GHz
Packet inter-arrival time	20 ms
Packet size	1000 bytes
Bandwidth of evolved nodes	20 Mbps
Bandwidth of next generation nodes	100 Mbps
Simulation time	100 s

gathered from an Internet connected Linux machine running a 5G core network with open-source 5G core software. The network traffic on the 5G core machine limits was captured via Wireshark.

Normal data are categorized into two groups: one involving a single-user equipment simulation and the other involving two user equipment simulations. Malicious data consist of ten distinct attacks, classified into three primary categories: reconnaissance, denial of service (DoS), and network reconfiguration.

Reconfiguration attacks include unified data management, get all network functions, get user data, automatic redirect with a timer, and random data dump. Network reconfiguration attacks are divided into false access and mobility management function insert and delete attacks, as well as random access and mobility management function insert and delete attacks.

The DoS category includes the crash network repository function attack. The data set, covering a total of 50 000 records, is divided using the following proportions: 80% for training and 20% for testing. Data are exported in the CSV format and are used in the proposed research. The analysis considers such attributes as time, source, destination, protocol, length, sequence amounts, acknowledgment amounts, window size, length, timestamp echo reply field, and timestamp value field.

Table 2 presents the parameters used in the simulation setup, which define the characteristics of the simulated network environment necessary to evaluate the proposed mechanisms and algorithms. Together, these parameters create a realistic simulated environment that allows to effectively evaluate the proposed solutions.

**Tab. 3.** Comparative analysis of proposed and existing HO mechanisms.

Handover mechanism	Number of users				
	100	200	300	400	500
Number of handovers					
Event A3	693	821	1004	1158	1321
FLDH	570	698	881	1035	1198
FLDHDT	447	575	758	912	1075
BRFO+IGNN	324	452	635	789	952
Ping-pong ratio [%]					
Event A3	8.34	8.76	9.09	9.26	9.47
FLDH	5.97	6.40	6.72	6.89	7.11
FLDHDT	3.60	4.03	4.35	4.52	4.74
BRFO+IGNN	1.23	1.66	1.98	2.15	2.37
System throughput [Mbps]					
Event A3	93.72	155.96	242.86	355.83	459.63
FLDH	146.88	209.13	296.02	408.99	512.79
FLDHDT	200.05	262.29	349.19	462.16	565.96
BRFO+IGNN	253.21	315.46	402.35	515.33	619.12

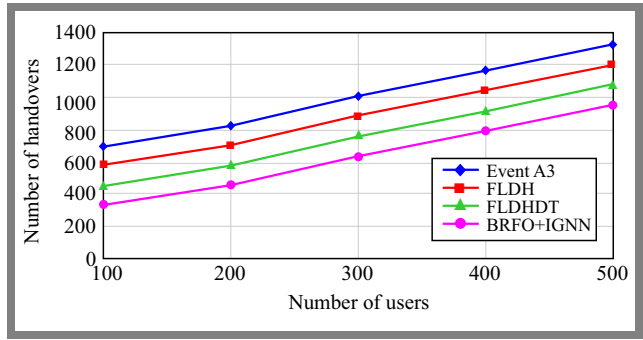
### 4.2. Comparison

Table 3 presents a comparative evaluation of the proposed HO mechanism against the existing approaches. Figure 2 illustrates the number of handovers corresponding to the varying number of users in different HO mechanisms. The data reveal distinct patterns in handovers as the number of users increases from 100 to 500. Event A3 shows a steady rise in handovers, experiencing a 90.7% increase from 693 at 100 users to 1321 at 500 users. Similarly, FLDH shows a continuous increase, with a 70.2% rise from 570 to 1198 HO. FLDHDT follows a similar trend, showing a significant 139.3% increase from 447 to 1075 handovers.

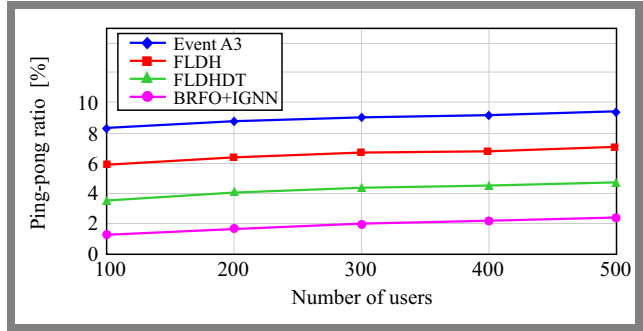
On the contrary, BRFO + IGNN shows a decreasing trend in HOs as the number of users grows, although the change still reflects a 66.7% drop from 324 to 952 handovers. Event A3, FLDH, and FLDHDT handover mechanisms show a positive relationship between user count and HOs, with increases ranging from 90.7% to 139.3%. However, BRFO + IGNN shows a negative correlation, with a decrease in HOs by 66.7% despite a growing user base.

These results suggest varying levels of efficiency and scalability across the mechanisms, underlining the need to select the most suitable approach based on specific network conditions and user requirements.

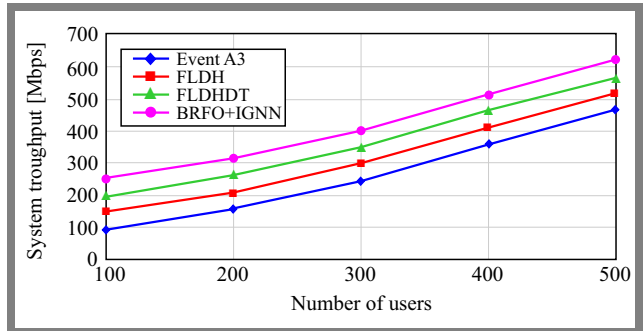
Figure 3 illustrates the ping-pong ratio across various user numbers for different HO mechanisms. As the number of users increases from 100 to 500, noticeable differences in the ping-pong ratios are observed among the mechanisms. Event A3 shows a consistent upward trend, with the ratio rising from 13.1% to 13.9%. Similarly, FLDH shows a steady increase in the ping-pong ratio, as it climbs from 18.7% to 19.0%. FLDHDT follows a similar pattern, with values ranging from 24.4% to 31.5%.



**Fig. 2.** Handovers against the number of users.



**Fig. 3.** Ping-pong ratio against the number of users.



**Fig. 4.** System throughput against the number of users.

On the contrary, RFO + IGNN reveals an opposite trend, where the ping-pong ratio decreases as the number of users increases. However, the change still varies, with increases from 91.7% to 91.1% within the user range. In general, Event A3, FLDH and FLDHDT exhibit a positive correlation between the number of users and the ping-pong ratio, increasing from 13.1% to 31.5%.

On the other hand, BRFO+IGNN demonstrates a negative correlation, even with increases of 91.1% to 91.7%. These results highlight the differing scalability and performance of HO mechanisms, stressing the importance of selecting the most suitable mechanism based on the specific network demands and user conditions.

Figure 4 presents the system throughput across different HO mechanisms as the number of users increases. As the user count increases from 100 to 500, distinct patterns in system throughput may be observed for each HO mechanism. Event A3 shows a steady increase in throughput, with improvements ranging from 390.1% to 391.8% over the user range. Similarly, FLDH shows a gradual rise in throughput, ranging from

**Tab. 4.** Comparative analysis of proposed and existing HO authentication mechanisms

HO authentication mechanism	Number of users				
	100	200	300	400	500
Authentication latency [s]					
TLS	8.52	14.32	20.15	25.64	31.25
Fuzzy	7.12	12.02	18.64	21.46	28.56
F-TLS	5.07	10.35	15.64	18.25	25.03
CLSTM	4.23	8.12	9.94	15.35	21.48
BRFO+IGNN	3.53	5.12	7.54	13.32	18.78
Number of unsuccessful handover authentications					
TLS	24	53	105	185	231
Fuzzy	21	46	101	142	195
F-TLS	18	40	98	120	174
CLSTM	14	35	75	85	152
BRFO+IGNN	8	24	55	62	112
Handover delay [ms]					
TLS	18.18	19.60	21.14	23.57	27.15
Fuzzy	14.52	15.95	17.49	19.92	23.50
F-TLS	10.87	12.29	13.83	16.27	19.84
CLSTM	7.22	8.64	10.18	12.61	16.19
BRFO+IGNN	3.56	4.99	6.52	8.96	12.53
Packet loss rate [%]					
TLS	14.40	14.48	14.60	14.67	14.77
Fuzzy	10.84	10.92	11.04	11.11	11.21
F-TLS	7.28	7.36	7.48	7.55	7.65
CLSTM	3.71	3.80	3.92	3.99	4.09
BRFO+IGNN	0.15	0.23	0.36	0.42	0.53

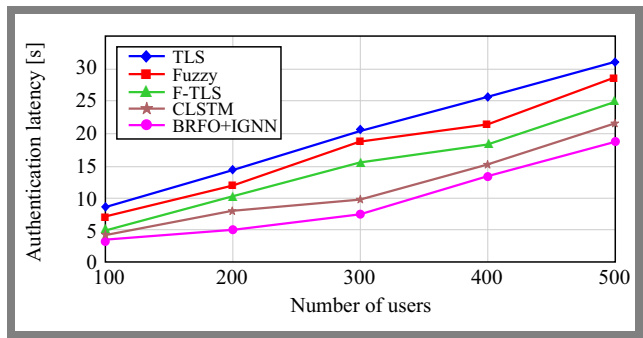
249.6% to 249.8%. FLDHDT follows a similar trend, with increases between 182.9% and 183.5%.

On the contrary, BRFO + IGNN shows a consistent increase in system throughput as the number of users increases, but the rate of change is smaller, fluctuating between 144.4% and 144.7%. Both Event A3, FLDH, and FLDHDT show a positive relationship between system throughput and the number of users, with increases varying from 182.9% to 391.8%.

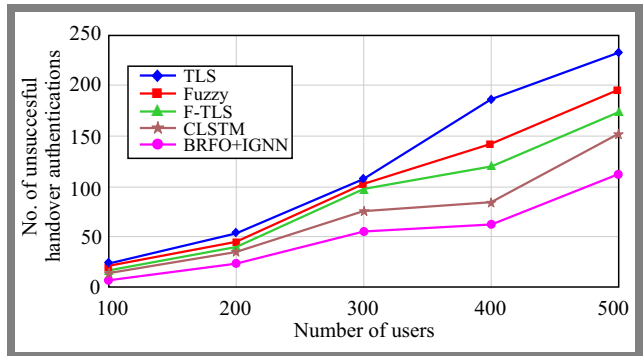
BRFO+IGNN also exhibits a positive trend, but with smaller increases, ranging from 144.4% to 144.7%. These observations highlight the different efficiencies and scalability of the HO mechanisms, underlining the importance of selecting the appropriate method depending on network needs and user scenarios.

### 4.3. Comparison of HO Authentication

Table 4 presents a comparison of the proposed HO authentication method with existing solutions. Figure 5 illustrates the authentication latency as the number of users changes for various HO authentication techniques. As the user count rises from 100 to 500, noticeable trends appear in authentication latency across the different methods. TLS authentication



**Fig. 5.** Authentication latency against the number of users.



**Fig. 6.** Amount of unsuccessful handover authentications against the number of users.

shows a steady increase in latency, with an improvement between 52% and 59% across the user range. Likewise, fuzzy authentication displays a gradual increase in latency, with an improvement between 35.23% and 40.12%.

In contrast, F-TLS authentication shows a reduction in latency as user numbers grow, with improvements ranging from 53.7% to 53.8%. CLSTM authentication also reveals an improvement in latency as the number of users increases, ranging from 61.6% to 61.5%. BRFO+IGNN authentication consistently improves latency even as the number of users increases, with an improvement from 46.9% to 46.8%. Both TLS and fuzzy authentication mechanisms exhibit a direct correlation between the number of users and authentication latency, with improvements from 62.6% to 70.125%.

However, the F-TLS, CLSTM, and BRFO+IGNN authentication methods show an inverse correlation, with enhancements ranging from 46.8% to 61.8%. These results highlight the varying performance and scalability of each authentication method, emphasizing the need to choose the most suitable method according to specific security demands and user conditions.

Figure 6 illustrates the number of unsuccessful handover authentications, as the number of users varies between different HO authentication methods. As the number of users increases from 100 to 500, trends in unsuccessful HO authentications for each mechanism become apparent. TLS authentication shows a steady increase in unsuccessful handovers, with a rate ranging from 12.54% to 15.12% as the number of users increases. Similarly, the fuzzy authentication method also experiences a gradual increase, with rates ranging from 25.24% to 28.62%. On the contrary, F-TLS authentication shows fluctuating

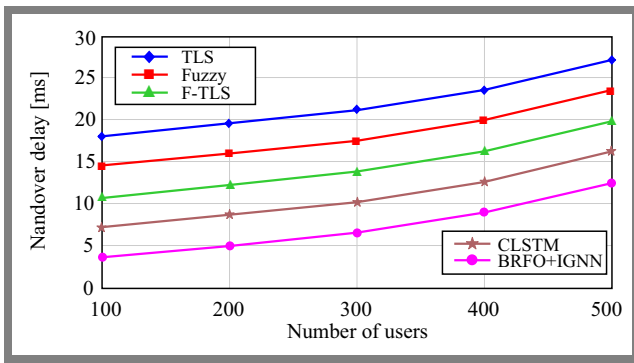


Fig. 7. Handover delay against the number of users.

tuating patterns, although there is a general increase, with rates varying from 32.15% to 36.53%. The CLSTM method follows a comparable trend, with its unsuccessful HO ratio ranging from 11.25% to 15.63%.

The BRFO+IGNN method also shows a continuous increase in unsuccessful handovers, even as the user count grows. Across all authentication methods (TLS, Fuzzy, F-TLS, CLSTM and BRFO+IGNN), there is a positive correlation between the number of users and the rate of unsuccessful HO authentications, with enhancements ranging from 23.51% to 28.62%. These results highlight the importance of evaluating the scalability and reliability of authentication mechanisms, stressing the need to address potential security vulnerabilities and optimize performance in different user contexts and security requirements.

Figure 7 illustrates the HO delay in relation to the varying user counts for different handover authentication methods. As the user count increases from 100 to 500, distinct trends are observed across the mechanisms. TLS authentication shows a gradual increase in the delay in HO, with improvements varying between 49.8% and 49.3%. Fuzzy authentication displays a steady increase in delay, with enhancements between 61.9% and 61.9%.

In contrast, F-TLS authentication reveals a reduction in delay as user numbers grow, with improvements ranging from 33.7% to 33.8%. CLSTM also shows a decrease in the delay with user count, with enhancements between 55.7% and 55.6%. BRFO + IGNN consistently reduces delay as the user count increases, with improvements ranging from 71.3% to 71.1%. TLS and fuzzy mechanisms exhibit a positive relationship between user numbers and HO delay, with enhancements of 49.3% to 61.9%. On the other hand, F-TLS, CLSTM, and BRFO+IGNN demonstrate a negative relationship, with enhancements ranging from 33.7% to 71.3%.

These results underscore the importance of optimizing HO mechanisms to reduce delays and improve network performance based on specific user requirements.

Figure 8 illustrates the packet loss rate as a function of the number of users for various HO authentication mechanisms. As the number of users increases from 100 to 500, different trends are observed in the packet loss rate. TLS authentication shows a steady increase in packet loss, with improvements ranging from 2.6% to 2.6% across the user range. Similarly, fuzzy authentication demonstrates a gradual increase in packet

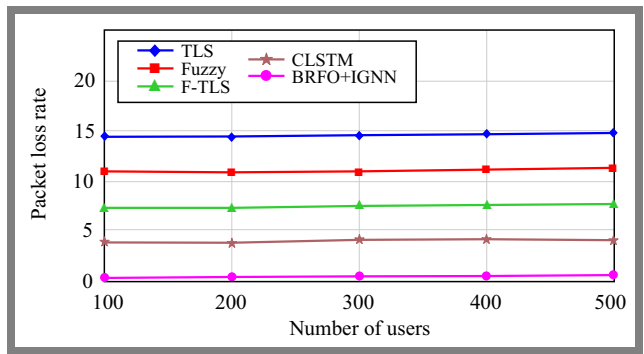


Fig. 8. Packet loss rate against the number of users.

Tab. 5. Comparative analysis of the proposed and benchmark HO authentication mechanisms [%].

HO decision model	Accuracy	Precision	Recall	F-measure
RF	54.66	53.24	53.67	53.46
DT	61.64	60.23	60.66	60.44
NB	68.63	67.21	67.64	67.43
LR	75.61	74.20	74.63	74.41
SVM	82.60	81.18	81.61	81.40
XGBoost	89.58	88.17	88.60	88.38
BRFO+IGNN	96.57	95.15	95.58	95.37

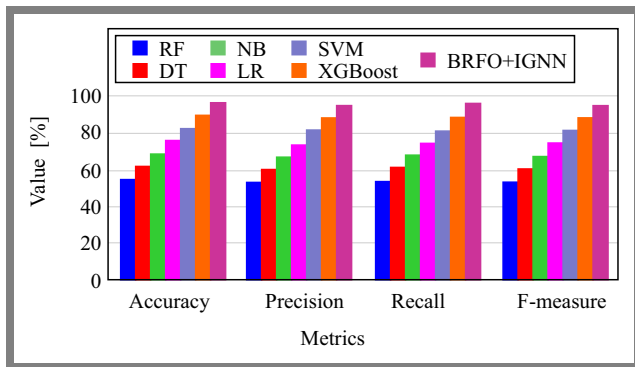
loss, with improvements varying from 3.4% to 3.4%, while the F-TLS authentication method exhibits a reduction in packet loss as the number of users grows, with improvements ranging from 1.3% to 1.3%.

The CLSTM authentication method also shows a decrease in packet loss in the user range, with improvements of 5.4% to 5.3%. BRFO + IGNN authentication consistently reduces the packet loss rate, even as the number of users increases, with improvements from 98.7% to 98.7%. TLS and fuzzy methods show a positive correlation between the number of users and packet loss, with improvements ranging from 2.6% to 3.4%.

On the other hand, the F-TLS, CLSTM and BRFO+IGNN methods demonstrate a negative correlation, with enhancements varying from 1.3% to 98.7%. These results highlight the critical need to fine-tune authentication mechanisms to reduce packet loss and improve network stability, tailored to specific user scenarios and requirements.

#### 4.4. Comparison of HO Decision Making Mechanisms

Table 5 presents a comparison of the results between the proposed and existing HO decision-making mechanisms. BRFO+IGNN consistently surpasses the benchmark models in all performance metrics, demonstrating higher accuracy, precision, recall, and F-measure values. Among the benchmark models, RF shows the poorest performance, with accuracy, precision, recall, and F-measure values of 54.65%, 53.24%, and 53.67%, respectively. DT and NB models exhibit moderate performance, with improvements in all metrics over RF. LR and SVM models show further performance gains, surpassing DT and NB in all metrics. The XGBoost mod-



**Fig. 9.** Comparison of proposed and benchmark HO authentication mechanisms.

el, a gradient boosting algorithm, demonstrates even better performance, outperforming LR and SVM on all metrics.

However, the proposed BRFO+IGNN stands out, achieving an impressive accuracy of 96.57%, a precision of 95.15%, a recall of 95.58%, and an F-measure of 95.37%.

Compared to the top-performing benchmark model (XG-Boost), BRFO+IGNN shows significant improvements in all areas, confirming its effectiveness in handover decision making (Fig. 9). This analysis highlights the superior performance, suggesting it has strong potential for practical deployment, enhancing both network reliability and performance.

## 5. Conclusions

This paper proposes a method that uses artificial intelligence (AI) techniques to improve handover triggering and management in wireless networks, specifically focusing on HO authentication. The approach applies Boruta random forest optimization (BRFO) to fine-tune the handover parameters, allowing to calculate optimal HO trigger points by adjusting the handover margins in order to strengthen supporting reliable authentication during vertical handovers. Additionally, an IGNN acts as the decision-making entity, predicting unwanted handovers and minimizing unnecessary handover events in small cell networks.

Performance of the proposed model is evaluated through simulation experiments which demonstrate its effectiveness in optimizing handover processes, authentication, and defense against potential attacks in 5G ultra-dense networks (UDNs).

The results show that BRFO + IGNN outperforms existing methods such as Event A3, FLDH, and FLDHDT in several key metrics.

## References

- [1] M.A. Adedoyin and O.E. Falowo, "Combination of Ultra-dense Networks and Other 5G Enabling Technologies: A Survey", *IEEE Access*, vol. 8, pp. 22893–22932, 2020 (<https://doi.org/10.1109/ACCESS.2020.2969980>).
- [2] R. Torre *et al.*, "Power Efficient Mobile Small Cell Placement for Network-coded Cooperation in UDNs", *Computer Networks*, vol. 201, art. no. 108559, 2021 (<https://doi.org/10.1016/j.comnet.2021.108559>).
- [3] V. Stoyanov *et al.*, "Ultra-dense Networks: Taxonomy and Key Performance Indicators", *Symmetry*, vol. 15, 2022 (<https://doi.org/10.3390/sym15010002>).
- [4] V. Stoyanov, A. Ivanov, and D. Mihaylova, "Flexible Access Network Design for Futuristic Mobile 5D Communications and Services", *AIP Conference Proceedings*, vol. 2570, art. no. 020009, 2022 (<https://doi.org/10.1063/5.0100110>).
- [5] T.M. Shami, D. Grace, A. Burr, and M.D. Zakaria, "Joint User-centric Clustering and Multi-cell Radio Resource Management in Coordinated Multipoint Joint Transmission", *Wireless Personal Communications*, vol. 124, pp. 2983–3011, 2022 (<https://doi.org/10.1007/s11277-022-09499-z>).
- [6] A. Mughees, M. Tahir, M.A. Sheikh, and A. Ahad, "Energy-efficient Ultra-dense 5G Networks: Recent Advances, Taxonomy and Future Research Directions", *IEEE Access*, vol. 9, pp. 147692–147716, 2021 (<https://doi.org/10.1109/ACCESS.2021.3123577>).
- [7] S. Sönmez, I. Shayea, S.A. Khan, and A. Alhammad, "Handover Management for Next-generation Wireless Networks: A Brief Overview", *2020 IEEE Microwave Theory and Techniques in Wireless Communications (MTTW)*, Riga, Latvia, 2020 (<https://doi.org/10.1109/MTTW51045.2020.9245065>).
- [8] D. Xenakis, N. Passas, L. Merakos, and C. Verikoukis, "Mobility Management for Femtocells in LTE-advanced: Key Aspects and Survey of Handover Decision Algorithms", *IEEE Communications Surveys & Tutorials*, vol. 16, pp. 64–91, 2013 (<https://doi.org/10.1109/SURV.2013.060313.00152>).
- [9] M. Emran *et al.*, "The Handover and Performance Analysis of LTE Network with Traditional and SDN Approaches", *Wireless Communications and Mobile Computing*, 2022 (<https://doi.org/10.1155/2022/7387737>).
- [10] R.A. Paropkari, A. Thantharate, and C. Beard, "Deep-mobility: A Deep Learning Approach for an Efficient and Reliable 5G Handover", *2022 International Conference on Wireless Communications Signal Processing and Networking (WiSPNET)*, Chennai, India, 2022 (<https://doi.org/10.1109/WiSPNET54241.2022.9767158>).
- [11] S.A. Khan, I. Shayea, M. Ergen, and H. Mohamad, "Handover Management over Dual Connectivity in 5G Technology with Future Ultra-dense Mobile Heterogeneous Networks: A Review", *Engineering Science and Technology, an International Journal*, vol. 35, 2022 (<https://doi.org/10.1016/j.jestch.2022.101172>).
- [12] A. Rammohan and D.K. R., "Revolutionizing Intelligent Transportation Systems with Cellular Vehicle-to-everything (C-V2X) Technology: Current Trends, Use Cases, Emerging Technologies, Standardization Bodies, Industry Analytics and Future Directions", *Vehicular Communications*, vol. 43, art. no. 100638, 2023 (<https://doi.org/10.1016/j.vehcom.2023.100638>).
- [13] S.S. Sefati and S. Halunga, "Ultra-reliability and Low-latency Communications on the Internet of Things Based on 5G Network: Literature Review, Classification, and Future Research View", *Transactions on Emerging Telecommunications Technologies*, vol. 34, art. no. e4770, 2023 (<https://doi.org/10.1002/ett.4770>).
- [14] Y. Ullah *et al.*, "A Survey on Handover and Mobility Management in 5G HetNets: Current State, Challenges, and Future Directions", *Sensors*, vol. 23, art. no. 5081, 2023 (<https://doi.org/10.3390/s23115081>).
- [15] R. Shafin *et al.*, "Artificial Intelligence-enabled Cellular Networks: A Critical Path to Beyond-5G and 6G", *IEEE Wireless Communications*, vol. 27, pp. 212–217, 2020 (<https://doi.org/10.1109/MWC.001.1900323>).
- [16] B. Ma, W. Guo, and J. Zhang, "A Survey of Online Data-driven Proactive 5G Network Optimization Using Machine Learning", *IEEE Access*, vol. 8, pp. 35606–35637, 2020 (<https://doi.org/10.1109/ACCESS.2020.2975004>).
- [17] C. Seródio *et al.*, "The 6G Ecosystem as Support for IoE and Private Networks: Vision, Requirements, and Challenges", *Future Internet*, vol. 15, art. no. 348, 2023 (<https://doi.org/10.3390/fi15110348>).
- [18] J. Wang, J. Liu, J. Li, and N. Kato, "Artificial Intelligence-assisted Network Slicing: Network Assurance and Service Provisioning in 6G", *IEEE Vehicular Technology Magazine*, vol. 18, pp. 49–58, 2023 (<https://doi.org/10.1109/MVT.2022.3228399>).

- [19] E. Esenogho, K. Djouani, and A.M. Kurien, "Integrating Artificial Intelligence, Internet of Things, and 5G for Next-generation Smart Grid: A Survey of Trends, Challenges, and Prospects", *IEEE Access*, vol. 10, pp. 4794–4831, 2022 (<https://doi.org/10.1109/ACCESS.2022.3140595>).
- [20] N. Haider, M.Z. Baig, and M. Imran, "Artificial Intelligence and Machine Learning in 5G Network Security: Opportunities, Advantages, and Future Research Trends", *arXiv*, 2020 (<https://doi.org/10.48550/arXiv.2007.04490>).
- [21] R. Karmakar, G. Kaddoum, and S. Chattopadhyay, "Mobility Management in 5G and Beyond: A Novel Smart Handover with Adaptive Time-to-trigger and Hysteresis Margin", *IEEE Transactions on Mobile Computing*, vol. 22, pp. 5995–6010, 2022 (<https://doi.org/10.1109/TMC.2022.3188212>).
- [22] W.S. Hwang, T.Y. Cheng, Y.J. Wu, and M.H. Cheng, "Adaptive Handover Decision Using Fuzzy Logic for 5G Ultra-dense Networks", *Electronics*, vol. 11, art. no. 3278, 2022 (<https://doi.org/10.3390/electronics11203278>).
- [23] Q. Liu *et al.*, "A Fuzzy-clustering Based Approach for MADM Handover in 5G Ultra-dense Networks", *Wireless Networks*, pp. 965–978, 2022 (<https://doi.org/10.1007/s11276-019-02130-3>).
- [24] V.O. Nyangaresi, A.J. Rodrigues, and S.O. Abeka, "Machine Learning Protocol for Secure 5G Handovers", *International Journal of Wireless Information Networks*, vol. 29, pp. 14–35, 2022 (<https://doi.org/10.1007/s10776-021-00547-2>).
- [25] S.V. Manjaragi and S.V. Saboji, "An Efficient Handover Authentication Mechanism Using Deep Learning in SDN-based 5G HetNets", *International Journal of Intelligent Engineering & Systems*, vol. 16, pp. 753–770, 2023 (<https://doi.org/10.22266/ijies2023.1231.63>).
- [26] J. Divakaran, A. Chakrapani, and K. Srihari, "Fuzzy Logic Based Handover Authentication in 5G Telecommunication Heterogeneous Networks", *Computer Systems Science and Engineering*, vol. 46, pp. 1141–1152, 2023 (<https://doi.org/10.32604/csse.2023.028050>).
- [27] V.O. Nyangaresi *et al.*, "Optimized Hysteresis Region Authenticated Handover for 5G HetNets", *Artificial Intelligence and Sustainable Computing: Proceedings of ICSIS CET 2021*, pp. 91–111, 2022 ([https://doi.org/10.1007/978-981-19-1653-3\\_9](https://doi.org/10.1007/978-981-19-1653-3_9)).
- [28] V.O. Nyangaresi, A.J. Rodrigues, S.O. Abeka, "ANN-FL Secure Handover Protocol for 5G and Beyond Networks", *Towards New e-Infrastructure and e-Services for Developing Countries: 12th EAI International Conference, AFRICOMM 2020*, pp. 99–118, 2020 ([https://doi.org/10.1007/978-3-030-70572-5\\_7](https://doi.org/10.1007/978-3-030-70572-5_7)).
- [29] A. Haghrah, J.M. Niya, and S. Ghaemi, "Handover Triggering Estimation Based on Fuzzy Logic for LTE-A/5G Networks with Ultra-dense Small Cells", *Soft Computing*, vol. 27, pp. 17333–17345, 2023 (<https://doi.org/10.1007/s00500-023-08063-6>).
- [30] A. Priyanka, P. Gauthamarayathirumal, and C. Chandrasekar, "Machine Learning Algorithms in Proactive Decision Making for Handover Management from 5G & Beyond 5G", *Egyptian Informatics Journal*, vol. 24, art. no. 100389, 2023 (<https://doi.org/10.1016/j.eij.2023.100389>).
- [31] M. Jamei *et al.*, "Developing Hybrid Data-intelligent Method Using Boruta-random Forest Optimizer for Simulation of Nitrate Distribution Pattern", *Agricultural Water Management*, vol. 270, art. no. 107715, 2022 (<https://doi.org/10.1016/j.agwat.2022.107715>).
- [32] A.A.M. Ahmed *et al.*, "LSTM Integrated with Boruta-random Forest Optimizer for Soil Moisture Estimation Under RCP4.5 and RCP8.5 Global Warming Situations", *Stochastic Environmental Research and Risk Assessment*, vol. 35, pp. 1851–1881, 2021 (<https://doi.org/10.1007/s00477-021-01969-3>).
- [33] S. Bera, S. Gupta, and A.S. Majumdar, "Device-independent Quantum Key Distribution Using Random Quantum States", *Quantum Information Processing*, vol. 22, art. no. 109, 2023 (<https://doi.org/10.1007/s11128-023-03852-2>).
- [34] J. Qadir *et al.*, "Mitigating Cyber Attacks in LoRaWAN via Lightweight Secure Key Management Scheme", *IEEE Access*, vol. 11, pp. 123456–123467, 2023 (<https://doi.org/10.1109/ACCESS.2023.3291420>).
- [35] V.P. Dwivedi *et al.*, "Benchmarking Graph Neural Networks", *arXiv*, 2023 (<https://doi.org/10.48550/arXiv.2003.00982>).
- [36] A. Tsitsulin, J. Palowitch, B. Perozzi, and E. Müller, "Graph Clustering with Graph Neural Networks", *arXiv*, 2023 (<https://doi.org/10.48550/arXiv.2006.16904>).
- [37] Y.-S. Chen, Y.-J. Chang, M.-J. Tsai, and J.-P. Sheu, "Fuzzy-logic-based Handover Algorithm for 5G Networks", *2021 IEEE Wireless Communications and Networking Conference (WCNC)*, Nanjing, China, 2021 (<https://doi.org/10.1109/WCNC49053.2021.9417298>).

---

### P. Rajesh, Research Scholar

Department of ECE

 <https://orcid.org/0000-0003-1217-5054>

E-mail: rajeshraipatnam@gmail.com

Vels Institute of Science, Technology & Advanced Studies,  
India

<https://vistas.ac.in>

### A. Vijayalakshmi, Ph.D.

Department of ECE

 <https://orcid.org/0000-0003-3594-6691>

E-mail: vijayalakshmi.se@velsuniv.ac.in

Vels Institute of Science, Technology & Advanced Studies,  
India

<https://vistas.ac.in>

### Ebenezer Abishek B., Ph.D.

Department of ECE

 <https://orcid.org/0000-0003-2908-7069>

E-mail: ebenezeraabishek@gmail.com

KCG College of Technology, India

<https://kcgcollege.ac.in>

# Synthesizing Wide-beam Array Patterns Using Phase-only Control and Trapezoidal Amplitudes for Satellite-based Internet Access

Zahraa Turki Hassan and Jafar Ramadhan Mohammed

*Ninevah University, Mosul, Iraq*

<https://doi.org/10.26636/jtit.2025.2.2067>

**Abstract** — Low Earth orbit satellite systems are capable of providing global Internet access due to their high downlink rate and low link budget. In such systems, wide beam array patterns are used to efficiently cover the required areas. In this paper, two efficient methods based on phase-only element excitation control for designing antenna arrays with required broaden beams are introduced. The first method, which is a simple algebraic approach, uses quadratic phase excitation while the amplitudes are chosen to be trapezoid. In the second method, an optimization algorithm is used to optimize the phase excitations of the array elements, while the amplitudes are still kept as a trapezoidal taper. Moreover, the use of trapezoidal-based amplitude excitations in both presented methods provides many desirable features compared to other conventional tapers. This is mainly due to the unique geometrical shape of the trapezoid taper, where the central coefficients have magnitudes of ones and the sided coefficients have decayed magnitudes. Simulations are presented to validate the proposed methods in which the beam width and maximum level of the radiated field were compared with those obtainable from the conventional standard Woodward-Lawson array.

**Keywords** — *antenna array, pattern synthesis, satellite application, wide beams*

## 1. Introduction

Future global Internet access requires low Earth orbit (LEO) satellite communication systems due to their ability to provide higher downlink capacity and a smaller link budget. In this application, the need for antenna arrays which have wide beam patterns are of a great interest. When the element excitation amplitudes and/or phases of an array are properly chosen, the shape of its radiation pattern can be achieved with required width. Thus, beamforming is an essential process in the antenna arrays of satellite communication systems to achieve a higher downlink capacity that is needed to succeed such systems.

Generally, the beam widths of the array patterns are inversely proportional to the apertures of the antenna array. Consequently, the beam widths become narrow for larger array apertures. Larger satellite arrays are essential to provide greater array

directivities and gains that help to achieve higher downlink capacity.

However, the satellite coverage areas decrease as the array aperture increases, and at the same time, widened beams are required to cover specific service areas.

The novelty of this paper is to introduce two new methods to efficiently synthesize widened beams for LEO satellite communication systems. In widened beams satellite applications, the flat-top level of the radiated fields is assumed to be uniform to ensure equal received power density within the coverage areas [1]–[3].

Many techniques have been proposed to synthesize wide beams [4]–[10]. In [11], [12], simple analytical techniques were introduced for the synthesization of widened beams. They are based on the quadratic and random selection of the phase-only element excitation control with fixed uniform amplitudes. However, these methods were not successful enough and their results were not promising when there were significant fluctuations in the obtained beams. Thus, the power density of these methods will not be equally received within the service areas. Moreover, a random selection of phase-only element excitations is not an effective approach, and it is mainly dependent on the trial-and-error process.

In [13]–[16], more powerful techniques based on evolutionary algorithms were used to synthesis widened beams.

In all of these aforementioned techniques, the element excitation amplitudes and phases are optimized jointly or separately to produce the required widened beams. Joint amplitude and phase excitation control methods are the most complicated [17], while separate control of amplitude or phase excitations is more simplified [18]. Phase-only control methods have been found to be more preferable than amplitude-only control [19].

In this paper, two new methods based on phase-only element excitation control are presented. The first method is based on a simple analytical approach where quadratic phase and trapezoid amplitude excitations are used to synthesis the widened beams. In the second method, a genetic optimization algorithm is used to optimize the elements of the phase excitations of the array instead of its quadratic values.

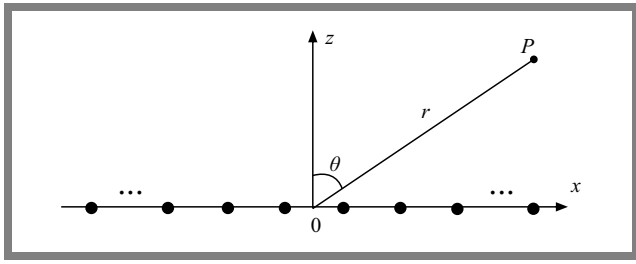


Fig. 1. Linear antenna array with isotropic elements.

Moreover, the use of fixed trapezoid amplitude excitations in the proposed two methods provides many desirable features in the radiation characteristics, such as uniform received power density across the served areas and low sidelobe levels.

## 2. The Proposed Method

The far-field radiation pattern of the linear antenna array with isotropic elements used can be represented mainly by the array factor. If the elements are placed symmetrically along the  $x$ -axis as shown in Fig. 1, then it is broadside array factor on the  $x - z$  plane is expressed as follows:

$$AF(\theta, x_n, I_n, \varphi_n) = \sum_{n=1}^N I_n e^{j \frac{2\pi}{\lambda} x_n \sin \theta + \varphi_n}, \quad (1)$$

where  $\lambda$  is wave length,  $N$  is the total number of array elements, and  $\theta$  is the direction of arrival angle from the broadside. From Eq. (1), it is clear that the factor of the array depends on three variable parameters that can be used to control the radiation patterns.

These design parameters are  $x_n$  which are the element locations  $x_n = [x_1, x_2, \dots, x_N]^T$ ,  $I_n$  which are the element excitation amplitudes  $I_n = [I_1, I_2, \dots, I_N]^T$ , and  $\varphi_n$  which are the element excitation phases  $\varphi_n = [\varphi_1, \varphi_2, \dots, \varphi_N]^T$ .

This three-dimensional-variables problem requires an efficient optimization algorithm to optimally determine element locations, amplitude excitations, and phase excitations. Usually, the locations were fixed to avoid iterative changes in the mechanical positions of the elements of the matrix. In this work, the locations were uniformly distributed at multiple integers of  $\frac{\lambda}{2}$ . Therefore, the elements are separated equally and evenly around the center of the array and Eq. (1) becomes:

$$AF(\theta, x_n, I_n, \varphi_n) = 2 \sum_{n=1}^{\frac{N}{2}} I_n e^{j \varphi_n} \cos \left[ \frac{2n-1}{2} \pi \sin \theta \right]. \quad (2)$$

The element excitation amplitudes  $I_n$  can be chosen according to the newly introduced trapezoidal taper window [20], [21]. The trapezoid taper is unique, and it has two different amplitude excitations. The uniform amplitudes with  $M$  elements in the center of the array, and two decayed amplitudes with  $N - M$  elements at the array sides. Thus, the  $I_n$  can be given

by [21]:

$$I_n = \begin{cases} \frac{n + \frac{N}{2}}{-\frac{M}{2} + \frac{N}{2}} - \frac{N}{2} \leq n \leq -\frac{M}{2} \\ 1 - \frac{M}{2} \leq n \leq \frac{M}{2} \\ \frac{N}{2} - n \\ \frac{N}{2} - \frac{M}{2} \leq n \leq \frac{N}{2} \end{cases}. \quad (3)$$

From Eq. (2), it is clear that there are only  $\frac{N}{2}$  variable excitation phases that must be determined instead of the original three-dimensional variables  $x_n, I_n, \varphi_n$  that were presented in Eq. (1). Furthermore, the  $\frac{N}{2}$  variable excitation phases are reduced to only  $\frac{N-M}{2}$  when using unit-amplitudes and zero-phasing with  $M$  central trapezoidal elements.

In the first proposed method, these  $\frac{N}{2}$  variable excitation phases are chosen according to the quadratic distribution, while in the second proposed method, they are taken as the optimization variables. Here in this research work, the peak sidelobe levels (SLL) along with the beam width constraints serve as the optimization objectives.

The objective function can be written as:

$$\begin{aligned} \text{Cost} = & \frac{\max(|AF|)_{\theta \in A}}{\max(|AF|)} \\ & + \max(FNBW - FNBW_D) \\ & + \sum_{i=1}^I \frac{\max(|AF(\theta_{null}^i)|)}{\max(|AF|)}, \end{aligned} \quad (4)$$

where  $A$  is the sidelobe area which is located outside of the main beam.

The first term on the right side of the equation is the normalized peak sidelobe level, the second term is the first-null-to-null-beamwidth where  $FNBW_D$  is the desired one. The third term is the required null directions toward the interfering signals, where  $I$  is the total number of the required null placement.

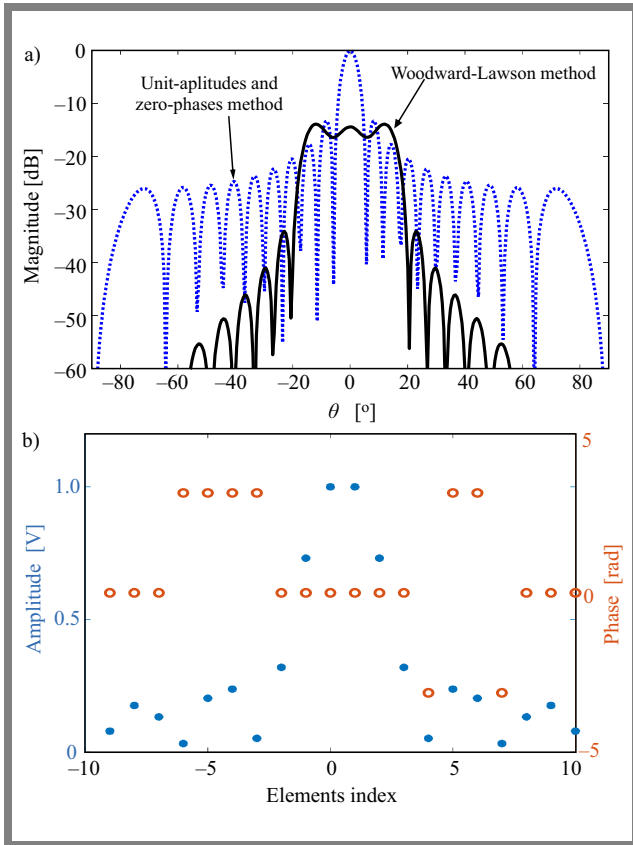
Then, the optimization problem can be expressed as:

$$\left. \begin{aligned} & \text{find } \varphi_n = [\varphi_1, \varphi_2, \dots, \varphi_N]^T \\ & \min(\text{cost}) \\ & \text{subject to } \frac{-\pi}{2} \leq \varphi_n \leq \frac{\pi}{2} \text{ for } n = 1, 2, \dots, \frac{N}{2} \end{aligned} \right\} \quad (5)$$

## 3. Simulation Results

Consider a linear symmetric antenna array that has  $N = 20$  elements with an interelement spacing of  $\frac{\lambda}{2}$ . In the following simulations, the optimization parameters of the genetic algorithm are chosen referring to [20], [21]. The used trapezoid taper for element excitation amplitudes has  $M = 4$  elements with unit-amplitudes and zero-phases at the center, while the remaining phases that need to be determined is  $N - M = 16$  elements which they are located at both array sides.

Comparisons are made with other non-optimization methods by using the same example and appropriate parameters setting.



**Fig. 2.** Beam patterns a) and their corresponding Woodward-Lawson amplitudes and phases b) for  $FNBD = 40^\circ$ .

For the classical unit-amplitude and quadratic-phases method,  $I_n = 1$  and

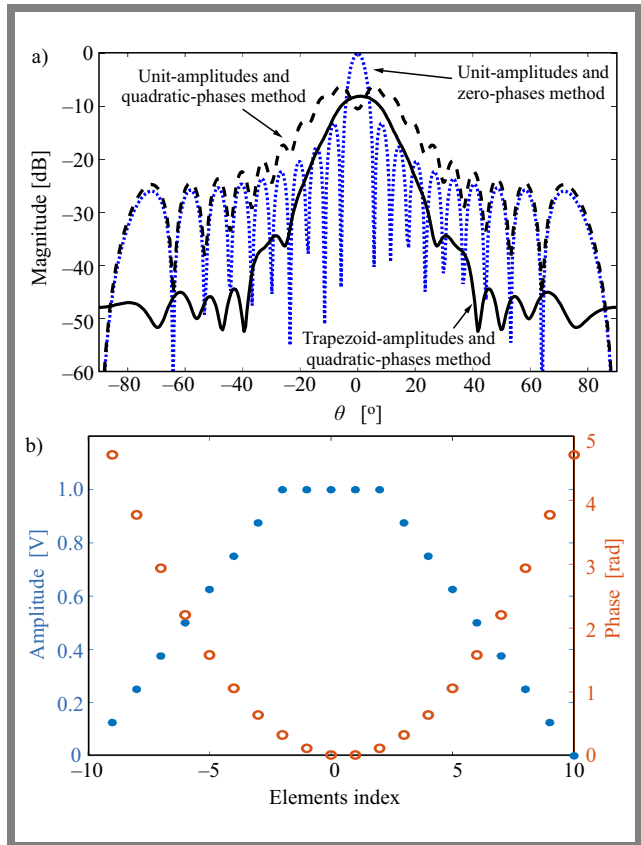
$$\varphi_n = 4 \varphi_{max} \left( \frac{x_n}{A_L} \right)^2 \text{ for } n = 1, 2, \dots, \frac{N}{2}.$$

Here,  $\varphi_{max} = 3\pi$  which is the maximum allowed phase value at the two end elements and  $x_n$  is the location of the  $n$ -th element along the array aperture length  $A_L$ .

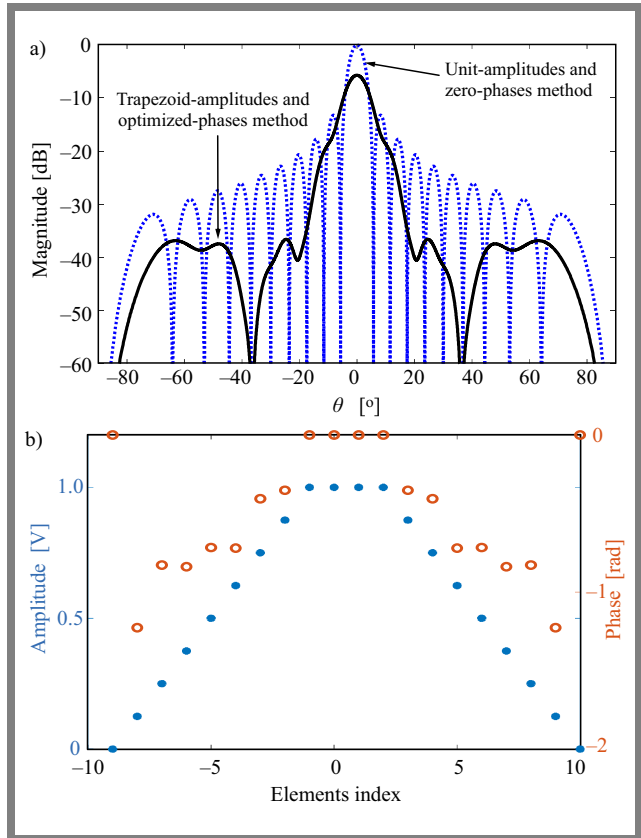
For the trapezoid-amplitude and quadratic-phases method,  $I_n$  values are computed according to Eq. (3),  $M = 4$ , and  $\varphi_n$  values are as mentioned in above. For the standard Woodward-Lawson method, the values of  $I_n$  and  $\varphi_n$  are chosen according to Woodward taper [22]. These aforementioned methods were studied and compared under different values of beam widths.

In the first example, the required beam width of the designed linear array is assumed to be equal to  $FNBD = 40^\circ$ . Figures 2–4 show the required amplitudes and phases of the element excitations along with their corresponding beam patterns for the Woodward-Lawson method, unit amplitudes and quadratic-phases method, trapezoid-amplitudes and quadratic phase method, and trapezoid-amplitudes and optimized phase method. From these three figures, it can be seen that the required widened beams have been achieved at the cost of lower directivities in the broadside directions.

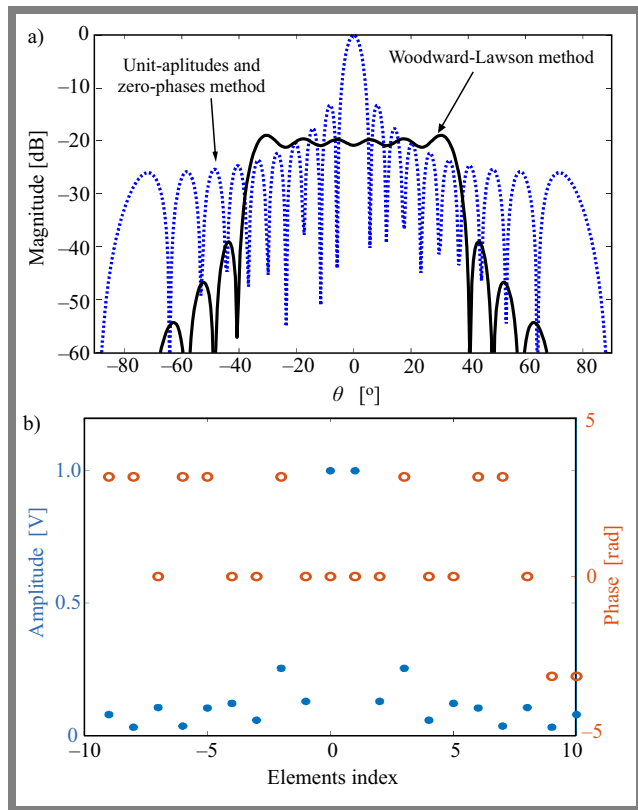
The level of the main beam, for the method of unit amplitude and zero phases, was normalized to 0 dB, while the beam patterns of other methods were normalized to the same value. As can be seen in Fig. 4, a minimum drop at  $\theta = 0^\circ$  occurs for



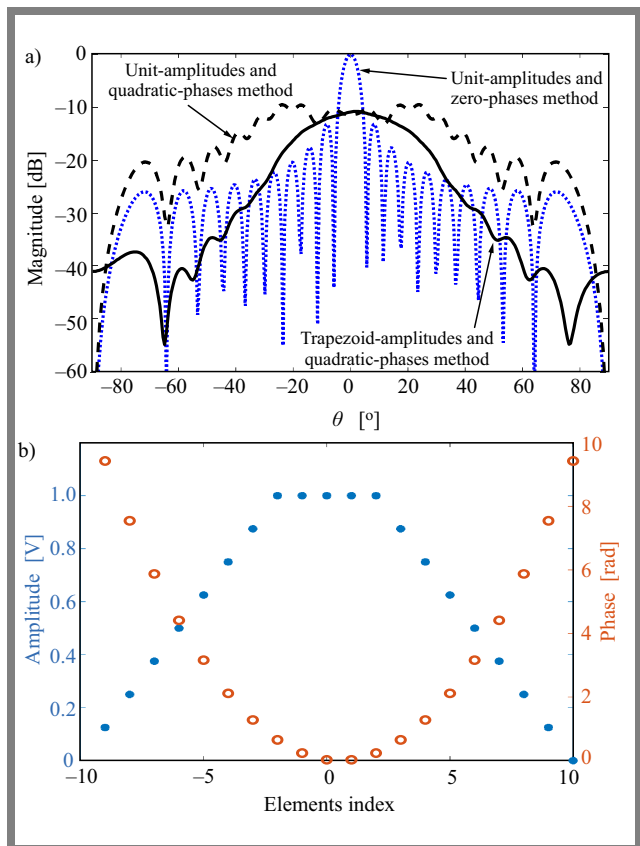
**Fig. 3.** Beam patterns a) and their corresponding trapezoid amplitudes and quadratic-phases b) for  $\varphi_{max} = 3\pi$ .



**Fig. 4.** Beam patterns a) and its corresponding trapezoid amplitudes and optimized phases b) for  $FNBD = 40^\circ$ .



**Fig. 5.** Beam patterns a) and its corresponding Woodward-Lawson amplitudes and phases b) for  $FNBD = 80^\circ$ .



**Fig. 6.** Beam patterns a) and their corresponding trapezoid amplitudes and quadratic-phases b) for  $\varphi_{max} = 6\pi$ .

the proposed method of trapezoid amplitudes and optimized-phases. These radiation characteristics were numerically computed and compared in the following example.

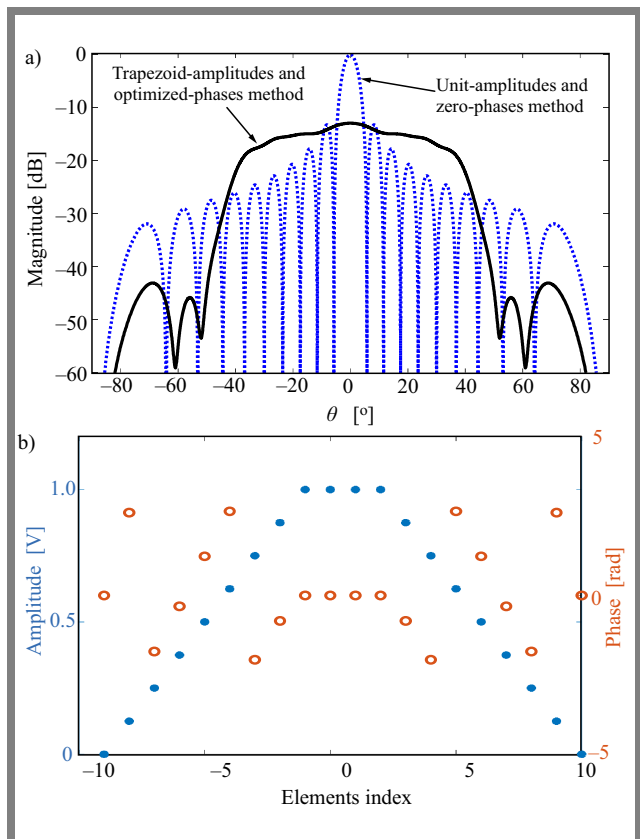
In the second example, the performances in terms of array complexity (i.e. needed RF components such as variable attenuators and phase shifters), first-null-to-null beam width (FNBD), directivity, aperture's taper efficiency, and peak sidelobe level (SLL) of these aforementioned methods were compared as shown in Tab. 1.

In the next example, the required first null-to-null beam width is assumed to be equal to  $FNBD = 80^\circ$  and its results are shown in Figs. 5–7 and Tab. 2.

From Figs. 5–7, it can be seen that the maximum levels of the widened beams further drop as the  $FNBD$  are increased. The directivities were also significantly reduced with compared to the classic method of unit amplitudes and zero-phases.

However, the proposed methods still provide a lower reduction with compared to that of Woodward-Lawson method. This is evident when comparing the magnitudes of Fig. 7 with that of Fig. 5 at  $\theta = 0^\circ$ .

Finally, the proposed method of trapezoid amplitudes and optimized phases is extended to include the two-dimensional rectangular planar array instead of its linear counterpart. The results are shown in Fig. 8 for the required null-to-null beam width  $FNBD = 80^\circ$  and an array size of  $N \times N = 20 \times 20$  elements.



**Fig. 7.** Beam patterns a) and its corresponding trapezoid amplitudes and optimized phases b) for  $FNBD = 80^\circ$ .

**Tab. 1.** Performance measures of the methods tested methods for the required  $FNBW_D = 40^\circ$ .

Method	Feeding network complexity	First null-to-null beam width (FNBW) [°]	Element excitations		Directivity [dB]	Aperture's taper efficiency	Peak sidelobe level [dB]
			Amp.	Phase [°]			
Classical unit-amplitudes and zero-phases	Zero transducer and zero phase shifters	The FNBW value is $11.42^\circ$ . This is narrower than the required one	1	0	26.04	1	-13.2
			1	0			
			1	0			
			1	0			
			1	0			
			1	0			
			1	0			
			1	0			
			1	0			
Woodward-Lawson method [22]	$N$ transducers and $N$ phase shifters	The FNBW value is $40^\circ$ . This is the same as the required FNBW	0.07	0	13.28	0.23	-34.4
			0.17	0			
			0.13	0			
			0.03	180			
			0.20	180			
			0.23	180			
			0.05	180			
			0.31	0			
			0.73	0			
			1.00	0			
Unit-amplitudes and quadratic-phases method [11]	$N/2$ phase shifters	The FNBW value is $34^\circ$ . This is narrower than the required one	1	270.00	13.76	0.24	-13.5
			1	216.14			
			1	168.28			
			1	126.39			
			1	90.49			
			1	60.58			
			1	36.64			
			1	18.69			
			1	6.73			
1	0.74						
Proposed trapezoid-amplitude and quadratic-phases method	$(N - M)/2$ transducers and $N/2$ phase shifters	The obtained FNBW value is $51.4^\circ$ . This is wider than the required one and mainly depending on the value of $\varphi_{max}$	0.12	270.00	16.38	0.33	-34.3
			0.25	216.14			
			0.37	168.28			
			0.50	126.39			
			0.62	90.49			
			0.75	60.58			
			0.87	36.64			
			1.00	18.69			
			1.00	6.73			
			1.00	0.74			
Proposed trapezoid-amplitudes and optimized-phases method	$(N - M)/2$ transducers and $(N - M)/2$ phase shifters	The FNBW value is $40^\circ$ . This is the same as the required FNBW	0.00	0.00	18.37	0.64	-36.7
			0.12	-70.27			
			0.25	-47.43			
			0.37	-48.02			
			0.50	-41.05			
			0.62	-41.27			
			0.75	-23.26			
			0.87	-20.16			
			1.00	0.00			
			1.00	0.00			

From Fig. 8, it can be seen that the amplitude excitations are exactly as the trapezoid taper, where it has three unit amplitudes on both sides of the array center and then decaying toward the array ends in four array quadrants. While phase excitations are optimized according to the cost function to

obtain wider beam that extends from  $-40^\circ$  (i.e., corresponds to a value of  $-0.64$ ) up to  $40^\circ$  on both  $u - v$  planes. The magnitude of the resultant array pattern is wide enough as required at the cost of little reduction in antenna array directivity.

**Tab. 2.** Performance measures of the methods tested methods for the required  $FNBW_D = 80^\circ$ .

Method	Feeding network complexity	First null-to-null beam width (FNBW) [°]	Element excitations		Directivity [dB]	Aperture's taper efficiency	Peak sidelobe level [dB]
			Amp.	Phase [°]			
Classical unit-amplitudes and zero-phases	Zero transducer and zero phase shifters	The FNBW value is $11.42^\circ$ . This is narrower than the required one	1	0	26.04	1	-13.2
			1	0			
			1	0			
			1	0			
			1	0			
			1	0			
			1	0			
			1	0			
			1	0			
Woodward-Lawson method [22]	$N$ transducers and $N$ phase shifters	The FNBW value is $80^\circ$ . This is the same as the required FNBW	0.07	180	7.13	0.11	-40.0
			0.03	180			
			0.10	0			
			0.03	180			
			0.10	180			
			0.12	0			
			0.05	0			
			0.25	180			
			0.12	0			
1.00	0						
Unit-amplitudes and quadratic-phases method [11]	$N/2$ phase shifters	The obtained FNBW value is $57.2^\circ$ . This is narrower than the required one and mainly depending on the value of $\varphi_{max}$	1	540	6.97	0.11	-10.0
			1	432			
			1	336			
			1	252			
			1	180			
			1	121			
			1	73			
			1	37			
			1	13			
1	1.4						
Proposed trapezoid-amplitude and quadratic-phases method	$(N - M)/2$ transducers and $N/2$ phase shifters	The obtained FNBW value is $75.4^\circ$ . This is narrower than the required one and mainly depending on the value of $\varphi_{max}$	0.125	540	10.69	0.17	-30.0
			0.250	432			
			0.375	336			
			0.500	252			
			0.625	180			
			0.750	121			
			0.875	73			
			1.000	37			
			1.000	13			
1.000	1.4						
Proposed trapezoid-amplitudes and optimized-phases method	$(N - M)/2$ transducers and $(N - M)/2$ phase shifters	The FNBW value is $80^\circ$ . This is the same as the required FNBW	0.125	0.00	11.63	0.25	-43.0
			0.250	-23.31			
			0.375	69.75			
			0.500	76.06			
			0.625	-1.35			
			0.750	50.77			
			0.875	112.96			
			1.000	44.02			
			1.000	0.00			
1.000	0.00						

### 4. Conclusions

It has been shown that the wide beam patterns with required first null-to-null beam widths and low sidelobe levels can be efficiently generated by controlling the phase-only excitations of the array elements either algebraically by using a simple quadratic phase method or optimally by using a genetic op-

timization method. In both methods, the amplitudes were constraint as a trapezoidal taper.

Results of using the first proposed method of trapezoid-amplitudes and quadratic-phases showed significant improvements in terms of reducing the sidelobe level, improving the taper's efficiency, and enhancing the array directivity compared to other conventional methods. For the case of

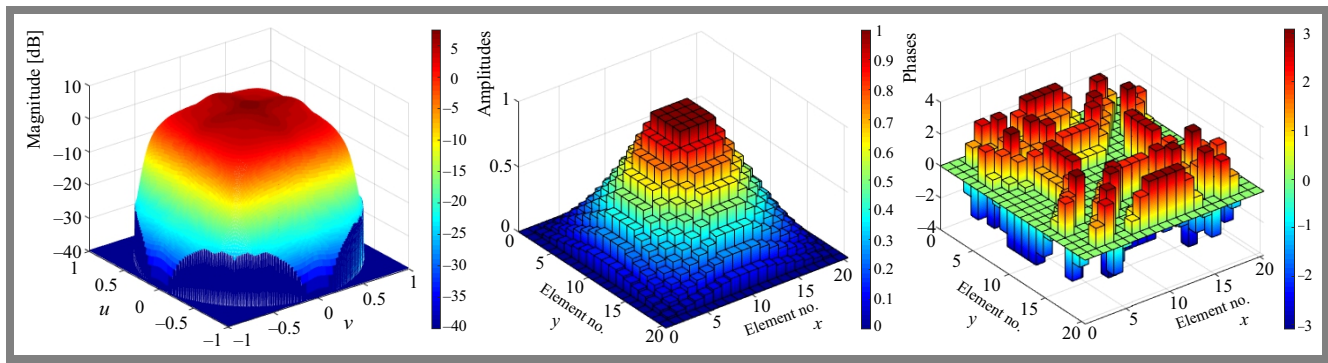


Fig. 8. Beam pattern and its corresponding trapezoid amplitudes and optimized phases for planar array with  $FNBW_D = 80^\circ$ .

generating beam width with  $FNBW = 40^\circ$ , the SLL improvement is more than  $-20$  dB, directivity improvement is more than  $2.6$  dB, while for the case of  $FNBW = 80^\circ$ , these improvements were  $-24$  dB and  $4$  dB respectively.

The results of using the second proposed method of trapezoid amplitudes and optimized phases showed significant improvements in the radiation characteristics. Thus, the two proposed methods are the way of using the widened beams in LEO satellite systems to successfully provide global Internet access applications.

## References

- [1] W. Lin, Y. Wu, and B. Su, "Broadened-beam Uniform Rectangular Array Coefficient Design in LEO SatComs Under Quality of Service and Constant Modulus Constraints", *IEEE Access*, vol. 12, pp. 184909–184928, 2024 (<https://doi.org/10.1109/ACCESS.2024.3513330>).
- [2] J.R. Mohammed "Optimizing Multiple Beam Patterns for 5G mmWave Phased Array Applications", *Periodica Polytechnica Electrical Engineering and Computer Science*, vol. 67, pp. 369–375, 2023 (<https://doi.org/10.3311/PPee.22111>).
- [3] S. Dai, M. Li, Q.H. Abbasi, and M.A. Imran, "A Zero Placement Algorithm for Synthesis of Flat Top Beam Pattern with Low Sidelobe Level", *IEEE Access*, vol. 8, pp. 225935–225944, 2020 (<https://doi.org/10.1109/ACCESS.2020.3045287>).
- [4] K.H. Sayidmarie and J.R. Mohammed, "Performance of a Wide Angle and Wideband Nulling Method for Phased Arrays", *Progress in Electromagnetics Research M*, vol. 33, pp. 239–249, 2013 (<https://doi.org/10.2528/PIERM13100603>).
- [5] J.R. Mohammed, A.J. Abdulkadeer, and R. Hamdan, "Array Pattern Recovery under Amplitude Excitation Errors Using Clustered Elements", *Progress In Electromagnetics Research M*, vol. 98, pp. 183–192, 2020 (<https://doi.org/10.2528/PIERM20101906>).
- [6] D. Schwartzman, R.D. Palmer, M. Herndon, and M.B. Yearly, "Enhanced Weather Surveillance Capabilities with Multiple Simultaneous Transmit Beams", *IEEE Transactions on Radar Systems*, vol. 3, pp. 272–289, 2025 (<https://doi.org/10.1109/TRS.2025.3527882>).
- [7] B.K. Daniel and A.L. Anderson, "Phase-only Beam Broadening of Contiguous Uniform Subarrayed Arrays", *IEEE Transactions on Aerospace and Electronic Systems*, vol. 56, pp. 4001–4013, 2020 (<https://doi.org/10.1109/TAES.2020.2987115>).
- [8] J.R. Mohammed and K.H. Sayidmarie, "Sensitivity of the Adaptive Nulling to Random Errors in Amplitude and Phase Excitations in Array Elements", *International Journal of Telecommunication, Electronics, and Computer Engineering*, vol. 10, pp. 51–56, 2018 (<https://jtec.utem.edu.my/jtec/article/view/2023>).
- [9] J.R. Mohammed, "Phased Array Antenna with Ultra-low Sidelobes", *Electronics Letters*, vol. 49, pp. 1055–1056, 2013 (<https://doi.org/10.1049/el.2013.1642>).
- [10] Y. Su *et al.*, "Broadband LEO Satellite Communications: Architectures and Key Technologies", *IEEE Wireless Communications*, vol. 26, pp. 55–61, 2019 (<https://doi.org/10.1109/MWC.2019.1800299>).
- [11] K.H. Sayidmarie and Q.H. Sultan, "Synthesis of Wide Beam Array Patterns Using Quadratic-phase Excitations", *International Journal of Electromagnetics and Applications*, vol. 3, pp. 127–135, 2013 (<https://doi.org/10.5923/j.ijea.20130306.01>).
- [12] K.H. Sayidmarie and Q.H. Sultan, "Synthesis of Wide Beam Array Patterns Using Random Phase Weights", *International Conference on Electrical, Communication, Computer, Power and Control Engineering (ICECCPCE)*, Mosul, Iraq, 2013 (<https://doi.org/10.1109/ICECCPCE.2013.6998734>).
- [13] M.C. Viganò *et al.*, "Sparse Antenna Array for Earth-coverage Satellite Applications", *Fourth European Conference on Antennas and Propagation*, Barcelona, Spain, 2010 (<https://ieeexplore.ieee.org/document/5504930>).
- [14] B. Samantaray, K.K. Das, and J.S. Roy, "Designing Smart Antennas Using Machine Learning Algorithms", *Journal of Telecommunications and Information Technology*, no. 4, pp. 46–52, 2023 (<https://doi.org/10.26636/jtit.2023.4.1329>).
- [15] Z. Qu, "LEO Satellite Constellation for Internet of Things", *IEEE Access*, vol. 5, pp. 18391–18401, 2017 (<https://doi.org/10.1109/ACCESS.2017.2735988>).
- [16] P. Callaghan and P.R. Young, "Beam- and Band-width Broadening of Intelligent Reflecting Surfaces Using Elliptical Phase Distribution", *IEEE Transactions on Antennas and Propagation*, vol. 70, pp. 8825–8832, 2022 (<https://doi.org/10.1109/TAP.2022.3199451>).
- [17] X. He, Y. Zhang, Z.H. Jiang, and W. Hong, "A Generalized Flat-topped Beam Synthesis Approach for Uniform Linear Array With Arbitrary Beam Directions", *IEEE Open Journal of Antennas and Propagation*, vol. 3, pp. 709–721, 2022 (<https://doi.org/10.1109/OJAP.2022.3184823>).
- [18] A.J. Abdulkadeer, J.R. Mohammed, and R. Hamdan, "Phase-only Nulling with Limited Number of Controllable Elements", *Progress In Electromagnetics Research C*, vol. 99, pp. 167–178, 2020 (<https://doi.org/10.2528/PIERC20010203>).
- [19] A.F. Morabito, A. Massa, P. Rocca, and T. Isernia, "An Effective Approach to the Synthesis of Phase-only Reconfigurable Linear Arrays", *IEEE Transactions on Antennas and Propagation*, vol. 60, pp. 3622–3631, 2012 (<https://doi.org/10.1109/TAP.2012.2201099>).
- [20] J.R. Mohammed, "Array Pattern Synthesis Using a New Adaptive Trapezoid Window Function for Sidelobe Suppression and Nulls Control", *Progress In Electromagnetics Research M*, vol. 129, pp.83–90, 2024 (<https://doi.org/10.2528/PIERM24083102>).
- [21] J.R. Mohammed, "An Optimized Phase-only Trapezoid Taper Window for Array Pattern Reshaping", *Progress in Electromagnetics Research C*, vol. 152, pp. 245–251, 2025 (<https://doi.org/10.2528/PIERC24122101>).
- [22] R.S. Elliott, "More on the Woodward-Lawson Method", *IEEE Antennas and Propagation Society Newsletter*, vol. 30, pp. 28–29, 1988 (<https://doi.org/10.1109/MAP.1988.6086101>).

**Zahraa Turki Hassan, Student**

College of Electronics Engineering

 <https://orcid.org/0009-0009-3956-7892>

E-mail: zahraa.hassan2017@stu.uoninevah.edu.iq

Ninevah University, Mosul, Iraq

<https://uoninevah.edu.iq>

**Jafar Ramadhan Mohammed, Professor**

College of Electronics Engineering

 <https://orcid.org/0000-0002-8278-6013>

E-mail: jafar.mohammed@uoninevah.edu.iq

Ninevah University, Mosul, Iraq

<https://uoninevah.edu.iq>

# Analysis of Pyramidal Microwave Absorbers for Enhanced Performance in 1 – 10 GHz Frequency Range

Aya Raad Thanoon and Khalil H. Sayidmarie

*Ninevah University, Mosul, Iraq*

<https://doi.org/10.26636/jtit.2025.2.2092>

**Abstract** — One of the main applications of microwave absorbers is in anechoic chambers, where the walls are lined with pyramidal foam impregnated with a lossy material. This paper investigates the impact that various design parameters of pyramidal microwave absorbers exert on their performance, with the aim of finding the best design values that ensure better operational properties. Typical pyramid absorbers were investigated by conducting simulations with the use of the CST Microwave Suite simulator, across the frequency range of 1 – 10 GHz, at various angles of the incident wave. The investigations also considered absorbers backed by conducting plates that are used in shielded anechoic chambers. The study shows that higher permittivity leads to higher reflection, while increased loss tangent improves absorption, and the same applies to magnetic materials. Larger pyramid heights lead to lower reflection, but only in the case of thicker absorbers. A pyramidal absorber with the height of 16 cm, designed using lossy material with permittivity and permeability of 1.5 and loss tangent of 0.5 achieved a reflection coefficient that was lower than  $-60$  dB for frequencies between 3 and 10 GHz. The results are useful in designing absorbers relying on materials that offer only dielectric or magnetic properties, or that combine both of them to achieve enhanced performance.

**Keywords** — *anechoic chambers, dielectric loss, microwave absorber, reflectivity*

## 1. Introduction

Microwave absorbers are commonly used in many applications, such as radar, military defense systems, and compliance testing of electronic devices. Their primary aim is to reduce or eliminate the reflection, transmission, and scattering of microwave radiation, thereby minimizing interference or reducing the risk of detection in many sensitive applications.

However, newly built RF and microwave devices and systems must be experimentally tested to assess their performance and evaluate their radiation leakage. Testing of these devices must be conducted either in an open area that is free of any surrounding objects and other interference generating devices, or inside an anechoic chamber. The former solution is rather expensive and may not always be available due to weather conditions, while anechoic chambers are readily available.

Currently, commercially available absorbers are mostly constructed of polyurethane and polystyrene foam impregnated with materials that absorb electromagnetic energy. Many of

these absorbers use carbon or its derivatives as a lossy material. Other types rely on ferrites, polymers, and lossy dielectric materials [1].

The efforts to improve the performance of microwave absorbers can be divided into the following stages:

- selecting a suitable and cheap lossy material,
- choosing a proper profile of the absorber's surface and its thickness,
- determining the frequency range across which acceptable reflection levels are achieved,
- deciding on the weight of the absorber.

Certain agricultural leftovers, such as coconut shells, banana peels, sugar cane, water hyacinth, and other byproducts have been suggested as alternatives to foam in absorber construction. The author of [2] mixed rubber tire dust and rice husk in three different ratios (50:50, 25:75, and 75:25) to build a microwave absorber. For a 15 cm thick absorber, the average reflection coefficient obtained equaled  $-22.03$  dB,  $-21.54$  dB, and  $-32.51$  dB, respectively, in the frequency range of 7 to 13 GHz.

In [3], sugar cane bagasse was used to build pyramidal microwave absorbers operating in the frequency range of 0.1 – 20 GHz, with an average reflection of  $-44$  dB. The authors of [4] used rice husk and a mixture of rice husk and coal to construct pyramid microwave absorbent structures which produced mean reflections of  $-30$  dB and  $-40$  dB, respectively. Furthermore, the authors of [5] proposed the use of water hyacinth to fabricate microwave absorbers, achieving reflectivity values ranging from  $-30$  to  $-10$  dB for a 13 cm thick absorber. In [6], a mixture of two materials, such as polyurethane and carbon, was proposed for use in the frequency range of 1 – 10 GHz. The fabricated absorber, being 30 cm thick, achieved a reflection coefficient of  $-30$  dB at 3 GHz.

The shape of the absorber is another important factor in determining the level of the reflection, since it forms the boundary between the air and the lossy material of the absorber where the electromagnetic wave is incident. At this interface, the law of reflection is derived by applying the boundary conditions. Pyramidal absorbers were extensively used, as they usually offer a smooth and gradual change from the air to the absorbing material. The performance of pyramidal and truncated

pyramidal absorbers with relative permittivity values of 2.5, 2.9 and 3.3, made of solid, hollow or coated materials were investigated in [7]. The collected data, which were presented as the best or average values, demonstrated that slightly higher permittivity corresponded to somewhat higher reflection. This is because as the absorber's permittivity increases along with the increase in the reflection of the incident wave from its surface.

In [8], a triangular pyramid, an isosceles pyramid, and an equilateral pyramid were investigated, with rice husks as the absorbing medium, at frequencies of 1–20 GHz, and the results obtained results showed that the shape of the pyramid base could affect the performance of the entire absorber. Paper [9] examined a hexagonal pyramid made from banana leaves, rice husks, and rice straws, at frequencies ranging from 0.01 to 20 GHz. The results obtained for 13 cm thick absorbers showed reflection levels of –35 dB, –39 dB, and –38 dB, respectively. The same agricultural residues were also applied to truncated hexagonal pyramids which offered lower reflection coefficients of –35 dB, –35 dB, and –37 dB, respectively.

Another shape, such as wall tiles made of kenaf and coconut coir was investigated in [10] as an absorber for the frequency range of 1 to 12 GHz. The kenaf had the highest absorption in the C band and the coconut coir showed better results in the X band. In [11], carbon from biomaterials was used and its reflectance was typically lower than –25 dB.

Although the pyramidal shape offers a gradual introduction of the absorbing material, planar multilayer microwave absorbers have also been employed. Their configuration comprises several layers (all characterized by different thicknesses) of various types of materials that differ in terms of their electrical and magnetic properties. Such a solution is effective in multiple frequency bands and offers better performance provided it is well optimized at the design stage [12].

Metamaterial surfaces can be designed to work as electromagnetic wave absorbers characterized by decent efficiency, and their geometric structures may be used to improve the performance of traditional absorbers. A metamaterial absorber with three layers of different resistive films and a central via was presented in [13]. A 90% absorption rate was obtained across the frequency range of 3.2 to 35.5 GHz or 167% relative bandwidth. However, the level of the reflection coefficient was not provided, as the emphasis was placed on bandwidth and absorption. In [14], lossy carbon paint was applied to a frequency-selective surface (FSS) to enhance the absorption rate. While this absorber is very thin, the bandwidth achieved was limited.

Hollow pyramidal absorbers (HPA) can offer an additional degree of freedom in optimizing broadband performance by furnishing hollow sections of various sizes. Thus, they can absorb waves within a wider range of frequencies. The authors of [15] designed a slotted HPA that contains slots of isosceles triangles, which achieved absorptivity of –26.32 dB, in the L, S, C, and X frequency bands. [16] introduced a triangular-slotted HPA using the Sierpinski principle. The conclusion was that an increase in the number of smaller slots resulted

in the highest level of absorption, the most stable impedance, and the widest bandwidth.

The former studies focused on the suitability of certain materials, with less emphasis placed on the preferred characteristics of the material that are required in order to ensure better performance. Moreover, the pyramidal shape was extensively used, but the dimensions that yield the lowest reflection remained unclear.

This study explores the performance of the pyramid absorber and the impact that its height, base size, permittivity, permeability,  $\tan \delta$ , and angle of incidence exert on the electromagnetic wave (EM) applied. The aim is to find a combination of dimensions and material properties capable of achieving better performance. Furthermore, it investigates the case where the absorber is backed by a conducting plate that is employed in shielded anechoic chambers.

The paper is organized as follows. Section 2 analyzes the reflection of the EM wave from lossy dielectric surfaces. Section 3 presents the modeling of the pyramidal absorber using CST software, while Section 4 explores the effect of the various absorber parameters. Section 5 investigated the case of absorbers supported by a conducting plate, while Section 6 contains considerations on the type of plane surface material. Section 7 provides a comparison with other published works. The conclusions drawn are presented in Section 8.

## 2. Analysis of Reflection at the Air-dielectric Boundary

The impact that microwave absorbers exert on incident EM waves depends on the reflection coefficient, properties of the dielectric material of the absorber and its geometry. Dielectric and magnetic properties of the absorber are important factors determining absorption performance. The relative complex permittivity is typically written as:

$$\varepsilon_r = \varepsilon'_r - j\varepsilon''_r, \quad (1)$$

where  $\varepsilon'_r$  represents the material's energy storage capacity and  $\varepsilon''_r$  indicates the losses suffered due to the electric field. The electric loss tangent is the ratio between the imaginary and real parts of permittivity, and is expressed as:

$$\tan \delta_d = \frac{\varepsilon''_r}{\varepsilon'_r}. \quad (2)$$

Similarly, the relative permeability of the complex is:

$$\mu_r = \mu'_r - j\mu''_r, \quad (3)$$

where  $\mu'_r$  represents the material's ability to store magnetic energy and  $\mu''_r$  indicates the losses suffered due to the magnetic field. The magnetic loss tangent is the ratio between the imaginary and real parts of permeability, and is expressed as:

$$\tan \delta_m = \frac{\mu''_r}{\mu'_r}. \quad (4)$$

The three basic scenarios for the reflection of an EM wave that arrives at a specific material from the air are shown in Fig. 1. The reflection coefficient for a normally incident wave

on a planar medium can be given by [17]:

$$\Gamma = \frac{\eta - \eta_0}{\eta + \eta_0}, \quad (5)$$

where:

$$\eta_0 = \sqrt{\frac{\mu_0}{\varepsilon_0}} \quad \text{and} \quad \eta = \sqrt{\frac{j\omega\mu}{\sigma + j\omega\varepsilon}}$$

are the intrinsic impedances of air and the absorber material, respectively. By substituting these in Eq. (5), it can be shown that for a general material, the reflection coefficient can be written as:

$$\Gamma = \frac{1 - \sqrt{\frac{\varepsilon_r}{\mu_r} - \frac{j\sigma}{\omega\varepsilon_0\mu_r}}}{1 + \sqrt{\frac{\varepsilon_r}{\mu_r} - \frac{j\sigma}{\omega\varepsilon_0\mu_r}}}, \quad (6)$$

which for lossless material ( $\sigma = 0$ ) reduces to:

$$\Gamma = \frac{1 - \sqrt{\frac{\varepsilon_r}{\mu_r}}}{1 + \sqrt{\frac{\varepsilon_r}{\mu_r}}}, \quad (7)$$

where  $\varepsilon_r$  and  $\mu_r$  are the relative complex dielectric constant of the complex and the relative permeability of the complex given by Eq. (1) and Eq. (3), respectively. After substitution for  $\varepsilon_r$  and using Eq. (2), the reflection coefficient for a nonmagnetic material  $\Gamma$  can be expressed using two forms:

$$\Gamma = \frac{1 - \sqrt{\varepsilon_r' - j\varepsilon_r''}}{1 + \sqrt{\varepsilon_r' - j\varepsilon_r''}}, \quad (8)$$

$$\Gamma = \frac{1 - \sqrt{\varepsilon_r' \sqrt{1 - j \tan \delta_d}}}{1 + \sqrt{\varepsilon_r' \sqrt{1 - j \tan \delta_d}}}. \quad (9)$$

For the absorber limiting case when the dielectric material of the absorber is lossless ( $\varepsilon_r'' = 0$ , or  $\tan \delta_d = 0$ ), Eq. (8) and (9) will be:

$$\Gamma = \frac{1 - \sqrt{\varepsilon_r'}}{1 + \sqrt{\varepsilon_r'}}. \quad (10)$$

Equation (10) indicates that a higher relative permittivity  $\varepsilon_r'$  leads to larger reflection values, while Eq. (9) indicates that a higher value of the loss tangent  $\tan \delta_d$  also leads to higher reflection.

Similarly, it can be shown that the reflection coefficient for an absorber made of magnetic material with  $\varepsilon_r = 1$ , and  $\tan \delta_d = 0$  is:

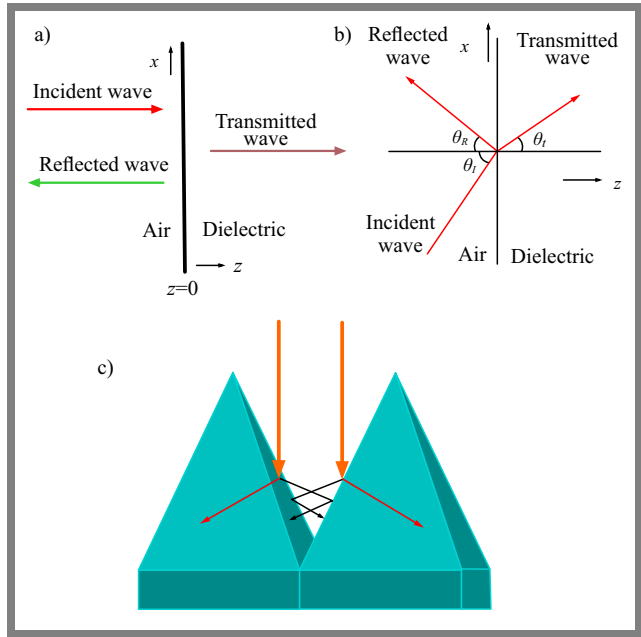
$$\Gamma = \frac{\sqrt{\mu_r' - j\mu_r''} - 1}{\sqrt{\mu_r' - j\mu_r''} + 1}, \quad (11)$$

For the limiting case where the magnetic material of the absorber is lossless,  $\mu_r'' = 0$ , or  $\tan \delta_m = 0$ , Eq. (11) will be:

$$\Gamma = \frac{\sqrt{\mu_r'} - 1}{\sqrt{\mu_r'} + 1}. \quad (12)$$

Equations (11), (12) indicate that higher relative permeability  $\mu_r'$  leads to larger reflection values. Higher value of the loss tangent  $\tan \delta_m$  also leads to increased reflection.

The above relations have been derived for a plane air-dielectric interface. Thus, they cannot be applied directly to the pyramidal shape of the absorber, but offer a clear insight into the reflection generated by the absorber. Moreover, Eq. (7) shows that for a loss-free material, when the relative dielectric



**Fig. 1.** Basic scenario for the reflection of an EM wave from a dielectric surface: a) normal incidence, b) oblique incidence, and c) incidence on the surface of a pyramidal absorber.

constant approaches the relative permeability, the reflection coefficient trends towards zero.

When the EM wave is incident obliquely on the dielectric material, then the reflection coefficient will be influenced by the angle of incidence and polarization with respect to the interface. The reflection coefficient for parallel  $\Gamma_{\parallel}$  and perpendicular  $\Gamma_{\perp}$  polarization is given, respectively, by [17]:

$$\Gamma_{\parallel} = \frac{\eta \cos \theta_t - \eta_0 \cos \theta_i}{\eta \cos \theta_t + \eta_0 \cos \theta_i}, \quad (13)$$

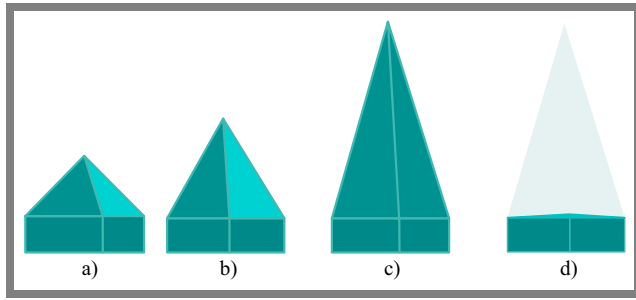
$$\Gamma_{\perp} = \frac{\eta \cos \theta_i - \eta_0 \cos \theta_t}{\eta \cos \theta_i + \eta_0 \cos \theta_t}, \quad (14)$$

The angle of incidence  $\theta_i$  is another influencing factor. When a pyramidal absorber is considered, oblique incidence is inevitable, even if the wave is normally incident on the absorber, as seen in Fig. 1c. Moreover, the wave may reflect from one pyramid's surface towards an adjacent pyramid. Such a case is difficult to analyze using the ray approach, and therefore the use of CST simulation can be a powerful technique to evaluate the reflection's properties.

### 3. Analysis of a Pyramidal Microwave Absorber

Figure 2 shows the construction of a pyramidal microwave absorber that was modeled using the CST simulator. The pyramidal shape was chosen by many designers and manufacturing companies since it offers a gradual change from the air to the absorbing material. Moreover, the sides of neighboring pyramids generate multiple reflections, with a fraction of the incident wave's power being absorbed in each reflection.

Computer simulations assumed a typical pyramidal absorber with a base of  $10 \times 10$  cm with a thickness of 5 cm, and



**Fig. 2.** Investigated pyramidal microwave absorber with heights of: a) 8 cm b) 16 cm, and c) 24 cm. The base of the pyramid is shown in d).

pyramid height values of 8 cm, 16 cm, and 24 cm. Various relative permittivity and loss tangent values were selected as parameters in the simulation.

### 4. Results of the Simulations

The pyramidal absorber was modeled using the unit cell approach with two ports, as shown in Fig. 3. The EM wave was incident on port 1, and the reflection coefficient  $S_{11}$  was determined at port 1, while the transmission coefficient  $S_{21}$  was determined at port 2. The unit cell approach was adopted in the modeling of the absorber, and the frequency domain solver was employed in the calculations.

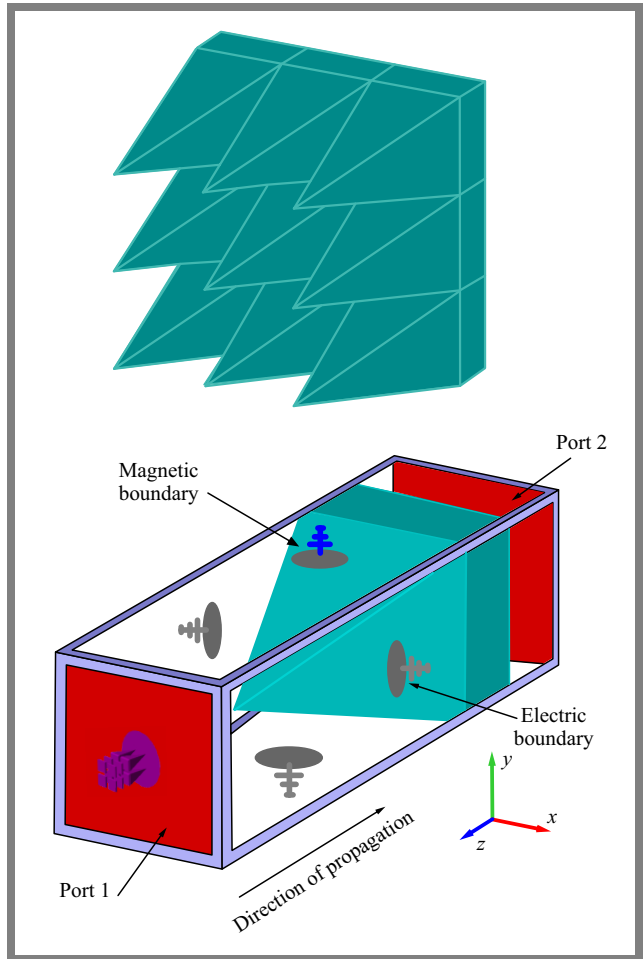
The influence of the various absorber parameters, such as relative permittivity, relative permeability, loss tangent, pyramid height, and angle of the incident wave, were studied by a parameter sweep. In these investigations, one of the parameters was varied, while the remaining variables remained fixed, and the obtained results were used to draw conclusions concerning the impact the variable parameter exerted on the absorber’s performance. The results obtained are presented and discussed in the following subsections.

#### 4.1. Effect of Relative Permittivity

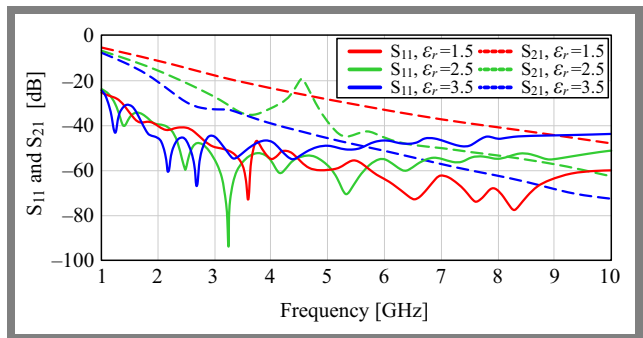
Relative permittivity of the absorber material is an important factor in determining the reflection coefficient. The reflection coefficient for an EM wave propagating from the air onto a plane dielectric material was discussed in Section 2. While that section offered a basic idea about the reflection mechanism, a closed-form relation is difficult to derive for the pyramidal case. Thus, the investigation using CST modelling is capable of offering a faster and easier insight into the investigation.

Figure 4 shows the variation of the reflection coefficient of a non-magnetic absorber with frequency for three values of relative permittivity (1.5, 2.5, and 3.5), with the loss tangent fixed at 0.5, and the height of the pyramid set at 16 cm. The results show that for frequencies above approximately 5 GHz, the reflection increases along with relative permittivity, as indicated by Eq. (10), for the plane dielectric material. However, this trend is not evident at lower frequencies.

This finding may be attributed to the fact that for frequencies above 5 GHz, where the wavelength in the air is smaller



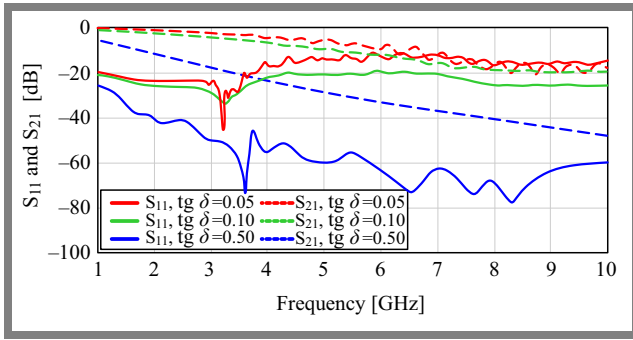
**Fig. 3.** Simulation model used in the CST Microwave software.



**Fig. 4.** Variation of the reflection coefficient  $|S_{11}|$  and transmission coefficient  $|S_{21}|$  with frequency for pyramid height of 16 cm and tan  $\delta_d$ , with various values of  $\epsilon_r$ .

than 6 cm, behavior of the absorber’s surface is closer to that of a plane dielectric layer. However, for lower frequencies, the absorber size is comparable to the wavelength, and the approximation of the absorber by a plane surface departs from the actual case. The variation of the reflection coefficient with frequency shows a faster decline at lower frequencies, compared to a scenario involving higher frequencies.

The same figure illustrates the variation of the transmission coefficient  $|S_{21}|$  with changes in frequency, indicating decreasing values at higher frequencies. Moreover, as the reflection increases with increasing relative permittivity, less power



**Fig. 5.** Variations of reflection coefficient  $|S_{11}|$  and transmission coefficient  $|S_{21}|$  along with frequency changes, for pyramid height of 16 cm and  $\epsilon_r = 1.5$ , for various values of  $\tan \delta_d$ .

is introduced into the absorber, leading to lower transmitted power.

#### 4.2. Effect of Loss Tangent

The loss tangent  $\tan \delta_d$  of the absorber material is an important factor in absorbing the EM wave that penetrates into the absorber material, as it influences the level of the reflected wave that leaves the absorber. Figure 5 shows the level of the reflection coefficient for a non-magnetic material having  $\tan \delta_d$  values of 0.05, 0.1, and 0.5, while the pyramid height was kept at 16 cm, and  $\epsilon_r = 1.5$ . One may clearly see those higher values of loss tangent  $\tan \delta_d$  lead to lower reflection levels. The effect is more pronounced at higher frequencies.

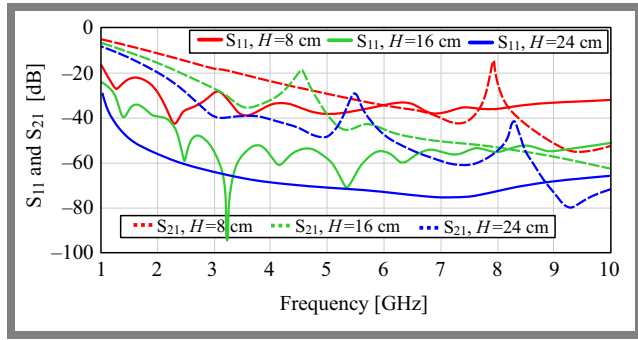
An increase in the loss tangent indicates an increase in the material's ability to convert electromagnetic waves into heat, hence reducing the reflected energy. This explains why the reflection coefficient decreases as the loss tangent increases. On the other hand, as the loss tangent increased, the transmission coefficient  $S_{21}$  decreased, as the wave penetrating the absorber experienced more losses.

#### 4.3. Effect of Pyramid Height

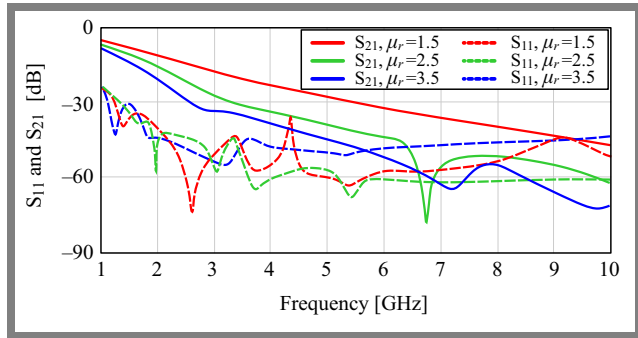
The height of the pyramid is a determinant factor in specifying the thickness of the absorber, and thus the total volume of the absorbing material and the overall shape of the absorber. For a fixed base size, a larger height of the pyramid means a narrower tip and consequently smoother introduction of the absorber material into the air, which will lead to reduced reflection. However, a larger height means a thicker absorber, thus resulting in more weight and cost.

Figure 6 shows the effect of pyramid height  $H$  on the achieved reflection of the absorber, where  $H$  was varied from 8 to 16 cm, and 24 cm, with the parameters of the non-magnetic material remaining fixed at  $\epsilon_r = 2.5$ ,  $\tan \delta_d = 0.5$ .

The results obtained show that a larger pyramid height leads to a lower reflection. An average improvement of approximately 20 dB in the reflection coefficient is noticed when the height increases from 8 to 16 cm, while an improvement between 15 to 20 dB can be noticed when the height is further increased to 24 cm. This finding can be explained by the fact that a larger height, with fixed base dimensions, means larger



**Fig. 6.** Variation of reflection coefficient  $|S_{11}|$  and transmission coefficient  $|S_{21}|$  along with frequency changes, for pyramid heights of 8, 16, and 24 cm, for  $\epsilon_r = 2.5$  and  $\tan \delta_d = 0.5$ .



**Fig. 7.** Variation of reflection coefficient  $|S_{11}|$  and transmission coefficient  $|S_{21}|$  along with frequency changes, for pyramid height 16 cm and  $\tan \delta_m = 0.1$ , with various values of  $\mu_r$ ,  $\epsilon_r = 1$  and  $\tan \delta_d = 0$ .

absorbing volume and thus more losses and, consequently, a lower reflection.

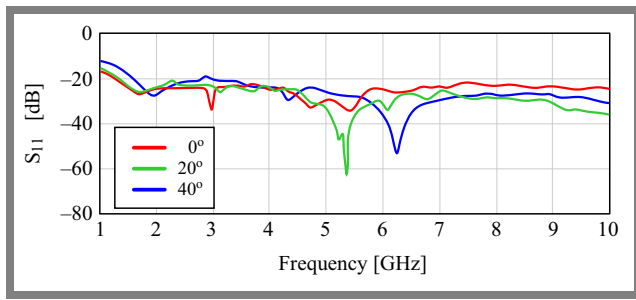
Moreover, a larger height at a fixed base means a smaller tip angle of the pyramid. This results in a more gradual introduction of the absorber material and, consequently, a lower reflection.

Figure 6 also shows the variation of the transmission coefficient  $|S_{21}|$  with frequency for various heights of the pyramid. It is noted that larger heights lead to lower transmission, since the loss caused by the absorber is proportional to its thickness.

#### 4.4. Effect of Permeability

The permeability of the absorber material also affects its performance as it influences both the reflection and absorption inside the material due to the imaginary part of  $\mu$ . Figure 7 illustrates the performance of the investigated absorber when only the magnetic properties are considered, while the dielectric properties are excluded ( $\epsilon_r = 1$ ,  $\tan \delta_d = 0$ ).

One may notice that both reflection and transmission coefficients decrease along with the increase in frequency. Furthermore, each of the coefficients increases as the relative permeability  $\mu_r$  is increased. The general trend of the variation is similar to that illustrated in Fig. 4, when relative permittivity was varied for a non-magnetic material. The similarity in the behavior can be understood by comparing Eq. (10) with Eq. (12), indicating similar relations between re-



**Fig. 8.** Variation of reflection coefficient  $|S_{11}|$  with frequency changes, for pyramid height of 8 cm,  $\epsilon_r = 1.5$  and  $\tan \delta_d = 0.1$ , at various incidence angles of  $0^\circ$ ,  $20^\circ$ ,  $40^\circ$ .

flection coefficient  $\Gamma$  and the real parts of permittivity and permeability, respectively.

Further simulations performed to evaluate the effect of the loss tangent and pyramid’s height generated similar results to those shown in Figs. 5, 6, respectively and the results are not presented here for brevity.

#### 4.5. Effect of Angle of Incidence

In the cases investigated in the previous sections, the EM wave was normally incident on the absorber. However, oblique incidence at a specific angle is encountered in the majority of scenarios. The impact of the angle of incidence stems from the fact that the reflection coefficient depends on the incidence angle with respect to the absorber’s surface and the properties of the absorber’s material.

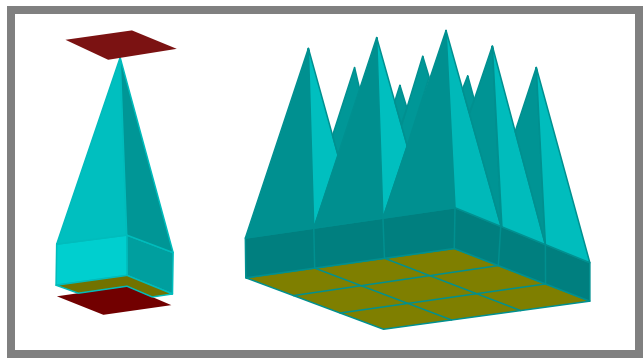
The relations of Eqs. (13) and (14) show the effect that the angle of incidence and the absorber’s parameters have on the reflection coefficient when the EM wave is incident on a half-plane dielectric material. However, closed-form relations for pyramidal absorbers will be very difficult to derive. The effect of the angle of incidence was investigated using proper settings introduced to the CST software.

Figure 8 shows the effect of changing the angle of incidence on a non-magnetic pyramid absorber while fixing the other parameters height of 8 cm,  $\epsilon_r = 1.5$ , and  $\tan \delta_d = 0.1$ . The variation in the reflection coefficient is approximately 10 dB for angles ranging from  $0^\circ$  to  $40^\circ$ . Moreover, lower variations may be observed at lower frequencies. For the case of inclined incidence, the EM wave will face the pyramid’s surface and may penetrate the absorber, exiting from the other side and thus contributing to the reflected wave.

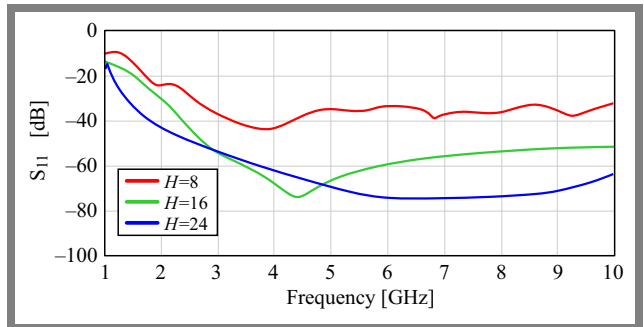
### 5. Reflection from Conductor-backed Absorbers

Many echoic chambers are also designed to shield EM waves. In these designs, the walls of the chamber are usually covered by conducting sheets and then the absorbers are placed on the conducting walls. These designs can be modeled by placing a conducting plane underneath the absorber.

The CST Microwave Suite was used in the modeling process, with a 1 mm thick copper sheet placed under the pyramidal absorber, as shown in Fig. 9. The results obtained for such



**Fig. 9.** Model simulated using the CST Microwave Suite with a copper sheet placed below the pyramidal absorber.



**Fig. 10.** Variation of reflection coefficient  $|S_{11}|$  with frequency changes for a pyramidal absorber backed by a 1 mm thick copper plate, where  $\epsilon_r = 2.5$ ,  $\tan \delta_d = 0.5$ , for  $H = 8, 16$ , and  $24$  cm.

a setup are shown in Fig. 10 for a normal incident wave. The reflection coefficient decreases as the thickness of the absorber increases from 8 to 16 cm and then to 24 cm.

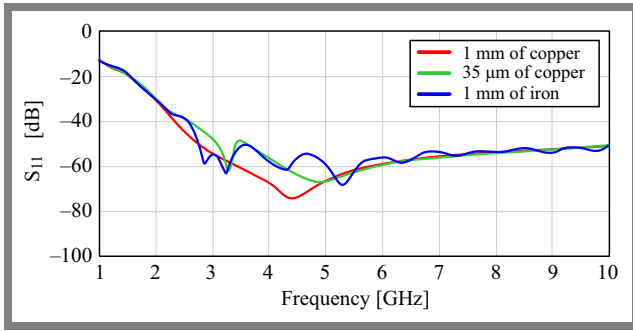
A thicker absorber means that the incident wave propagates through the absorbing material towards the back conductor, where it is fully reflected, and then travels back towards the direction of incidence. This means the wave is attenuated twice inside the absorber. Thus, for adequate losses in the absorber material, the reflected wave is mainly that was initially reflected at the absorber surface.

To show the effect of the back conducting plane, Fig. 11 compares the results obtained for three various backplanes: a  $35 \mu\text{m}$  thick copper plate, a 1 mm copper plate, and a 1 mm galvanized iron plate. The results show similar performance except for the frequency range of 2.5 to 5.5 GHz. The iron plate has shown a reflection of less than  $-50$  dB across 75% of the frequency range and can be considered the best choice in terms of cost and performance.

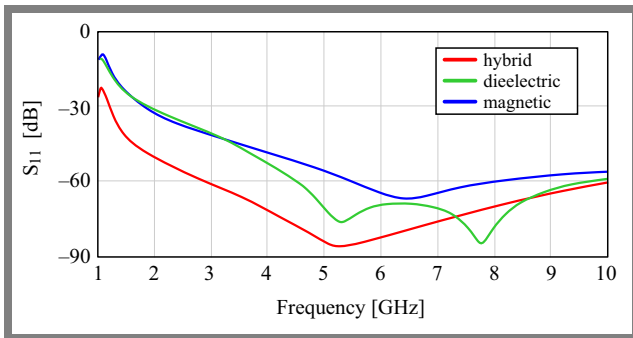
### 6. Electric and Magnetic Properties

Investigations described in the previous sections have shown that lossy dielectric and magnetic materials can be used to design microwave absorbers. The employment of a material having both properties will increase the absorption rate and may help reduce the reflection from the absorber.

As estimated by Eq. (7), the reflection from a plane surface of a certain material can be reduced if the relative values of its permittivity and permeability are equal.



**Fig. 11.** Variation of reflection coefficient  $|S_{11}|$  along with frequency changes for a pyramidal absorber with the height of 16 cm,  $\epsilon_r = 2.5$  and  $\tan \delta_d = 0.5$ , with the following back plate options: 35  $\mu\text{m}$  copper, 1 mm copper, and 1 mm iron.



**Fig. 12.** Changes in reflection coefficient  $|S_{11}|$  along with frequency changes for a 16 cm tall pyramidal absorber and a 1 mm iron backing plate for three scenarios: a)  $\epsilon_r = \mu_r = 1.5$  and  $\tan \delta_d = \tan \delta_m = 0.5$ , b)  $\epsilon_r = 1.5$ ,  $\tan \delta_d = 0.5$ ,  $\mu_r = 1$ , and  $\tan \delta_m = 0$ , c)  $\mu_r = 1.5$ ,  $\tan \delta_m = 0.5$ ,  $\epsilon_r = 1$ , and  $\tan \delta_d = 0$ .

The performance of a 16 cm high pyramidal absorber that is backed by a 1 mm thick iron plate is shown in Fig. 12 for three different cases. The first scenario is when the material has the same magnetic and electric properties  $\epsilon_r = \mu_r = 1.5$  and  $\tan \delta_d = \tan \delta_m = 0.5$ . The two other cases are for a non-magnetic dielectric material and a magnetic material with  $\epsilon_r = 1$  and  $\tan \delta_d = 0$ .

The results demonstrate that a material with dielectric- or magnetic-only properties offers similar performance and can be used while designing the absorber. However, employing a material that offers, simultaneously, both electrical and magnetic loss properties reduce the reflection coefficient appreciably. For the shown cases, the reduction in the reflection coefficient is greater than 15 dB for the majority of the frequency range involved.

## 7. Comparison with Other Published Papers

In order to further assess the results obtained with the help of the simulations, they are compared here with those obtained by the authors of other works. The general feature of the results obtained is the decrease of the reflection coefficient as the frequency of the incident wave changes. This trend was

also observed in [2], [4], [6], and [9], especially at frequencies below 5 GHz.

Table 1 shows various performance parameters identified in previously published papers and compares them with the outcomes of this work. It can be seen that most of the previous research ([5], [6], [10], and [11]) ignored the effect of the loss tangent, despite its significant impact on the reflection coefficient, as was demonstrated in Fig. 6. However, the remaining references made in Tab. 1 took into consideration small ranges of the loss tangent, thus its effect was not demonstrated clearly.

Most of the former literature has considered agricultural residuals and other cheap materials with the objective of reducing the cost of absorbers and find use for the leftover materials. However, it failed to consider the impact of moisture, as excessive moisture can cause degradation of the materials used in absorbers, affecting the permittivity value and, consequently, the absorption value. In other words, the performance of absorbers based on these materials becomes unpredictable.

As the absorbers are used mainly to cover the walls of anechoic chambers in order to prevent reflections generated in certain regions of the test environment, their thickness is a crucial parameter. This parameter determines the volume of the material used and, consequently, the cost of the solution. In [2]–[4], [8] and [9], pyramids having a total thickness of 15 cm were used (along with their base) and achieved average reflection coefficients between  $-21$  and  $-44$  dB. This meant a reflection level/cm of thickness ratio of  $-1.4$  to  $-3$  dB/cm. The experiments described in [6], [11] employed a material thickness of 25 cm with a base of 5 cm, or an overall thickness of 30 cm, but achieved reflection levels of  $-30$  dB and  $-25$  dB, respectively, thus obtaining a ratio of  $-1$  to  $-0.83$  dB/cm. However, in the three examples investigated in this paper, materials with a thickness of 8 cm, 16 cm, and 24 cm were used, with a 5 cm thick base, and achieved an average reflection coefficient of  $-52.6$  dB,  $-56.7$  to  $-66.5$  dB, or  $-2.5$  dB/cm,  $-2.7$  dB/cm, and  $-2.3$  dB/cm, respectively.

The fourth design example, which used magnetic and dielectric properties of the absorber material, achieved an average reflection of  $-67$  dB, which is equivalent to  $-3.2$  dB/cm. This shows that performance of the absorber is much better when the material is characterized by equal values of relative permittivity and permeability.

## 8. Conclusions

The choice of the appropriate absorber is influenced by the required level of reflection, the thickness of the absorber, and its weight and cost. Relative permittivity and permeability impact the initial reflection at the surface of the absorber, and it was demonstrated analytically and through computer simulations that lower permittivity or permeability result in lower reflection coefficients. Moreover, the reflection coefficient can be appreciably reduced if the absorber material has equal values of relative permittivity and permeability.

**Tab. 1.** Performance comparison of the results obtained with those of previous papers.

Ref.	Material	Dimensions [cm]	Frequency range [GHz]	$\epsilon_r$	$\tan \delta$	$S_{11}$ [dB]	dB/cm
[2]	Rubber tire dust and rice husk (75:25)	Base $5 \times 5 \times 2$ , pyramid's height 13	7 – 13	3.43	0.048	Average –32.51	–2.27
[2]	Rubber tire dust and rice husk (50:50)		7 – 13	2.67	0.076	Average –22.03	–1.47
[2]	Rubber tire dust and rice husk (25:75)		7 – 13	2.08	0.103	Average –21.54	–1.44
[3]	SCB		0.1 – 20	1.44	0.161	Average –44.39	–2.96
[4]	Rice husk		0 – 20	1.91	0.079	Average –31.93	–2.12
[4]	A mixture of rice husk and coal		0 – 20	2.50	0.086	Average –43.5	–2.9
[5]	Water hyacinth	Square pyramid shape $5 \times 5 \times 13$	0 – 20	NA	NA	Lowest about –30 dB, highest about –10 dB	–2.3 –0.77
[6]	Polyurethane	Base $10 \times 10 \times 5$ , pyramid height 25	1 – 10	3.5	NA	Better absorption at higher frequencies: below –30 dB from 3 GHz	–1
[6]	Carbon		1 – 10	2.6	NA		–1
[6]	Polyurethane with carbon		1 – 10	3.5 and 2.6	NA		–1
[8]	Rice husk	Triangle base of height 2, pyramid height 13, side lengths $5.6 \times 5$	1 – 20	2.9	0.084	Average –41.142	–2.7
[8]	Rice husk	Equilateral triangle, base height 2, pyramid height 13, side length 5.25 (3 sides)	1 – 20	2.9	0.084	Average –39.878	–2.6
[8]	Rice husk	Square base, height 2, pyramid height 13, side length 5 (4 sides)	1 – 20	2.9	0.084	Average –39.423	–2.6
[9]	Banana leaves	Hexagonal pyramid, base thickness 2, pyramid height of 13	0.01 – 20		1.014	Average –35.3	–2.4
[9]	Rice husk		0.01 – 20	3.22	0.827	Average –39.70	–2.6
[9]	Rice saw		0.01 – 20	3.10	0.956	Average –38.92	–2.6
[9]	Banana leaves	Truncated hexagonal pyramid, base thickness 2, pyramid height 13, with truncated miniature hexagon	0.01 – 20	2.49	1.014	Average –35.37	–2.4
[9]	Rice husk		0.01 – 20	3.22	0.827	Average –35.94	–2.4
[9]	Rice straw		0.01 – 20	3.1	0.956	Average –37.2	–2.5
[10]	Kenaf	$20 \times 20 \times 2$	1 – 12	2	NA	NA	NA
[10]	Coconut coir			2.6			
[11]	Carbon	Base $10 \times 10 \times 5$ , pyramid height 25, with a tip of $1 \times 1$	0.01 – 10	2.6	NA	Mostly below –25	–0.83
This work	Carbon	Base $10 \times 10 \times 5$ , pyramid height 16	1 – 10	2.5	0.5	Average –52.6	–2.5
		Base $10 \times 10 \times 5$ , pyramid height 16	1 – 10	1.5	0.5	Average –56.66	–2.7
		Base $10 \times 10 \times 5$ , pyramid height 24	1 – 10	2.5	0.5	Average –66.5	–2.3
	Carbon with iron back plate	Iron back plate $10 \times 10 \times 1$ mm base $10 \times 10 \times 5$ , pyramid height 16	1 – 10	$\mu_r=1, \epsilon_r=1.5$	$\delta_m=0, \delta_d=0.5$	Average –56.9	–2.7
			1 – 10	$\mu_r=1.5, \epsilon_r=1$	$\delta_m=0.5, \delta_d=0$	Average –51	–2.4
			1 – 10	$\mu_r=1.5, \epsilon_r=1.5$	$\delta_m=0.5, \delta_d=0.5$	Average –67	–3.2

The loss tangent of the absorber material is an important factor as well. It was demonstrated in the course of the simulations that larger values of the loss tangent reduce reflection. The shape of the absorber also influences the reflection, as pyramidal absorbers offer a gradual introduction of the absorbing material into the air, and they provide a greater degree of design freedom.

The results of this investigation may help designers choose the proper parameters of the absorber material to achieve better performance. If the proper material is not naturally available, combinations of various materials may be used to achieve the desired parameters for the absorber.

## References

- [1] A.A. Abu Sanad *et al.*, "The Prospect of Using Hollow Pyramidal Microwave Absorbers for 5G Anechoic Chamber Applications: A Review", *Journal of Applied Physics*, vol. 136, art. no. 230701, 2024 (<https://doi.org/10.1063/5.0244666>).
- [2] M.F. Bin Abd Malek *et al.*, "Rubber Tire Dust-rice Husk Pyramidal Microwave Absorber", *Progress In Electromagnetics Research*, vol. 117, pp. 449–477, 2011 (<https://doi.org/10.2528/PIER11040801>).
- [3] L. Zahid *et al.*, "Development of Pyramidal Microwave Absorber Using Sugar Cane Bagasse (SCB)", *Progress In Electromagnetics Research*, vol. 137, pp. 687–702, 2013 (<https://doi.org/10.2528/pier13012602>).
- [4] H. Kaur, G. Deep, and V. Chawla, "Enhanced Reflection Loss Performance of Square Based Pyramidal Microwave Absorber Using Rice Husk-coal", *Progress In Electromagnetics Research M*, vol. 43, pp. 165–173, 2015 (<https://doi.org/10.2528/PIERM15072603>).
- [5] A. Nuan-on *et al.*, "Design and Fabrication of Microwave Absorbers Using Water Hyacinth", *Engineering Access*, vol. 3, pp. 7–10, 2017 (<https://doi.org/10.14456/mijet.2017.2>).
- [6] S.I. Orakwue and I.P. Onu, "Pyramidal Microwave Absorber Design for Anechoic Chamber in the Microwave Frequency Range of 1 GHz to 10 GHz", *European Journal of Engineering and Technology Research*, vol. 4, pp. 1–3, 2019 (<https://doi.org/10.24018/ejers.2019.4.10.1409>).
- [7] H. Nornikman, P.J. Soh, A.A.H. Azremi, and M.S. Anuar, "Performance Simulation of Pyramidal and Wedge Microwave Absorbers", *2009 3rd Asia International Conference on Modelling and Simulation*, Bandung, Indonesia, 2009 (<https://doi.org/10.1109/AMS.2009.13>).
- [8] H. Nornikman *et al.*, "Reflection Loss Performance of Triangular Microwave Absorber", *International Symposium on Antennas and Propagation*, 2010 [Online] Available: [https://www.ieice.org/cs/isap/ISAP\\_Archives/2010/pdf/281.pdf](https://www.ieice.org/cs/isap/ISAP_Archives/2010/pdf/281.pdf).
- [9] H. Nornikman, F. Malek, P.J. Soh, and A.A.H. Azremi, "Reflection Loss Performance of Hexagonal Base Pyramid Microwave Absorber Using Different Agricultural Waste Material", *2010 Loughborough Antennas and Propagation Conference, LAPC 2010*, Loughborough, UK, 2010 (<https://doi.org/10.1109/LAPC.2010.5666029>).
- [10] L.M. Kasim *et al.*, "A Study of Electromagnetic Absorption Performance of Modern Biomass Wall Tile", *International Journal of Electrical and Electronic Engineering & Telecommunications*, vol. 9, pp. 429–433, 2020 (<https://doi.org/10.18178/IJEETC.9.6.429-433>).
- [11] H. Nornikman, P.J. Soh, and A.A.H. Azremi, "Potential Types of Biomaterial Absorber for Microwave Signal Absorption", *4th International Conference on X Rays and Related Techniques in Research and Industries 2008 (ICXRI 2008)*, Kota Kinabalu, Sabah, Malaysia, 2008.
- [12] Y.M. Zong, "Optimization of Multilayer Microwave Absorbers Using Multi-strategy Improved Gold Rush Optimizer", *ACES Journal*, vol. 39, pp. 708–717, 2024 (<https://doi.org/10.13052/2024.ACES.J.390806>).
- [13] R. Xu *et al.*, "An Ultra-wideband Metamaterial Absorber with Angular Stability", *ACES Journal*, vol. 39, pp. 675–682, 2024 (<https://doi.org/10.13052/2024.ACES.J.390802>).
- [14] M.B. Jasim and K. Sayidmarie, "Radar Cross-section Reduction of Planar Absorbers Using Resistive FSS Unit Cells", *Journal of Telecommunication and Information Technology*, no. 4, pp. 61–67, 2023 (<https://doi.org/10.26636/jtit.2023.4.1331>).
- [15] M.F. Asmadi *et al.*, "The Optimal Performance of a Geopolymer Hollow Pyramidal Microwave Absorber with Triangular Slotted", *Solid State Phenomena*, vol. 344, pp. 97–102, 2023 (<https://doi.org/10.4028/p-belmea>).
- [16] A.S. Yusof *et al.*, "Slotted Triangle on Hollow Pyramidal Microwave Absorber Characteristics", *2016 6th IEEE International Conference on Control System, Computing and Engineering (ICCSCE)*, Penang, Malaysia, 2016 (<https://doi.org/10.1109/ICCSCE.2016.7893639>).
- [17] D.M. Pozar, *Microwave Engineering*, John Wiley & Sons, 3rd ed., 720 p., 2005 (ISBN: 9780471448785).

### Aya Raad Thanoon, PG Student

Department of Communication Engineering

 <https://orcid.org/0009-0006-5000-4765>

E-mail: [ayah.dhunun2012@stu.uoninevah.edu.iq](mailto:ayah.dhunun2012@stu.uoninevah.edu.iq)

Ninevah University, Mosul, Iraq

<https://uoninevah.edu.iq/en/>

### Khalil H. Sayidmarie, Professor

Department of Communication Engineering

 <https://orcid.org/0000-0001-6525-0949>

E-mail: [kh.sayidmarie@uoninevah.edu.iq](mailto:kh.sayidmarie@uoninevah.edu.iq)

Ninevah University, Mosul, Iraq

<https://uoninevah.edu.iq/en/>

# UAV-BS-based Hybrid OMA-NOMA System with Multiple Antennas for Multi-user Communication

Ameer Y. Sadeeq and Mohamad A. Ahmed

*Ninevah University, Mosul, Iraq*

<https://doi.org/10.26636/jtit.2025.2.2045>

**Abstract** — In this paper, an unmanned aerial vehicle (UAV) using hybrid orthogonal multiple access (OMA) and non-orthogonal multiple access (NOMA) solutions aided by multiple input multiple output (MIMO) technology is proposed to provide wireless communication for ground users (GUs). The proposed OMA-NOMA-MIMO system aims to improve throughput and spectrum efficiency. Additionally, it also strives to maximize the sum rate while achieving good user fairness. UAVs are considered as base stations (BSs) that provide services to users in multiple real-life scenarios, e.g. during natural disasters. They also enable aerial surveillance and help establish BSs during mass events. A pairing algorithm is proposed for far-near NOMA users with an optimized power allocation (PA) mechanism to improve the performance of NOMA-UAV-BS. Channels between UAV-BS and GU are established as being of the line-of-sight (LoS) and non-line-of-sight (NLoS) varieties, taking into consideration the angle of departure (AoD) and the angle of arrival (AoA). The results obtained demonstrate that NOMA performs best in specific scenarios, while OMA overcomes NOMA in others. The outcomes of the project may be utilized to control the transmission performed by the UAV-BS by serving the GUs depending on the required quality of service.

**Keywords** — *non-orthogonal multiple access, pairing algorithms, successive interference cancellation, unmanned aerial vehicle*

## 1. Introduction

Unmanned aerial vehicles (UAVs) may be useful in providing access to wireless networks by acting as base stations (BS) to expand coverage at a location where a fixed radio BS is not available [1]– [3]. Commonly, in wireless cellular communication, such as 2G, 3G, and 4G systems that rely on orthogonal multiple access (OMA) techniques in the air interface, each user is assigned a single resource block, enjoys limited throughput, and suffers from insufficient user fairness [4], [5]. The problem becomes more acute when the number of users increases in congested areas.

The non-orthogonal multiple access (NOMA) technique is designed for 5G and beyond radio access cellular networks. This multiple access solution is capable of simultaneously utilizing a single resource block, i.e. the same frequency band, for several users. This mechanism leads to a significant improvement in throughput and spectrum efficiency, with satisfactory user fairness levels. Therefore, in this article, the

focus will be on using NOMA for UAV-BS to overcome the challenges faced by OMA radio access technologies [3], [6]. Serving multiple users in NOMA can be achieved by assigning a different power level to each user. On the contrary, serving users in OMA is implemented by assigning a fixed power level for each user. This primary difference leads to OMA outperforming NOMA in terms of spectral efficiency and throughput [7], [8]. To further enhance NOMA's performance, multiple input multiple output (MIMO) technology may be merged with NOMA, as suggested in [8], [9].

In this paper, a pairing algorithm is proposed for the NOMA technology which divides the multiple users into groups, i.e., clusters, with each group consisting of two users. In other words, the mechanism pairs a user with a good channel gain (stronger user) with another user with poor channel gain (weak user) [10]. The channel proposed between UAV-BS and the ground user (GU) is, in some cases, of the line of sight (LoS) and in other cases of the non-LoS (NLoS) variety, taking into consideration the probability of LoS and NLoS [11]. We assume that the superimposed NOMA signal from UAV-BS reaches GU over a Rician fading channel.

Some key parameters play an important role in determining the quality of the channel. These include the following: Rician factor, path loss exponent, distance between UAV-BS and GU, angle of departure (AoD), angle of arrival (AoA) and carrier wavelength [12].

The authors contribution to previous works is summarized as follows:

- In this work, we investigate unmanned aerial vehicles (UAV) relying on hybrid orthogonal multiple access (OMA) and nonorthogonal multiple access (NOMA) mechanisms, aided by the multiple input multiple output (MIMO) technology. Several metrics have been measured at different GUs using two channel types (LoS and NLoS) to estimate the sum rate GU with the overall rate of NOMA and OMA and determine the probability of  $P_{out}$  outage, with a comparison between NOMA and OMA.
- To enhance the overall performance of the system, UAV-BSs are equipped with the MIMO technology.
- The pairing algorithm is employed to significantly improve spectral efficiency and achieve better user fairness. All users who are underserved by the UAV-BS are divided

into multiple groups according to such parameters as user channel gain (weak or strong), user data rate, and distance of the user from the UAV-BS.

- To account for all potential channels that may be established between the UAV-BS and GU, we consider LoS and NLoS channels. To make the channel calculations more accurate, we consider the elevation and azimuth of AoD and AoA along with Rician fading and Rayleigh fading.
- In [13], non-linear energy harvesting in NOMA systems with a fixed base station has been proposed and investigated, while in this work, we propose a hybrid OMA-NOMA scheme intended specifically for UAV-based BSs equipped with MIMO antennas to serve multiple users in no coverage areas.
- In [14], the SIC error on near and far users in the NOMA technology is studied by estimating the near- and far-user BER. The outcomes prove the superiority of NOMA over OMA. Here, we estimated  $P_{out}$  and sum rate for each user in NOMA and OMA and we used SIC without the error effect on near and far users in NOMA.
- In [15], different antenna techniques are developed to improve the performance of wireless mobile communication using the MIMO technology. Our study models UAV-BS MIMO communication with explicit AoD and AoA considerations under Rician fading, thus ensuring more accurate channel estimation for aerial systems.

The rest of the paper is arranged as follows. A model of the system is presented in Section 2, and the pairing of users is described in Section 3. Section 4 outlines the detection process performed by the receiver. Simulation results and discussions are presented in Section 5, with conclusions contained in Section 6.

## 2. System Model

As illustrated in Fig. 1, we propose that several users, i.e.  $U_1, U_2, U_3, \dots, U_N$  present at a particular location take advantage of the services provided by a UAV-BS by applying either the NOMA or OMA technology. It is assumed that the users are located at different distances, i.e.  $d_1, d_2, d_3, \dots, d_N$ , from the UAV. According to these distances,  $U_1$  is the weakest user due to being located at the furthest distance from the base station (UAV) (this user has a poor channel gain).  $U_4$  is the strongest user due to being located the closest (strong channel gain) to the UAV.

According to the principles of NOMA, the power allocation coefficients for each user are denoted by  $a_1, a_2, a_3, \dots, a_N$ . The larger portion of the power is assigned to the weakest users (farthest), while a lower power level is assigned to the strongest users (nearest) to ensure user fairness [3]. The combined power allocation coefficients provided by the UAV-BS should equal 1, in accordance with the NOMA theory [16].

The UAV-BS is equipped with multiple antennas  $N_t$ , while each GU is equipped with antenna  $N_r$ . Such a MIMO structure

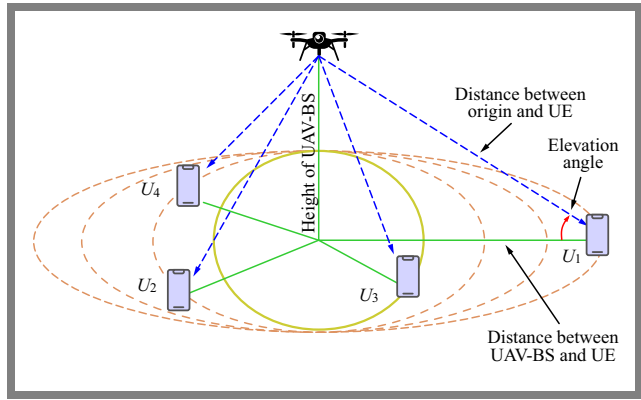


Fig. 1. Model of a UAV-BS communication system.

is intended to enhance the performance of the proposed NOMA system [16].

Three NOMA power allocation schemes may be used for UAV-based wireless communication: fixed power allocation (FPA), equal power allocation (EPA) and fractional transmit power allocation (FTPA) [3], [17].

- In FPA: the distance between the user and the UAV-BS is taken into consideration; low power is allocated to a given user if it is located close to the UAV base station, and vice versa.
- In EPA: constant power is allocated to all users and the distance between the user and the UAV-BS is also taken into consideration (a distinction is made whether the user is close to or far away from the UAV base station).
- In FTPA: the power is allocated dynamically, meaning that the distance between the user and the UAV-BS is taken into consideration. In other words, FTPA allocates low power to users who enjoy good channel conditions and high power to users with poor channel conditions. This approach ensures fairness between users in terms of power allocation [3].

Several scenarios are considered to evaluate the GUs rates, the UAV-BS's sum rate, error probability (bit error rate – BER), and spectral efficiency (SE) for the proposed system by relying on orthogonal multiple access (OMA) and NOMA techniques under different channels conditions.

In the first scenario, we assume that  $U_1$  is the user located the farthest from the UAV-BS. The highest power allocation coefficient is assigned to this user, considering  $a_2, a_3, \dots, a_N$  as interference. The achievable rate at the  $n$ -th user is [3]:

$$R_1 = B \log_2 \left( 1 + \frac{a_1 p_t g_1}{a_2 p_t g_1 + a_3 p_t g_1 + a_4 p_t g_1 + w_1} \right), \quad (1)$$

where  $p_t$  is the transmit power of the UAV-BS, and for any  $n$ -th user, power allocation factor (PAF), noise power and channel gain are denoted by  $a_n, w_n$ , and  $g_n$ , respectively.

At the second user,  $U_2$ , PAF is lower than the one assigned to  $U_1$ , as the user has better channel gain, i.e. ( $a_1 > a_2$ ), and this mechanism is applied to the rest of the users depending on their channel gain, i.e. ( $a_1 > a_2 > a_3 > a_4$ ). Furthermore, all users, except the farthest one, require the application of SIC to subtract the superposed signals of higher order users.

Therefore, the achievable rate at  $U_2$  after removing the signal of  $U_1$  is given as:

$$R_2 = B \log_2 \left( 1 + \frac{a_2 p_t g_2}{a_3 p_t g_2 + a_4 p_t g_2 + w_2} \right). \quad (2)$$

For the third user, the achievable rate at  $U_3$ , after removing interference caused by the signals of  $U_1$  and  $U_2$ , is given as:

$$R_3 = B \log_2 \left( 1 + \frac{a_3 p_t g_3}{a_4 p_t g_3 + w_3} \right). \quad (3)$$

The same procedure is applied to the fourth user, in which case the achievable rate at  $U_4$ , after applying SIC to remove all interfering signals caused by other users with higher PAFs, is expressed as:

$$R_4 = B \log_2 \left( 1 + \frac{a_4 p_t g_4}{w_4} \right). \quad (4)$$

On the other hand, when the OMA technique is used, instead of NOMA, to serve the same four users, the achievable rates for any user can be given as:

$$R_n^{OMA} = 0.5 B \log_2 \left( 1 + \frac{p_t g_n}{w_n} \right). \quad (5)$$

The average rate for a particular user served by NOMA, i.e. the  $n$ -th user, is:

$$\bar{R}_n = E\{R_n\}, \quad (6)$$

where  $E\{\cdot\}$  refers to the expectation process.

Similarly, the average rate achievable by any user served by OMA can be expressed as [18]:

$$\bar{R}_n^{OMA} = E\{R_n^{OMA}\}. \quad (7)$$

The sum rate  $SR$  of the NOMA system, i.e., the overall rate provided by the UAV, can be found by adding the individual rates of all users:

$$SR_{NOMA} = \sum_{n=1}^N R_n. \quad (8)$$

The sum rate of the OMA system can be determined by adding up the individual rates of users:

$$SR_{OMA} = \sum_{n=1}^N R_n^{OMA}. \quad (9)$$

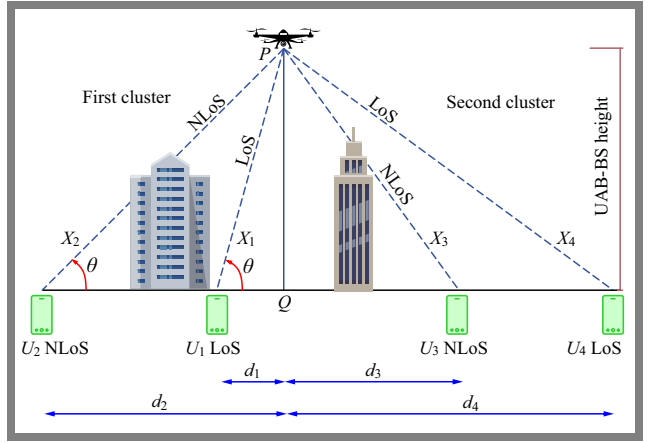
Furthermore, spectral efficiency  $SE$ , in bits/sec/Hz, regardless of the technique used, is:

$$SE = \frac{SR}{BW}. \quad (10)$$

## 2.1. Channel Model

As shown in Fig. 2, the channels between the GUs and UAV-BS are of the LoS and NLoS variety [19]. These channels are impacted by environmental factors prevailing within the coverage area, represented by the density and altitude of the buildings as well as the distance between the ground and the UAV-BS [11].

The total path loss, as shown in Fig. 2, of both GUs in cluster 1, i.e.  $U_1$  and  $U_2$ , is computed either directly, based on the free-space path loss, or based on the excessive losses that occurred



**Fig. 2.** Channel model with LoS and NLoS links, where  $\{X_1, X_2, X_3, X_4\}$  represent the actual physical distances between the GUs and the UAV-BS.

**Tab. 1.** Path loss exponent for different environments [20].

Type of environment	$\eta$
Free space	2
Urban area cellular radio	2.7 to 3.5
Shadowed urban cellular radio	3 to 5
In building line-of-sight	1.6 to 1.8
Obstructed by buildings	4 to 6
Obstructed by factories	2 to 3

along the NLoS paths due to reflections of the transmitted signals from obstacles within the coverage area.

The power received by  $j$ -th user is determined as in [11]:

$$Prx_j(\text{dB}) = Ptx(\text{dB}) - L_j(\text{dB}). \quad (11)$$

where  $Ptx$  indicates the power transmitted by the UAV-BS and  $L_j$  indicates the path loss of the air-to-ground (A2G) channel between the UAV-BS and any GU. The path loss for the A2G channel for the  $j$ -th link is represented by distance  $X_j$  and is evaluated similarly to [11] for LoS and NLoS links, as:

$$L_j = 10 \eta \log_{10}(X_j) + X_{LoS}, \quad (12)$$

$$L_j = 10 \eta \log_{10}(X_j) + X_{NLoS}, \quad (13)$$

respectively, where  $\eta$  represents the path loss exponent, which is considered one of the important parameters in wireless communication. The value of the loss-pass exponent is affected by environmental factors, such as interference, reflection, and diffraction. The value of the path loss exponent for different environments is presented in Tab. 1. In addition,  $X_{LoS}$  and  $X_{NLoS}$  represent the excessive path losses of both LoS and NLoS links, respectively.

The probability that a GU has a LoS link with a UAV-BS is given by [11]:

$$Pr_{LoS}(j) = \frac{1}{1 + \alpha e^{-\beta(\theta_j - \alpha)}}, \quad (14)$$

**Tab. 2.**  $\alpha$  and  $\beta$  values of various environment types [3].

Environment	$\alpha$	$\beta$
Suburban	0.1	750
Urban	0.3	500
Dense urban	0.5	300
Urban high-rise	0.5	300

where  $\alpha$  and  $\beta$  are constant values linked to a given environmental profile, such as urban area, dense urban area, etc. The range of  $\alpha$  and  $\beta$  varies from 0.1 to 0.8 and 100 to 750, respectively. In practical terms,  $\alpha$  is the ratio of the ground area covered by buildings to the total ground area (dimensionless) and  $\beta$  is the number of buildings per square kilometer [3].

In Tab. 2, some examples of  $\alpha$  and  $\beta$  values for various environment types are presented.

The probability of a ground user having an NLoS link with a UAV-BS is given by [11]:

$$Pr_j(NLoS) = 1 - Pr_j(LoS). \quad (15)$$

In general, the UAV is not aware of the type of terrain in its vicinity to specify the type of link (LoS or NLoS). Therefore, Eq. (12) is reworked as [11]:

$$Pr_{x,j}(\text{dB}) = Ptx(\text{dB}) - \bar{L}_j(Rc, H). \quad (16)$$

where  $\bar{L}_j(Rc, H)$  determines the mean path loss including probabilities for a LoS and NLoS link between the UAV-BS and ground user. It is calculated as follows [11]:

$$\bar{L}_j(Rc, H) = Pr_j(LoS)L_j(LoS) + Pr_j(NLoS)L_j(NLoS). \quad (17)$$

As mentioned above, the signal from the UAV-BS reaches the GU directly by LoS and not directly by NLoS over the Rician fading channel. In this paper, we considered both AoD and AoA cases.

The channel between the UAV-BS and the ground user with LoS and NLoS is expressed as follows [12]:

$$H = \sqrt{\xi \frac{K}{K+1}} b(\theta_x, \theta_y)_{AoA} a(\theta_x, \theta_y)_{AoD} + \sqrt{\frac{\xi}{K+1}} \bar{H}. \quad (18)$$

where  $\bar{H}$  is the NLoS component,  $\xi = d^{-\eta}$  represents large-scale fading,  $K$  is the Rician factor, while  $a$  and  $b$  are the steering vectors of the UAV-BS and the ground user, respectively.

The Rician factor  $K$  represents the ratio between the power of the LoS component and the power of the scattered (NLoS) [21], i.e., it is the parameter representing the strength of the LoS component:

$$K = \frac{P_{LoS}}{P_{NLoS}}. \quad (19)$$

This factor measures fading severity, for example  $K = 0$  represents NLoS (Rayleigh channel), and  $K = \infty$  represents LoS (absent fading case). In other words, when  $K \gg 1$ , strong LoS dominance is present. The different values of the Rician factor for various environments are shown in Tab. 3.

**Tab. 3.** Rician factor  $K$  for different environments [18].

Environment	Rician $K$ factor [dB]
Lake	13.10
Hilly	13.61
Rural	9.28
Suburban	6.93

On the UAV-BS side, the elevation and azimuth (AoD) along  $x$  and  $y$  axes are as follows [12]:

$$\theta_{xAoD} = -\frac{2\pi d_{BS}}{\lambda} \cos \theta_{BS} \cos \phi_{BS}, \quad (20)$$

$$\theta_{yAoD} = -\frac{2\pi d_{BS}}{\lambda} \cos \theta_{BS} \sin \phi_{BS}. \quad (21)$$

On the GU side, the elevation and azimuth (AoA) along  $x$  and  $y$  axes are [12]:

$$\theta_{xAoA} = \frac{2\pi d_{GU}}{\lambda} \cos \theta_{GU} \cos \phi_{GU}, \quad (22)$$

$$\theta_{yAoA} = \frac{2\pi d_{GU}}{\lambda} \cos \theta_{GU} \sin \phi_{GU}. \quad (23)$$

Channel gain is:

$$g_n = |H_n * w_n|^2, \quad (24)$$

where  $H_n$  and  $w_n$  are the channel matrix and the beamforming vector for the  $n$ -th user, respectively.

Beamforming techniques can be employed in wireless communications systems to eliminate inter-user interference. One of the most powerful and efficient approaches is the zero-force (ZF) beamforming method. It is considered a linear beamforming precoder that applies weight to users' signals at the UAV-BS in the downlink phase of the transmission. The ZF precoding vector for the  $n$ -th user that passes through channel  $H_n$  can be expressed as:

$$w_n = \rho H_n^H (H_k H_n^H)^{-1}, \quad (25)$$

with

$$\rho = \frac{1}{\sqrt{Tr(H H^H)^{-1}}}, \quad (26)$$

where  $\rho$  is a normalization value that satisfies transmitting signals within the available transmitted power of the UAV-BS. Furthermore,  $(A)^{-1}$ ,  $Tr(A)$ , and  $A^H$  are used to denote the operations of matrix inversion, a trace of a matrix, and the Hermitian transport of matrix  $A$ , respectively.

### 3. Pairing Users

The pairing algorithm is used to improve spectral efficiency and achieve better user fairness for systems utilizing the NOMA technique with many users. By using the pairing technique, the system obtains the required information about the users' circumstances, i.e. the channel gain, and utilizes this information for dividing the coverage area into clusters, with each of them comprising a specific number of users. The pairing algorithm plays a significant role in selecting users with different channel gains to be served by the UAV-BS. In

each cluster, a high power level is allocated to the user who has a weak channel gain, i.e., the far user (FU), while a lower power level is allocated to a user who has a strong channel gain, i.e., the near user (NU) [13]. Figure 3 shows the paired users within clusters and the SIC operation associated with NOMA, relied upon to detect the signal for each user from the superposed NOMA signal.

In this paper, we consider all users within the coverage area of the UAV-BS, divided into groups based on a pairing algorithm. Each group  $M$  is to serve  $2M$  users, meaning each cluster contains two paired users. The pairing algorithms perform tests to arrange channel gains  $g_n$  of all users in a descending order, i.e.,  $g_1 \geq g_2 \geq \dots \geq g_{2M}$ . The UAV-BS establishes the first group by pairing the nearest user characterized by channel gain  $g_1$  with the farthest user having channel gain  $g_{2M}$ . Then, the second user in the list  $g_2$  is paired with users  $g_{2M-1}$ . This operation is continued in the same manner until all users are paired within their clusters.

It is worth noting that this way of pairing is called near-far pairing (NF) [10]. Algorithm 1 shows the NF pairing algorithm for NOMA technology users [21], [22].

---

#### Algorithm 1 Near-far (NF) pairing algorithm.

---

**Input:** Channel gain for users  $g_n = [g_1 + g_2 + \dots + g_{2M}]$ , number of groups, number of users  $2M$

- 1: Arrange channel gain of all users in downward order:  
 $g_1 \geq g_2 \geq \dots \geq g_{2M}$
- 2: Determine the group of channel gains:  
 $Q = \{g_1, g_2, \dots, g_{2M}\}$
- 3: **for**  $P = 1$  to  $M$  **do**
- 4:      $g_n = \{ \}$
- 5:      $g_{max} = \max\{Q\}, g_{min} = \min\{Q\}$
- 6:      $g_n = g_{max} \cup g_{min}$
- 7:      $Q = Q(p + 1 : \text{end} - 1)$
- 8: **end for**

**Output:** Group of pairs

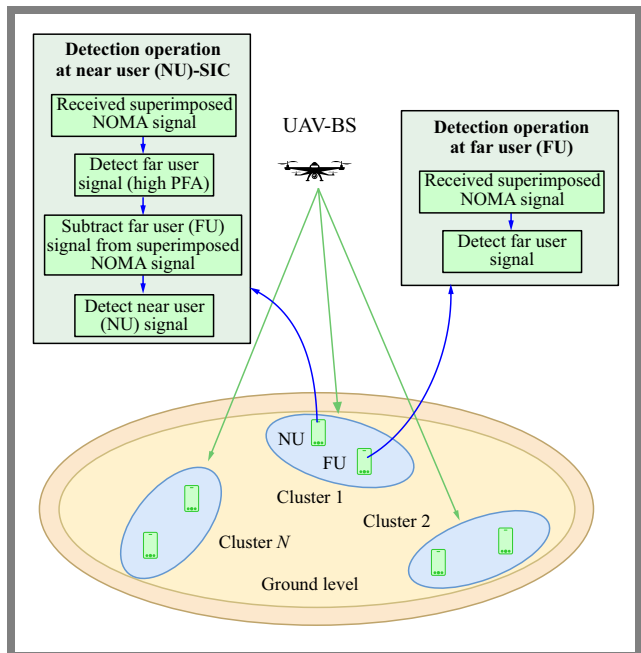
---

It is noteworthy that their other pairing strategies, such as near-near (NN) pairing and far-far (FF) pairing, which can be used for this purpose, are available as well.

## 4. Detection Operation

In the receiver, the superimposed NOMA signal from the UAV-BS is received by multiple users. In accordance with the NOMA principle, two users (near and far) are grouped within one cluster with the use of a pairing technique. As mentioned above, a lower PAF is provided to the NU, i.e. the user with strong channel gain, and a high PAF is given to the FU, the one with weak channel gain.

Furthermore, the NU needs the SIC to be implemented to eliminate interference caused by the power of the far user, which is considered significantly high. The SIC process is implemented at the NU terminal by detecting the strong signal of the FU, which is subtracted afterward from the



**Fig. 3.** Detecting the NOMA signal in the NU and FU receivers after the pairing operation.

superimposed NOMA signal to obtain an interference-free signal of the NU [20], [21].

Meanwhile, on the FU side, the detection of the FU signal is realized directly without applying SIC, by considering the interference caused by the NU, which is considered, due to low PAF, to be additive noise. This process for the two users is illustrated in Fig. 3.

## 5. Simulation Results and Discussions

The simulations described in this paper were conducted using the Matlab programming suite, with the aim of evaluating such performance metrics as the achievable rate and the  $P_{out}$ . The performance of the proposed system relying on the NOMA technique is compared with the traditional OMA approach over different scenarios.

In the first scenario, the MIMO-NOMA technology is considered, in which the UAV-BS is equipped with  $N_t$  and each user has  $N_r$ . We assume that the UAV-BS provides services to users in suburban areas with obstructions in the form of buildings. Therefore, the Rician fading channel between the UAV-BS and GUs is considered with a Rician factor of  $K = 6.93$  dB, and the path loss exponent is assumed to be  $\eta = 4$ . All the parameters used in the simulation are listed in Tab. 4.

In MIMO-NOMA, the pairing algorithm has been applied to choose every two users in one ground. On the GU side, two scenarios need to be considered for detection of the MIMO receiver: in the first case, the NU has a low PAF and it is necessary to apply the SIC operation to remove the interference. The second case, the FU has a high PAF, meaning that no SIC is required. The distances between the UAV-BS and the FU and NU as GUs are  $d_F = 300$  and  $d_N = 150$ , respectively. The power allocation coefficients

**Tab. 4.** Simulation parameters.

Notation	Parameter	Value
$K$	Rician $K$ factor	6.93 dB
$d_F$	Distance between UAV-BS and FU	300 m
$d_N$	Distance between UAV-BS and NU	150 m
$a_F$	Power allocation coefficient for FU	0.88
$a_N$	Power allocation coefficient for NU	0.12
$\eta$	Path loss exponent	4
$d_\lambda$	Antenna spacing	$0.5 \lambda$
$N_t \times N_r$	MIMO antenna	$2 \times 2, 8 \times 2$
$N_s$	Number of transmitted symbols	$10^5$
$(\varphi, \theta)$ UAV-BS	Azimuth and elevation angles for FU and NU	$(30^\circ, 45^\circ)$ and $(90^\circ, 45^\circ)$
$(\varphi, \theta)$ GU	Azimuth and elevation angles for FU and NU	$(60^\circ, 30^\circ)$ and $(60^\circ, 30^\circ)$
$f_o$	Operation frequency	410 MHz to 7.125 GHz
$BW$	Bandwidth	20 MHz
$P_t$	Transmitted power by UAV-BS	0 – 40 dBm
$R_n$	The target rate for NU	4 bit/sec/Hz
$R_F$	The target rate for FU	3 bit/sec/Hz

assigned by the UAV-BS to the FU and NU are  $a_F = 0.88$  and  $a_N = 0.12$ , respectively. In addition, the azimuth and elevation angles equal  $\{30^\circ, 45^\circ\}$  and  $\{90^\circ, 45^\circ\}$  for the FU and NU, respectively. On the GU side, the azimuth and elevation angles are  $\{60^\circ, 30^\circ\}$  and  $\{60^\circ, 30^\circ\}$  towards the FU and NU, respectively. The bandwidth is 20 MHz.

The individual data rates for the FU and NU, with MIMO-NOMA and MIMO-OMA deployed at the UAV-BS, are shown in Fig. 4. The two techniques are assumed to provide services to all users under the same conditions, with changes to the number of antennas used by the UAV-BS and the two users. Two scenarios for the MIMO approach for  $N_t \times N_r$  as  $2 \times 2$  and  $8 \times 2$  are presented. One may notice from this figure that NOMA offers better performance for users with weak channel gains, while OMA is a better choice for users with strong channel gains. Moreover, depending on the data rate required for a specific user within the cluster, the UAV-BS may switch from NOMA to OMA, may rely on a hybrid NOMA-OMA approach or make the required changes on the PAFs of the NOMA system to improve the rate of a particular user. It can also be noticed that increasing the number of antennas at the UAV-BS may significantly improve the rate for all users.

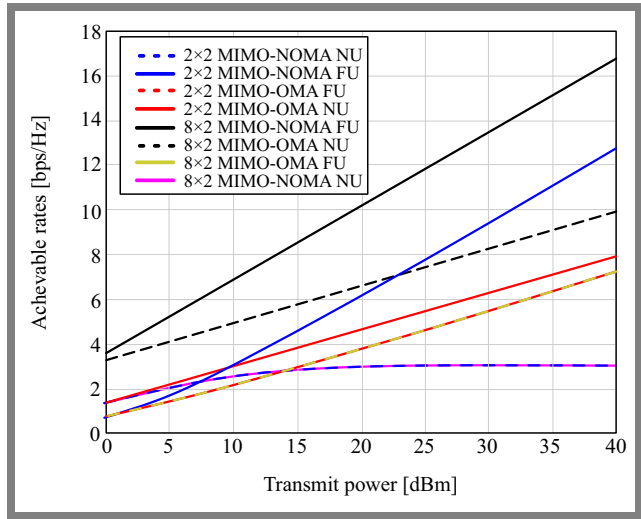
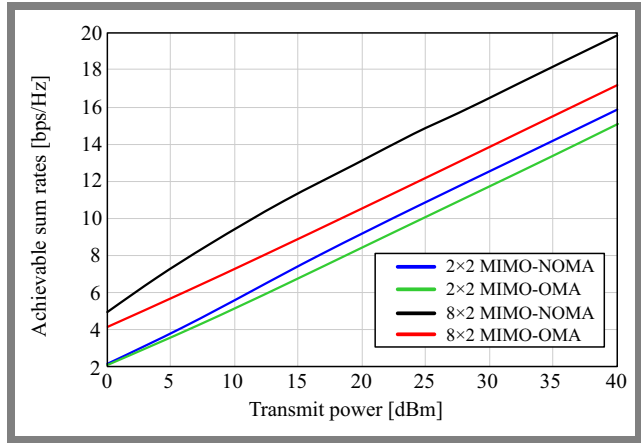

**Fig. 4.** Individual data rates of two users with NOMA and OMA applied for different numbers of  $N_t$  and  $N_r$ .

**Fig. 5.** Sum rates of NOMA and OMA technologies with different numbers of  $N_t$  and  $N_r$ .

Figure 5 shows the sum rates for the proposed system under the same conditions mentioned above and for different numbers of  $N_t$  and  $N_r$ . The transmitted power of the UAV-BS is 40 dBm. NOMA outperforms OMA with approximately 3 bps/Hz and 0.7 bps/Hz for  $N_t$  being equal to 8 and 2, respectively, and for a fixed  $N_r = 2$ .

Figure 6 shows the  $P_{out}$  for the proposed NOMA and OMA systems. The  $P_{out}$  represents the probability that a user obtains a data rate  $R$  that is below the target rate required for an acceptable quality of service (QoS). For this simulation, the required rate for the NU was  $R_n = 4$  bps and the FU was  $R_f = 3$  bps.

One may notice that a lower  $P_{out}$  may be obtained in the FU's receiver in the case of  $8 \times 2$  MIMO-NOMA, since this user has been allocated a higher PAF to overcome its weak channel gain. The second and third best  $P_{out}$  occurred for the cases of  $8 \times 2$  MIMO-OMA in the NU and FU, respectively, due to the entire available power being served to these users by OMA.

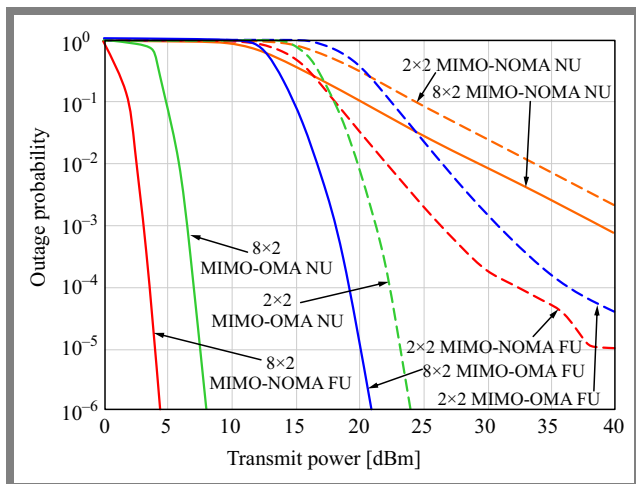
On the other hand, the worst  $P_{out}$  is obtained for the cases of  $2 \times 2$  MIMO-NOMA and  $8 \times 2$  MIMO-NOMA at the NU

**Tab. 5.** Comparison of results obtained with outcomes of related works, for 35 dBm.

Ref.	BS type	No. of users	Path-loss exponent	MIMO ( $N_t \times N_r$ )	Sum data rate users [bps/Hz]		Data rate NOMA [bps/Hz]		Data rate OMA [bps/Hz]		$P_{out}$ NOMA $\times 10^{-6}$		$P_{out}$ OMA $\times 10^{-6}$	
					NOMA	OMA	FU	NU	FU	NU	FU	NU	FU	NU
[18]	Fixed ground BS	4	4	SISO	15.8	NA	NA	NA	NA	NA	NA	NA	NA	NA
				$2 \times 2$	18	NA	8.2	9.5	NA	NA	7943	1584	NA	NA
[24]	Fixed ground BS	2	4	$2 \times 1$	NA	NA	0.5	2.6	0.4	2.2	31	50	158	251
[25]	UAV-RIS	2	3	NA	NA	NA	NA	NA	NA	NA	1	1	1	1
[26]	Fixed ground BS	2	NA	$3 \times 3$	5.9	3.9	NA	NA	NA	NA	NA	NA	NA	NA
		3			3.7	3	NA	NA	NA	NA	NA	NA	NA	NA
[27]	Fixed ground BS	2	Not specified	Not specified	NA	NA	2.2	13.8	NA	NA	1584	125	NA	NA
[22]	Fixed ground BS	2	Not specified	$2 \times 2$	6.1	NA	1.7	10.4	NA	NA	1000	1259	NA	NA
This work	UAV-BS	2	4	$2 \times 10$	NA	NA	NA	NA	NA	NA	200	251	NA	NA
				$2 \times 2$	14.152	13.36	12.9	3	7	7.9	39	6309	125	0
				$2 \times 8$	18.163	15.47	16.9	3	7	9.9	0	2000	0	0

terminal, due to the lower PAF applied to these two scenarios. Additionally, increasing the number of antennas in the UAV-BS can significantly reduce the  $P_{out}$ .

As shown in Tab. 5, the proposed UAV-BS system demonstrates superior performance compared to solutions described in previous works, both in fixed ground BSs and UAV-RIS configurations. At 35 dBm, it achieves significantly higher sum data rates and lower outage probabilities. In the case of a  $2 \times 8$  MIMO configuration, the system reaches a NOMA sum data rate of 18.163 bps/Hz and an OMA sum data rate of 15.47 bps/Hz, outperforming traditional setups. Individual user data rates are also enhanced, particularly in the case of NOMA, where the UAV-BS achieves a result of up to 16.9 bps/Hz for the far user. Furthermore, the proposed system shows a notable reduction in outage probability, reaching zero for both NOMA FU and FU-NU OMA, highlighting improved reliability of the link. These results confirm that UAV-BS systems with large MIMO arrays and beamforming are capable of greatly enhancing network capacity and reliability levels for future wireless communication solutions.



**Fig. 6.** Outage probabilities for FU and NU after applying NOMA and OMA techniques and utilizing different numbers of antennas.

## 6. Conclusions

In this paper, a UAV-BS equipped with multiple antennas has been proposed to serve multiple users. The hybrid OMA-NOMA technique is assumed to provide services to all users in the UAV-BS coverage area. Users with strong and weak channel gains are assigned to the same clusters by near-far pairing algorithms.

The UAV-BS provides services by applying the OMA technique between clusters, while the NOMA approach has been applied to users in each cluster by superimposing their signal after allocating each one of them with a suitable PAF. Different scenarios involving channels between the UAV-BS and the GUs have been considered, depending on whether LoS component of the wave existed or not. AoD and AoA angles for all terminals within this system have been given consideration as well.

The results showed that NOMA outperformed OMA in the cluster under some conditions, depending on the choice of suitable PAF levels for the paired users. In addition, OMA demonstrated good performance in terms of  $P_{out}$  and throughput when compared with NOMA for strong channel gain users. because of lower PAF levels assigned to these users. Furthermore, the UAV-BS can apply a hybrid mechanism by switching from OMA to NOMA, or the other way round, depending on specific situations and circumstances (i.e. depending on the target throughput and  $P_{out}$ ).


For future work, several promising projects are envisaged, such as tracking users by the UAV-BS to enhance the performance of the entire system or optimizing the power level by automatically selecting PAF for each user depending on their status. The use of the mmWave technology may be considered as well to boost the overall capacity of the system. Additionally, mechanisms switching between NOMA and OMA may be developed to avoid service interruptions when one of the techniques remains unavailable.

## References

- [1] S. Hayat, E. Yanmaz, and R. Muzaffar, "Survey on Unmanned Aerial Vehicle Networks for Civil Applications: A Communications Viewpoint", *IEEE Communications Surveys & Tutorials*, vol. 18, pp. 2624–2661, 2016 (<https://doi.org/10.1109/COMST.2016.2560343>).
- [2] N.R. Zema, E. Natalizio, and E. Yanmaz, "An Unmanned Aerial Vehicle Network for Sport Event Filming with Communication Constraints", *First International Balkan Conference on Communications and Networking (Balkancom 2017)*, Tirana, Albania, 2017.
- [3] L.M. Mei, M.S. Johal, F. Idris, and N. Hashim, "Spectrally Efficient UAV Communications Using Non-orthogonal Multiple Access (NOMA)", *Journal of Advanced Research in Applied Sciences and Engineering Technology*, vol. 43, pp. 227–242, 2025 (<https://doi.org/10.37934/araset.43.1.227242>).
- [4] J. Ghosh *et al.*, "Performance Investigation of NOMA versus OMA Techniques for mmWave Massive MIMO Communications", *IEEE Access*, vol. 9, pp. 125300–125308, 2021 (<https://doi.org/10.1109/ACCESS.2021.3102301>).
- [5] W. Jaafar *et al.*, "Multiple Access in Aerial Networks: From Orthogonal and Non-orthogonal to Rate-splitting", *IEEE Open Journal of Vehicular Technology*, vol. 1, pp. 372–392, 2020 (<https://doi.org/10.1109/OJVT.2020.3032844>).
- [6] S.K. Zaidi *et al.*, "Exploiting UAV as NOMA Based Relay for Coverage Extension", *2019 2nd International Conference on Computer Applications & Information Security (ICCAIS)*, Riyadh, Saudi Arabia, 2019 (<https://doi.org/10.1109/CAIS.2019.8769542>).
- [7] A. Li, Y. Lan, X. Chen, and H. Jiang, "Non-orthogonal Multiple Access (NOMA) for Future Downlink Radio Access of 5G", *China Communications*, vol. 12, pp. 28–37, 2015 (<https://doi.org/10.1109/CC.2015.7386168>).
- [8] G. Liu *et al.*, "Hybrid Half-duplex/full-duplex Cooperative Non-orthogonal Multiple Access with Transmit Power Adaptation", *IEEE Transactions on Wireless Communications*, vol. 17, pp. 506–519, 2017 (<https://doi.org/10.1109/TWC.2017.2767601>).
- [9] K. Higuchi and Y. Kishiyama, "Non-orthogonal Access with Random Beamforming and Intra-beam SIC for Cellular MIMO Downlink", *2013 IEEE 78th Vehicular Technology Conference (VTC Fall)*, Las Vegas, USA, 2013 (<https://doi.org/10.1109/VTCFall.2013.6692307>).
- [10] A. Abd Saeed and M.A. Ahmed, "Cognitive Radio Based NOMA for the Next Generations of Wireless Communications", *2022 International Conference on Electrical Engineering and Informatics (ICELTICs)*, Banda Aceh, Indonesia, 2022 (<https://doi.org/10.1109/ICELTICs56128.2022.9932105>).
- [11] M.F. Sohail, C.Y. Leow, and S. Won, "Non-orthogonal Multiple Access for Unmanned Aerial Vehicle Assisted Communication", *IEEE Access*, vol. 6, pp. 22716–22727, 2018 (<https://doi.org/10.1109/ACCESS.2018.2826650>).
- [12] X. Hu, C. Zhong, and Z. Zhang, "Angle-domain Intelligent Reflecting Surface Systems: Design and Analysis", *IEEE Transactions on Communications*, vol. 69, pp. 4202–4215, 2021 (<https://doi.org/10.1109/TCOMM.2021.3064328>).
- [13] N. Ben Halima, "Non Orthogonal Multiple Access (NOMA) Using a Nonlinear Energy Harvesting Model", *Wireless Personal Communications*, vol. 135, pp. 2165–2175, 2024 (<https://doi.org/10.1007/s11277-024-11129-9>).
- [14] M. Jain, S. Soni, N. Sharma, and D. Rawal, "Performance Analysis at Far and Near User in NOMA Based System in Presence of SIC Error", *AEU-International Journal of Electronics and Communications*, vol. 114, art. no. 152993, 2020 (<https://doi.org/10.1016/j.aeu.2019.152993>).
- [15] J.S. Roy, "Multiple-antenna Techniques in Wireless Communication Technical Aspects", *International Journal of Information Communication Technology and Digital Convergence*, vol. 1, pp. 24–32, 2016 (<https://doi.org/10.17577/IJERTCONV4IS01013>).
- [16] X. Luo, H. Li, Y. Bai, and S. Wei, "Research on Power Allocation Algorithm in Non-orthogonal Multiple Access Systems", 2019 14th IEEE Conference on Industrial Electronics and Applications (ICIEA), Xi'an, China, 2019 (<https://doi.org/10.1109/ICIEA.2019.8834152>).
- [17] A. Al-Hourani, S. Kandeepan, and A. Jamalipour, "Modeling Air-to-ground Path Loss for Low Altitude Platforms in Urban Environments", *2014 IEEE Global Communications Conference*, Austin, USA, 2014 (<https://doi.org/10.1109/GLOCOM.2014.7037248>).
- [18] A.A. Saeed and M.A. Ahmed, "Algorithms and Throughput Analysis of Cognitive Radio-based MIMO-NOMA for the Next Generation of Wireless Communications", *Transactions on Emerging Telecommunications Technologies*, vol. 35, art. no. e4988, 2024 (<https://doi.org/10.1002/ett.4988>).
- [19] J. Miranda *et al.*, "Path Loss Exponent Analysis in Wireless Sensor Networks: Experimental Evaluation", *2013 11th IEEE International Conference on Industrial Informatics (INDIN)*, Bochum, Germany, 2013 (<https://doi.org/10.1109/INDIN.2013.6622857>).
- [20] A. Doukas and G. Kalivas, "Rician K Factor Estimation for Wireless Communication Systems", *2006 International Conference on Wireless and Mobile Communications (ICWMC'06)*, Bucharest, Romania, 2006 (<https://doi.org/10.1109/ICWMC.2006.81>).
- [21] Z. Ding, Z. Yang, P. Fan, and H.V. Poor, "On the Performance of Non-orthogonal Multiple Access in 5G Systems with Randomly Deployed Users", *IEEE Signal Processing Letters*, vol. 21, pp. 1501–1505, 2014 (<https://doi.org/10.1109/LSP.2014.2343971>).
- [22] A.A. Saleh and M.A. Ahmed, "Performance Enhancement of Cooperative MIMO-NOMA Systems Over Sub-6 GHz and mmWave Bands", *Journal of Telecommunications and Information Technology*, no. 2, 2023 (<https://doi.org/10.26636/jtit.2023.170023>).
- [23] Z. Cui *et al.*, "Wideband Air-to-ground Channel Characterization for Multiple Propagation Environments", *IEEE Antennas and Wireless Propagation Letters*, vol. 19, pp. 1634–1638, 2020 (<https://doi.org/10.1109/LAWP.2020.3012889>).
- [24] R. Chandrasekhar *et al.*, "Performance Evaluation of MIMO-NOMA for the Next Generation Wireless Communications", *2021 3rd International Conference on Signal Processing and Communication (ICSPC)*, Coimbatore, India, 2021 (<https://doi.org/10.1109/ICSPC51351.2021.9451780>).
- [25] S. Li, X. Liu, J. Gaber, and G. Pan, "Modeling Analysis for Downlink RIS-UAV-assisted NOMA over Air-to-ground Line-of-sight Rician Channels", *Drones*, vol. 8, art. no. 659, 2024 (<https://doi.org/10.3390/drones8110659>).
- [26] M. Zeng *et al.*, "Capacity Comparison Between MIMO-NOMA and MIMO-OMA with Multiple Users in a Cluster", *IEEE Journal on Selected Areas in Communications*, vol. 35, pp. 2413–2424, 2017 (<https://doi.org/10.1109/JSAC.2017.2725879>).
- [27] A. Abd Saeed and M.A. Ahmed, "Cognitive Radio Based NOMA for the Next Generations of Wireless Communications", *2022 International Conference on Electrical Engineering and Informatics (ICELTICs)*, Banda Aceh, Indonesia, 2022 (<https://doi.org/10.1109/ICELTICs56128.2022.9932105>).

**Ameer Y. Sadeeq, Student**

College of Electronics Engineering


 <https://orcid.org/0009-0006-5514-541X>

E-mail: ameer.yaseen@stu.uoninevah.edu.iq

Ninevah University, Mosul, Iraq

<https://uoninevah.edu.iq>**Mohamad A. Ahmed, Ph.D.**

College of Electronics Engineering

 <https://orcid.org/0000-0001-6412-2275>

E-mail: mohamad.alhabbar@uoninevah.edu.iq

Ninevah University, Mosul, Iraq

<https://uoninevah.edu.iq>

# Using Modified Gorti-enhanced Homomorphic Cryptosystem to Improve Security of ECG Signal

Fatma Zohra Besmi, Samia Belkacem, and Nouredine Messaoudi

University of Boumerdes, Boumerdes, Algeria

<https://doi.org/10.26636/jtit.2025.2.2002>

**Abstract** — While offering vast data storage capabilities, cloud computing poses numerous security- and privacy-related challenges. This requires robust security measures, particularly for sensitive data, such as electrocardiograms (ECG). Homomorphic encryption (HE) emerges as a promising solution by enabling secure computations to be performed directly on encrypted data. This study introduces a novel approach to enhance the security of ECG data. We modified the Gorti-enhanced homomorphic cryptosystem (MEHC) method by optimizing its key generation procedure and then applied the linear congruential generator (LCG) algorithm to create a list of huge prime integers. Furthermore, we increased the modulus value and enlarged the message space. These enhancements boosted overall security by substantially improving immunity to factorization attacks. We used quantization and fixed-point representation to enhance the encryption process. As an additional security layer, an evaluation process has been added to the proposed algorithm which performs various mathematical operations homomorphically on the encrypted data, rather than on the original data. This modified algorithm enables efficient and secure encryption of ECG data while preserving the ability to reliably identify arrhythmias, such as bradycardia and tachycardia. Using the MIT-BIH arrhythmia database, the proposed MEHC system demonstrated high accuracy (98.48%), sensitivity (99.10%) and positive predictive value (99.33%), while effectively safeguarding the ECG data. These results validate the efficacy of the MEHC system and confirm its suitability for secure and reliable ECG signal processing in healthcare applications.

**Keywords** — arrhythmia detection, cryptosystem, decryption, ECG, encryption, enhanced homomorphic cryptography

## 1. Introduction

An electrocardiogram (ECG) is an important method commonly used in medicine to track electrical changes in a patient's body linked to the beating of their heart. The signals are measured using electrodes placed on 12 standard wires (also known as channels) and applied to specific areas of the patient's chest, arms, and legs [1]. There are three basic components to an ECG signal: the P-wave, the QRS wave, and the T-wave [2]. ECG signals are studied by assessing the locations or magnitudes of PR and ST segments, QRS, PR, QT, and ST intervals, as well as additional data [3].

Comprehending these waves and intervals is essential in detecting problems affecting the cardiovascular system and assessing the overall condition of the heart. The Pan-Tompkins approach is the predominant technique for the diagnosis of ECG abnormalities [4].

Due to the fact that ECGs contain patient information, solutions (such as encryption) are required to protect patient data and maintain its quality, in order to avoid risks that lead to erroneous diagnoses or treatment plans. One of the recommended methods for encrypting ECGs is fully homomorphic encryption (FHE).

This type of encryption allows for an infinite number of operations to occur at any given moment; mixed operators can perform any number of operations on the ciphertext. FHE may be of the additive, multiplicative, or mixed variety [5]. Therefore, it offers greater security for cloud data and services.

The Gorti-enhanced homomorphic cryptosystem (EHC), introduced by Gorti and Garimella in [6], is a specific type of fully homomorphic encryption (FHE). EHC offers secure indistinguishability under a chosen ciphertext attack (IND-CCA), implying that an attacker cannot discern the keys for a selected ciphertext [7].

This scheme is characterized by better performance than previous systems, as it discovers and utilizes two distinct keys for encryption and decryption: a secret key  $q$ ,  $p$  and a public key  $m$  [8]. The security of the scheme relies on the difficulty of factoring the large number  $m$  (at least 1024, preferably 2048 bits), and the random number  $r$  adds an additional layer of security, making it harder to deduce the original message.

This research introduces a novel approach to enhancing the security of ECG signals in healthcare applications. Initially, the Pan-Tompkins approach is used to examine and process an electrocardiogram (ECG) signal in order to identify the QRS complex. This method provides strong detection performance when used with precise clinical ECG signal data. However, ECG recordings made at outpatient clinics, low quality of the signals and the noise present limit the ability of this algorithm to identify QRS complexes [9].

To address this limitation, we increase the bandpass filter's upper cutoff frequency to 25 Hz (resulting in a 5–25 Hz passband), removing noise while retaining significant signal components. To improve the efficiency of the encryption

process, signal ( $X_6$ ) is initially converted to an integer representation using a fixed-point representation. Subsequently, quantization is applied to the integer values, optimizing the data representation without sacrificing signal fidelity.

The resulting quantized signal is encrypted using the proposed MEHC algorithm which employs a robust key generation mechanism combining a linear congruential generator (LCG) and the Miller-Rabin primality test [10], where the sender applies the newly generated secret key  $p$ ,  $q$ ,  $Q$  to encrypt the ECG data. A key feature of this approach is the enabling of secure homomorphic operations, facilitated by a deterministic evaluation key derived via SHA-256 hashing of the secret key. This allows direct computation on the encrypted data through a linear polynomial function, a core characteristic of FHE.

After the encrypted evaluation is complete, the signal is decrypted using public key  $g$ . This ensures that the transmitted ECG data remains private and is protected against unauthorized access. The signal is re-quantized and converted back to its original floating-point representation after decryption, thereby preparing it for subsequent analysis and classification.

The proposed method improves robustness against factorization attacks, contributing to improved efficiency, reliability, and analysis of the transmitted ECG signal.

This study is structured as follows. Section 2 reviews related work, whereas Section 3 provides an in-depth description of the proposed algorithm. Section 4 analyzes the experimental findings and compares them with previous work, and Section 5 concludes the article.

## 2. Related Works

Before discussing the proposed approach, it is imperative that we first review the current methods used to protect the ECG signal. This section highlights relevant encryption techniques that allow calculations to be performed on encrypted data while maintaining secrecy, and discusses the strengths and limitations of these methods to provide context for our contributions.

To guarantee the security of medical data during their transmission, the authors of [11] suggest an encryption method for ECG signals based on the partial homomorphic encryption technique (PHE). To encrypt the signals, the study combines the Pan-Tompkins algorithms for QRS complex detection with the PHE-RSA method. This approach can detect cardiac abnormalities and calculating the heart rate while protecting sensitive data from unauthorized access.

The study was carried out using MIT-BIH arrhythmia data. Of 20 recordings, the results indicate that 18 of them had the same outcome. the method achieved a 90% accuracy rate and was quick, demonstrating its effectiveness in securing ECG signals. However, while focusing on enhanced security, the study lacks a deeper exploration of potential vulnerabilities or attack scenarios, which is crucial to understanding the system's robustness.

In article [12] a new system based on fully homomorphic encryption (FHE) methods was proposed. The strength of this approach lies in its ability to achieve a high sensitivity of 92.59% and a positive predictive accuracy of 90% for detecting arrhythmias, while simultaneously maintaining data privacy. This method faces some limitations, including computational complexity, in addition to a lack of detailed comparison with other FHE schemes, which makes it difficult to evaluate the advantages and disadvantages of the proposed approach. Addressing these limitations is critical to realize the full potential in secure healthcare data sharing.

In [13], a new fully homomorphic encryption algorithm with an advanced encryption standard (FHEAES) was introduced to encode electrocardiogram (ECG) signals for improved security and data privacy. The ECG measurements were processed using the Pan-Tompkins algorithms, followed by encoding using the FHEAES algorithm. The proposed algorithm integrates an evaluation process as an additional security layer, enabling mathematical operations to be performed on encrypted data without decryption.

The algorithm was evaluated using various ECG signals obtained from the MIT-BIH database and demonstrated a 95.8% accuracy rate in classifying heart rates and 100% accuracy in decryption, highlighting the effectiveness of homomorphic encryption in securing ECG signals. In general, while FHEAES offers robust security and privacy features, its computational complexity and key management requirements present challenges that need to be addressed for optimal implementation in real-world applications.

Paper [14] presents a new approach to secure transmission of ECG signals over the Internet by combining set partitioning in hierarchical trees (SPHIT) with RSA encryption. The approach aims to address bandwidth and security issues related to the transmission of sensitive medical data while maintaining the quality of the ECG signal. The use of RSA encryption provides a strong foundation for protecting confidentiality of patient data and ensuring secure communication. However, it is not appropriate for large signals or for real time use. Instead, the FHE encryption algorithm is capable of performing simultaneous addition and multiplication operations, thus adding further security to the signal.

The authors of [15] tested the feasibility of analyzing ECG data in the cloud using FHE, aiming to secure patient privacy while enabling remote healthcare services. The suggested technique is based on Gentry's FHE and BGV techniques. The findings indicate that while FHE is promising, it suffers from major performance constraints, especially since decryption time dominates addition operations and every bit of encryption necessitates a large amount of storage.

These limitations, as exemplified by the 800 000-fold increase in storage space needed by Gentry's FHE scheme for one hour of patient data, highlight the need for further exploration of alternative secure computation methods that are capable of balancing security with practical efficiency for real-world healthcare applications.

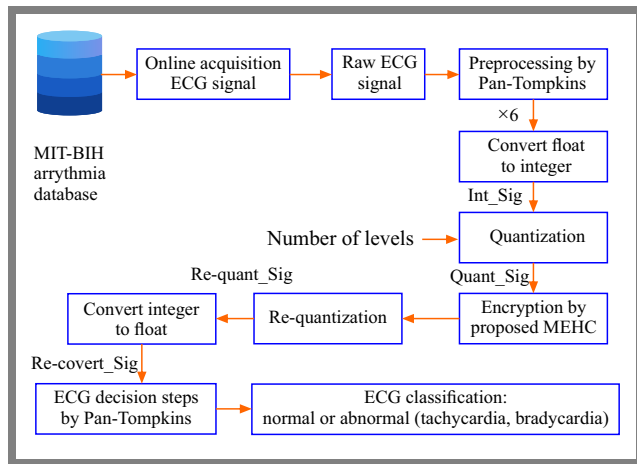


Fig. 1. Flow chart of the proposed approach.

### 3. Proposed Method

This research is divided into four steps: preprocessing, encryption, decision-making, and classification. Initially, the Pan-Tompkins technique was used for signal pre-processing. Subsequently, the signal (X6) was transformed into an integer representation and quantized to enhance the representation of the data while preserving signal fidelity.

In the second part, the proposed MEHC algorithm was used to encrypt the ECG signal, which improves security and makes it more difficult for attackers to gain access. Subsequently to decryption, the signal underwent re-quantization and was restored to its original format. In the decision-making stage, this resulting signal was analyzed to detect QRS complexes, a prominent feature in ECG assessments.

The characteristics of these features are used to calculate the heart rate, which is then classified. This ensures the privacy and integrity of sensitive medical information, improving the efficiency and reliability of ECG signal transmission and analysis.

Figure 1 shows the methodology of this study presented in the form of a flowchart. We use raw ECG data from the MIT-BIH database [16]. Details regarding the specific dataset used are provided in Subsection 4.1.

#### 3.1. Preprocessing Signal

The pre-processing stage (Fig. 2) involved several steps to prepare the ECG signal for subsequent encryption and analysis. The Pan-Tompkins algorithm was applied for initial detection of the QRS complex [17]. To improve the algorithm’s robustness against noise and low-quality signals, we increased the bandpass filter’s upper cutoff frequency to 25 Hz. This specific bandpass was chosen to ensure effective noise removal while preserving critical signal components, thus striking a necessary balance for reliable QRS detection under various recording conditions.

The processing of the ECG signal within this methodology is conducted in the following manner: the raw ECG signal (X1) is subjected to a bandpass filter, comprising a low pass (X2) and a high-pass (X3) stage, with a passband of 5–25 Hz.

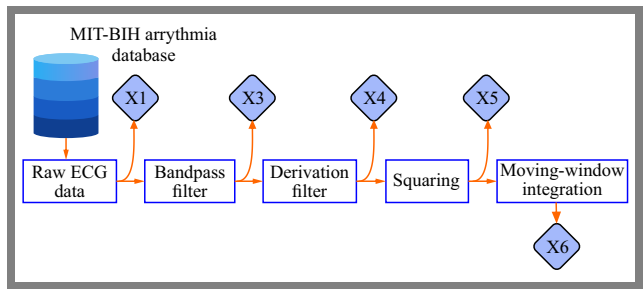


Fig. 2. Preprocessing ECG signal diagram.

Following the bandpass filter (X3), the signal-to-noise ratio is improved, consequently enhancing the overall sensitivity of the detector. Subsequently, the filtered signal undergoes a series of transformations, including differentiation (X4), a squaring function (X5), and integration of the moving window (X6).

#### 3.2. Fixed-point Representation and Quantization

After receiving the output of the moving window coordinates (X6), a two-step transformation (fixed-point representation and quantization) is performed to optimize it for the encryption process. By adopting these transformations, the system can efficiently manage the data without sacrificing its integrity.

It is a technique used to represent floating-point numbers using integers [18]. In this study, the signal (X6) is transformed into a discrete integer representation using a scaling factor. This requires selecting appropriate scale factors for each data point based on its magnitude. By scaling the values, we can preserve precision while decreasing the number of bits necessary for representation.

Once the relevant scale factors are determined, each floating-point value in the (X6) signal is multiplied by its corresponding scale factor and rounded to the nearest integer. This results in an integer signal (Int\_Sig) that precisely describes the (X6) signal.

To further optimize data representation, quantization is performed for integer values. This approach attempts to enhance the retention of essential signal properties while lowering distortion, thereby boosting performance of subsequent analytical tasks [19]. This technique is widely applied in audio and voice compression, image processing, and biological signal analysis, particularly in ECG signal evaluation. The following formulae were applied to determine the signal after quantization (Quant\_Sig) [20]:

$$Y_i = \left\lfloor \frac{X_i - X_{min}}{X_{max} - X_{min}} \right\rfloor \times 2^{n-1}, \quad (1)$$

where  $Y_i$  and  $X_i$  represent the  $i$ -th sample in quantized array  $Y$  and original signal  $X$ .  $X_{max}$  and  $X_{min}$  represent the maximum and minimum values of signal  $X$  (here the original signal is Int\_Sig).

Following acquisition of the quantization signal (Quant\_Sig), encryption is performed using the proposed MEHC system.

### 3.3. Modified Enhanced Homomorphic Cryptosystem (MEHC)

Based on the Gorti-enhanced homomorphic cryptosystem (EHC), we incorporate a modified MEHC scheme. This fully homomorphic public-key encryption scheme utilizes a pair of public and private keys.

Key modifications implemented to enhance security include using a larger modulus value (related to public key  $g$ ), which strengthens immunity to factorization attacks, and incorporating a prime number  $Q$  into the secret key to enlarge the message space. The robust key generation process integrates the Miller-Rabin primality test and the linear congruential generator. This modified scheme ensures data privacy by executing various computational operations (addition and multiplication) on the ciphertext without necessitating decryption, enabled by a deterministic evaluation key derived via SHA-256 hashing of the secret key. These changes aim to improve the performance and enhance security, making it particularly beneficial for applications demanding enhanced data privacy and security.

The approaches used in this study are presented below.

The most common approach for obtaining random numbers is a technique known as the linear congruential generator (LCG). We utilized it to create a massive list of odd integers for a subsequent primality test [21].

$$X_{n+1} = (aX_n + C) \pmod{m}, \quad n \in \mathbb{Z}^+. \quad (2)$$

The seed value of the series is  $X_0$ , while the multiplier is  $a$ , the increment is  $C$ , the created modulus is  $m$ , and the random numbers are denoted by  $X_n$ . The quality of the generated sequence depends heavily on the choice of  $a$ ,  $C$ , and  $m$  parameters. To ensure a long period and good statistical properties, the following conditions should be met:

$$\gcd(C, m) = 1, \quad (3)$$

$$a \equiv 1 \pmod{p} \quad \text{for every prime } p \text{ dividing } m, \quad (4)$$

$$a \equiv 1 \pmod{4} \quad \text{if } m \text{ is a multiple of } 4. \quad (5)$$

With  $m = 2^k$ ;  $a = 4b - 1$ ;  $C$  as an odd number ( $b; k > 0$ ). A satisfactory result can be obtained by setting the increment  $C$  to zero [22].

The Miller-Rabin primality-testing algorithm is a primality test that assesses if a specific number is likely to be prime. This test is based on Fermat's little theorem [10].

The proposed MEHC algorithm consists of four processes, as described below.

The proposed algorithm generates three keys: public key  $pk$ , private key  $sk$ , and evaluation key  $ek$ . The algorithm follows the following steps to create these keys:

- 1) Generate two large prime numbers  $p$  and  $q$  using LCG, such that  $p > q$ . The primality of these numbers is verified using the Miller-Rabin test.
- 2) Calculate modulus  $m = p \times q$
- 3) Calculate  $g = m^2 + 1$

**Algorithm 1** Miller-Rabin primality test to verify the primality of an odd number  $n$

- 1: Find integers  $r$  and  $m$  such that  $n - 1 = 2^r \times m$
- 2: Choose randomly any integer  $a \in [1, n - 1]$
- 3: Compute:  $b_0 = a^m \pmod{n}$
- 4: Compute  $b_i, k$  times  $b_i = b_{i-1} - 1$
- 5: The result must  $b_e = \pm 1$
- 6: **if** the result is 1 **then**
- 7:      $n$  is a composite number
- 8: **end if**
- 9: **if** the result is  $-1$  **then**
- 10:      $n$  is a prime number
- 11: **end if**

- 4) Select a prime number  $Q$  that satisfies the condition  $Q|_2^p|_-$ . The primality of this number is verified by the Miller-Rabin test.

Based on these generated values, the keys are defined as follows:

- Public key  $g$ . This value is designed to obscure the original value  $m$ .
- Private key  $p, q, Q$ .
- Evaluation key  $ek$  is derived by computing the SHA-256 hash of the combined parameters  $p, q$ , and  $Q$ , and then reducing the resulting hash value modulo  $g$ .

The encryption procedure takes the quantized ECG signal (Quant\_Sig) as input and generates the corresponding ciphertext  $C$ . This process utilizes scheme parameters  $Q$  and  $g$ , along with a random vector  $r$  for probabilistic security. The encryption is calculated using the following formula:

$$C = (\text{Quant\_Sig} + r \times Q + r \times g) \pmod{g}. \quad (6)$$

The evaluation process receives ciphertexts  $C$  that originate from the encryption process as input and creates new ciphertexts  $C'$  as output, which will be termed evaluated ciphertexts. An assessment procedure is performed on the server side before decryption. This fundamental procedure was suggested to be added to the MEHC algorithm in order to accomplish homomorphic encryption. It conducts mathematical operations homomorphically on the ciphertext.

Our approach converts the generic evaluation function  $f$  into a linear polynomial. This decision is driven by the simplicity of linear polynomials, which improves performance, as well as its ability to provide better security. A linear polynomial is characterized by constant real values and has a degree of one [23]. The suggested evaluation function generates evaluated ciphertexts  $C'$  using three parameters: random integer  $N$ , evaluation key  $ek$  and new modulo  $g$ . The following formula represents the procedure:

$$(C') \leftarrow \text{Eval}_{N,ek,g}(f, C), \quad (7)$$

$$C' = (N \times C + ek) \pmod{g}. \quad (8)$$

$N$  was chosen in the MEHC algorithm to be a variable number rather than a constant number to increase the level of security by making it more difficult to hack or break. Such a higher

degree of protection results from the fact that the variable number itself costs the hacker or third party. Before putting it in the encryption equation, the user must first ascertain whether the integer is variable or not.

This requires an additional effort to crack the encryption.

The proposed algorithm is fully homomorphic, which implies that it can handle both multiplicative and additive homomorphic properties, as demonstrated in Eq. (8). The correct decoding of the ciphertext is essential for the success of homomorphic encryption.

The decryption process involves using the inverse of the evaluation function (a linear polynomial function) followed by the MEHC decryption algorithm with the secret key to recover the plaintext from the ciphertext in the following manner:

$$C = ((C' - ek) \times N^{-1}) \pmod{g}, \quad (9)$$

$$Dec\_Sig = C \pmod{Q}. \quad (10)$$

### 3.4. Decision Making

Once decrypted, the signal is converted back to its original floating-point representation, preparing it for this decision phase. This step involves labeling signal peaks as QRS complexes by applying several thresholds, distinguishing them from noise peaks, which may include muscle noise and T waves. Figure 3 presents a diagram that illustrates the decision-making phase. To be identified as a QRS complex, a peak must be recognized as such a complex in both the decrypted signal and the bandpass-filtered waveforms. Specifically, the following threshold values were used [24]:

$$SPK = 0.125 \times Peak + 0.875 \times SPK, \quad (11)$$

$$NPK = 0.125 \times Peak + 0.875 \times NPK, \quad (12)$$

$$Threshold\ 1 = NPK + 0.125 \times (SPK - NPK), \quad (13)$$

$$Threshold\ 2 = 0.5 \times Threshold\ 1. \quad (14)$$

SPK and NPK are the peak values of the running estimates of the signal and noise, respectively. Value threshold 2 is used only if threshold 1 fails to find a QRS complex within a certain distance from the detected one.

The thresholds were automatically changed to adapt to changes in QRS shape and heart rate. Once the QRS complex is detected, the  $R$  position is specified, leading to the RR interval, and computing the heart rate (HR) according to the equation [21].

$$HR = \frac{\text{Number of } R \text{ peaks}}{\text{Time of ECG signal [s]}} \times 60. \quad (15)$$

The heart rate was classified as normal or abnormal HR. Normal heart rate varies from 60 to 100 bpm, but irregular heartbeats indicate arrhythmia (either bradycardia or tachycardia), which helps the physician to make the appropriate

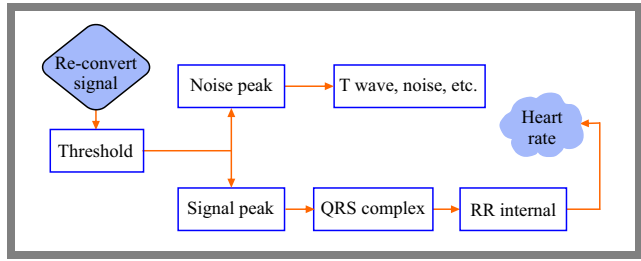


Fig. 3. Diagram of the decision-making phase.

inference. It also helps the analyst to perform further analysis and detection-related steps.

### 3.5. Parameter Estimation

The reference annotation and the WFDB toolbox provided by the Python suite were used for this evaluation [15]. Several values are introduced to assess the effectiveness of the presented algorithm, including:

- True positive (TP) – correctly identified QRS complexes.
- False positive (FP) – peaks incorrectly identified as QRS complexes.
- False negative (FN) – missed QRS complexes.
- True negative (TN) – peaks correctly identified non-QRS peaks.

Performance of the algorithm is evaluated by calculating its accuracy, sensitivity, positive prediction, and detection error rate, with all of the aforementioned terms defined as follows.

- Accuracy (Acc) assesses the algorithm's overall correctness in identifying both QRS complexes and non-QRS peaks.

$$Acc = \frac{TP + TN}{\text{Total beats}} \times 100 [\%]. \quad (16)$$

- Sensitivity (Se) is its ability to accurately identified QRS:

$$Se = \frac{TP}{TP + FN} \times 100 [\%]. \quad (17)$$

- Positive prediction represents the probability of QRS identified among the true QRS.

$$+P = \frac{TP}{TP + FP} \times 100 [\%]. \quad (18)$$

- Detection error rate (DER) is a measure that indicates the overall error rate of the algorithm.

$$DER = \frac{FP + FN}{\text{Total beats}} \times 100 [\%]. \quad (19)$$

## 4. Results and Discussion

This section presents the results obtained by applying the proposed MEHC scheme to process ECG signal in a secure manner. We detail the results achieved at various stages of our approach – from preprocessing, to encryption and secure classification of the decrypted signal. The effectiveness of the encryption method and the overall effectiveness in identifying cardiac abnormalities are evaluated using various indicators. The ECG signals were obtained from the MIT-BIH arrhythmia database containing 48 heartbeat signal recordings, with

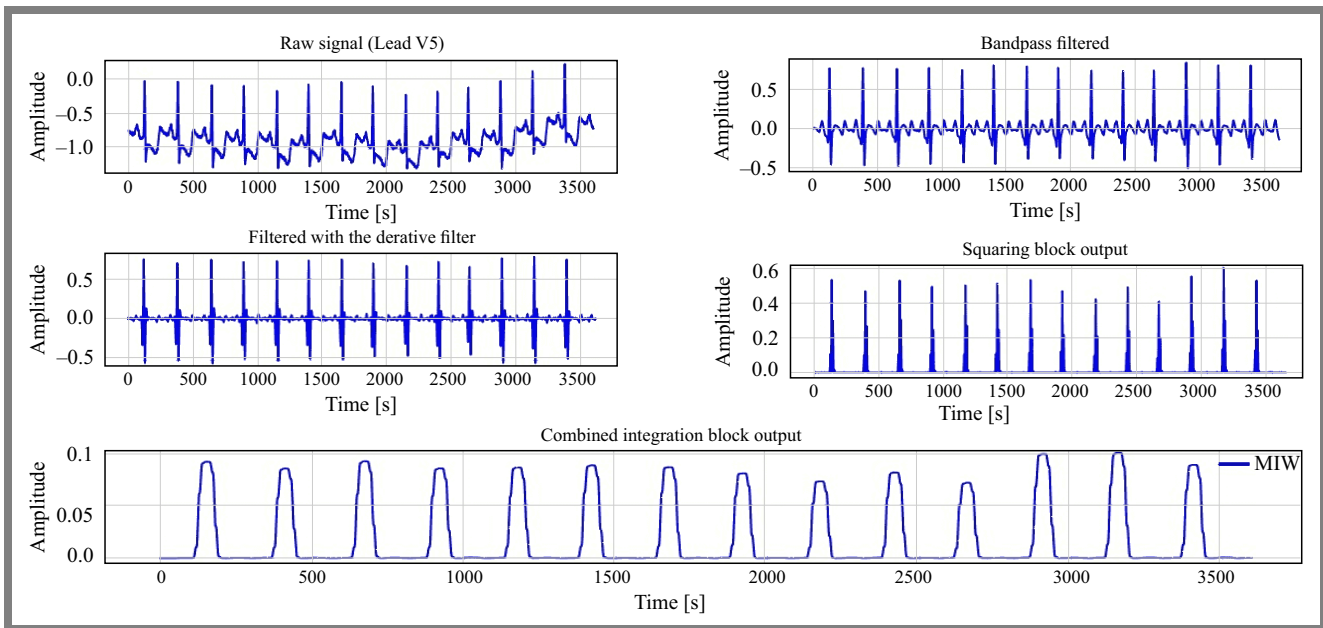


Fig. 4. ECG processing results involving record no. 112.

a sampling rate of 360 Hz. In each recording, the first channel represents the modified lead II (MLII), while the subsequent channel is designated as one of the leads v1, v2, v4, or v5, depending on a specific recording. Given that MLII is consistently present in all recordings and demonstrates great accuracy [25], we used MLII data for this investigation.

#### 4.1. Pre-processing Signal

This stage consists of two main processes: data acquisition and pre-processing. Figure 4 shows the procedures used in processing record no. 112. To remove the noise and existing artifacts, bandpass filtering was employed. The following phase involved the use of differentiation to determine the high slope to differentiate QRS complexes from other ECG waves. To make all the data positive and highlight the higher frequencies of the signal, the sample was first squared gradually. After this square, the waveform passed through the moving window integrator.

#### 4.2. Encryption/Decryption Process

To prepare the ECG signal for encryption, a two-step transformation procedure is applied. First, the signal is transformed into an integer representation using fixed-point arithmetic. This involves scaling the signal values to make sure that they fit within a certain integer range. Subsequently, adaptive quantization is performed on the scaled integer values, dividing the signal range into 256 discrete levels to maximize the representation of data.

Following quantization, the ECG signal is encrypted using the proposed MEHC scheme. The scheme relies on parameters derived from two large 512-bit prime numbers  $p$  and  $q$ , so the public key component  $g$  is approximately 2048 bits. The substantial size of these parameters, particularly of  $g$ , contributes significantly to the security of the scheme. Figure 5

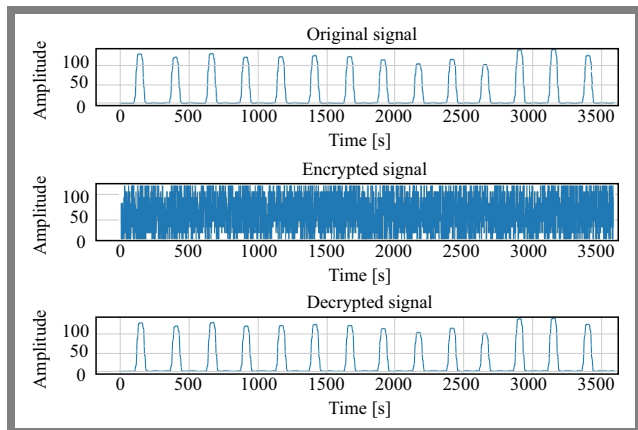


Fig. 5. ECG encryption/decryption results (record no. 112)

illustrates the signal at different stages, including the pre-processed signal (X6) and the resulting encrypted and decrypted signals after applying the MEHC scheme.

We noticed that the original signal is completely different from the encrypted signal. Nevertheless, upon decryption, we acquire a signal that is identical to the original signal. This observation highlights the effectiveness of the encryption process in protecting data while maintaining its integrity. To complement this visual assessment and further validate security against statistical attacks, we performed a detailed analysis.

#### 4.3. Analysis of Statistical Attacks

Using signal diffusion in the encrypted signal, intruders attempt to anticipate the original signal and secret keys in a statistical attack. The histogram, the correlation coefficient, and the mean square error (MSE) are analyzed for verification of the statistical attack. To evaluate the effectiveness of the encryption scheme in preventing statistical attacks, we ana-

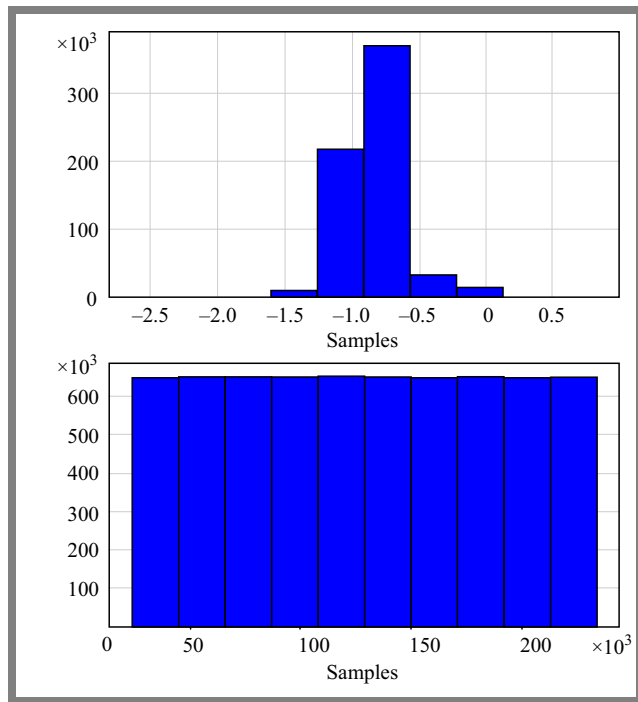


Fig. 6. Analysis of original and encrypted signal of sample no. 112.

alyzed the histogram of the encrypted signal. A well-encrypted signal should exhibit a uniform distribution that resembles random noise.

Figure 6 clearly demonstrates that the proposed approach meets this criterion. The histogram of the encrypted ECG signal is significantly different from that of the original signal and displays a uniform distribution. This indicates that the encrypted signal is statistically indistinguishable from random noise, making it resistant to statistical attacks.

The correlation coefficient is a statistical metric that assesses the degree and direction of the linear association between two variables. In cryptography, it can be used to evaluate the security of an encryption algorithm and its ability to fend off statistical attacks.

A correlation value of  $-1$  signifies a perfect negative linear connection,  $1$  signifies a perfect positive linear relationship, and  $0$  signifies the absence of a linear link.

We assessed the security of the proposed encryption scheme by calculating the correlation coefficients between the original and encrypted ECG signals. We analyzed all possible column and row pairings involving raw and encrypted ECG data [26]. The findings shown in Tab. 1 indicate that the correlation coefficients between the raw and encrypted signals are markedly low. This means that the encrypted signal is considerably different from the original signal, making it very resistant to statistical attacks. Consequently, the proposed technique is successful in ensuring the security and confidentiality of ECG data.

Mean squared error (MSE) was used to measure the distortion induced by encryption, homomorphic evaluation, and decryption stages [13]. The MEHC method produced a surprisingly low average MSE of roughly  $3.43 \cdot 10^{-8}$ . This tiny amount of error is a critical result, indicating that calculations per-

formed fully inside the encrypted domain generate only minor numerical errors. Such a high degree of signal fidelity after decryption is a therapeutic need. It guarantees that crucial ECG properties, including the exact amplitudes and durations of the P waves, QRS complexes, and T waves, as well as the critical PR, QRS, and QT intervals, are retained. This preservation directly affects the accuracy of automated arrhythmia detection systems and the reliability of future clinical interpretation. Therefore, the proposed MEHC scheme efficiently reconciles the need for robust data security in cloud-based ECG processing with the need to preserve high signal quality that is necessary for accurate and reliable medical diagnosis.

#### 4.4. Decision Making

After decryption, the ECG signal undergoes dequantization followed by conversion to its original floating-point form. The QRS complex identification process is then performed on the restored signal. Figure 7 illustrates the QRS complexes on a filtered ECG signal. The QRS complexes are marked by purple circles. The graph includes a background noise level

Tab. 1. Results of the correlation coefficients of encrypted signals.

Record no.	Correlation coefficient	Record no.	Correlation coefficient
100	0.001151854	201	0.000885001
101	-0.000578912	202	0.000250119
102	-0.000669556	203	0.000482145
103	0.001295191	205	0.000742561
104	0.001038385	207	0.001767004
105	0.000953979	208	-0.001236361
106	0.000304367	209	0.00297388
107	0.000532022	210	0.000755799
108	0.0000125535	212	0.000441197
109	0.003018805	213	0.000711752
111	0.001969152	214	0.000676252
112	-0.001360556	215	0.000885888
113	0.001442974	217	0.002380301
114	0.000284177	219	0.00088243
115	-0.000108819	220	0.000449283
116	0.002244413	221	-0.000336268
117	-0.000393281	222	0.000473388
118	0.002158905	223	0.001647422
119	0.00126128	228	0.003093744
121	0.00087087	230	-0.000017272
122	0.002106787	231	0.002384566
123	0.001947021	232	0.001998768
124	-0.000769574	233	0.000740326
200	0.000891936	234	0.001472958

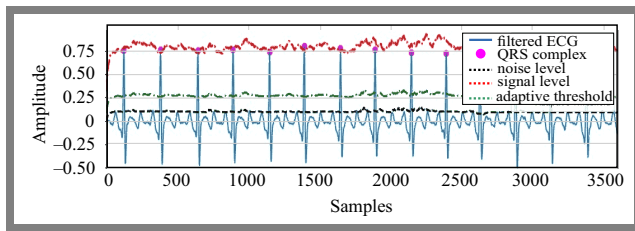


Fig. 7. Decision-making results (record no. 112).

line (black), a signal level line (red), and a dashed green line showing the adaptive threshold used to detect QRS complexes.

#### 4.5. Heart Rate and Classification

Table 2 highlights the analysis of heart rate results after using the MEHC algorithm. The 48 records were taken from the MIT-BIH arrhythmia database. We classified the results as normal or abnormal after comparing data collected after decryption with the heart rate range obtained from the database.

Our findings revealed that once the ECG signal was decrypted using the proposed method, 46 of 48 samples in the MIT-BIH database were classified correctly. For records 210 and 217, the results obtained were classified incorrectly because they contained significant amounts of noise. This suggests that the algorithm is mostly resilient and is capable of effectively classifying almost all samples despite the presence of noise in a couple of instances. Consequently, it underscores the possibility of practical implementation in real-world scenarios.

#### 4.6. Performance Comparison

To test the efficacy of the proposed system, we assessed each ECG sample in the MIT-BIH database using four essential metrics: accuracy, sensitivity, positive predictive value and detection error rate. The average values of these measures were calculated and compared with the findings of previous investigations, as shown in Tab. 3.

The results show competitive performance of the proposed solution when various metrics are compared with other ECG record encryption methods, demonstrating the effectiveness of the MEHC approach. Our method consistently achieves the highest levels of accuracy on various record counts, indicating reliable data preservation. We also obtained high rates of sensitivity, demonstrating the technique's decent efficacy in accurately detecting genuine positive instances of arrhythmia. Moreover, the proposed approach correctly classifies a vast majority of positive instances, as evidenced by the positive predictive values.

Moreover, the method boasts the lowest detection error rate (DER) among all studied approaches, demonstrating its exceptional stability and dependability.

## 5. Conclusions

The primary goal of signal security is to protect against unauthorized access, as well as tampering, disruption, alteration, and destruction of ECG signals. In this study, an altered ver-

Tab. 2. Heart rate results (in bpm) with classification.

Record	Heart rate range	HR after using MEHC	Classification	
100	70–89	75.5003	Normal	✓
101	55–79	62.0751	Normal	✓
102	72–78	72.6757	Normal	✓
103	62–92	69.2197	Normal	✓
104	69–82	71.6787	Normal	✓
105	78–102	86.1342	Normal	✓
106	49–87	65.7969	Normal	✓
107	68–82	70.5157	Normal	✓
108	44–78	64.9329	Normal	✓
109	77–101	83.9742	Normal	✓
111	64–82	70.5489	Normal	✓
112	74–91	84.3729	Normal	✓
113	48–87	59.616	Bradycardia	✓
114	51–82	59.4498	Bradycardia	✓
115	50–84	64.8997	Normal	✓
116	74–86	79.488	Normal	✓
117	48–66	51.0092	Bradycardia	✓
118	54–91	75.7329	Normal	✓
119	61–84	66.0295	Normal	✓
121	55–83	61.8757	Normal	✓
122	67–97	82.2794	Normal	✓
123	41–65	50.3446	Bradycardia	✓
124	47–64	53.4683	Bradycardia	✓
200	69–111	86.6658	Normal	✓
201	31–61	63.6702	Normal	✓
202	49–69	70.7483	Normal	✓
203	63–173	96.8345	Normal	✓
205	80–99	88.1612	Normal	✓
207	57–90	72.3434	Normal	✓
208	91–134	94.8738	Normal	✓
209	82–116	99.8585	Normal	✓
210	63–158	85.1372	Normal	×
212	63–108	91.3182	Normal	✓
213	101–113	107.4683	Tachycardia	✓
214	49–92	74.9354	Normal	✓
215	81–215	111.6222	Tachycardia	✓
217	69–103	73.2074	Normal	×
219	38–75	71.5126	Normal	✓
220	58–74	68.0566	Normal	✓
221	47–110	79.1225	Normal	✓
222	49–84	82.5785	Normal	✓
223	75–94	86.4	Normal	✓
228	54–80	69.6517	Normal	✓
230	63–99	74.9686	Normal	✓
231	49–69	52.2055	Bradycardia	✓
232	24–28	59.3169	Bradycardia	✓
233	98–110	102.0517	Tachycardia	✓
234	84–99	91.3514	Normal	✓

**Tab. 3.** Comparison with other studies.

Method	Software	Record	Same record	Acc [%]	Se [%]	+P [%]	DER [%]
RSA [9]	Matlab R2019a	20	18	90	NR	NR	NR
Our study	Python	20	19	99.11	99.19	99.9	0.89
Gentry FHE [10]	Matlab R2019a	27	25	NR	92.59	90	14.8
Our study	Python	27	26	98.58	98.89	99.45	1.65
FHEAES [11]	Matlab R2018b	48	46	95.8	NR	NR	NR
Our study	Python	48	46	98.48	99.10	99.33	1.52

sion of the EHC technique is used to build a secure ECG signal encryption system.

The effectiveness of the designed cryptosystem was assessed with the use of various signals taken from the MIT-BIH arrhythmia database.

Heart rate was calculated after using MEHC and classified as normal or abnormal, helping to diagnose arrhythmias, such as bradycardia and tachycardia. When compared with results from the original database, the suggested MEHC technique achieves a 98.48% degree of accuracy, a 99.10% sensitivity rate and a 99.33% positive prediction rate.

Furthermore, MEHC is considered a successful homomorphic encryption method, since it correctly decodes the encrypted data. This ensures the algorithm’s suitability and clarity. It also guarantees a high degree of security and privacy in practical applications, such as encryption and decryption of signal data.

## Acknowledgments

This study was supported by the Algerian Ministry of Higher Education and Scientific Research through funding for the PRFU Project.

## References

- [1] V.J. Naveen, K.M. Krishna, and K.R. Rajeswari, “Noise Reduction in ECG Signals for Bio-telemetry”, *International Journal of Electrical and Computer Engineering (IJECE)*, vol. 9, pp. 505–511, 2019 (<https://doi.org/10.11591/ijece.v9i1.pp505-511>).
- [2] M.A. Hashim, Y.W. Hau, and R. Baktheri, “Efficient QRS Complex Detection Algorithm Implementation on SOC-based Embedded System”, *Jurnal Teknologi*, vol. 78, pp. 49–58, 2016 (<https://doi.org/10.11113/jt.v78.9450>).
- [3] K.K. Patro and P.R. Kumar, “Effective Feature Extraction of ECG for Biometric Application”, *Procedia Computer Science*, vol. 155, pp. 296–306, 2017 (<https://doi.org/10.1016/j.procs.2017.09.138>).
- [4] H. Xiong, M. Liang, and J. Liu, “A Real-time QRS Detection Algorithm Based on Energy Segmentation for Exercise Electrocardiogram”, *Circuits, Systems, and Signal Processing*, vol. 40, pp. 4969–4985, 2021 (<https://doi.org/10.1007/s00034-021-01702-z>).
- [5] P. Dhiman *et al.*, “Secure Token-key Implications in an Enterprise Multi-tenancy Environment Using BGV-EHC Hybrid Homomorphic Encryption”, *Electronics*, vol. 11, art. no. 1942, 2022 (<https://doi.org/10.3390/electronics11131942>).
- [6] G. VNKV Subba Rao and G. Uma, “An Efficient Secure Message Transmission in Mobile Ad Hoc Networks Using Enhanced Homomorphic Encryption Scheme”, *Global Journal of Computer Science and Technology*, vol. 13, pp. 21–33, 2013 [Online]. Available: <https://computerresearch.org/index.php/computer/article/view/169>.
- [7] S. Zaineldeen and A. Ate, “Improve the Security of Transfer Data File on the Cloud by Executing Hybrid Encryption Algorithms”, *Indonesian Journal of Electrical Engineering and Computer Science*, vol. 20, pp. 521–527, 2020 (<https://doi.org/10.11591/ijeecs.v20.i1.pp521-527>).
- [8] H.M. Al-Mashhadi and A.A. Khalf, “Hybrid Homomorphic Cryptosystem for Secure Transfer of Color Image on Public Cloud”, *Journal of Theoretical and Applied Information Technology*, vol. 96, pp. 6474–6486, 2018 [Online]. Available: <https://www.jatit.org/volumes/Vol196No19/17Vo196No19.pdf>.
- [9] M.N. Imtiaz and N. Khan, “Pan-Tompkins++: A Robust Approach to Detect R-peaks in ECG Signals”, *2022 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, Las Vegas, USA, 2022 (<https://doi.org/10.1109/BIBM55620.2022.9995552>).
- [10] M. Ngendahimana and W. Shen, “RSA Cryptosystem Speed Security Enhancement (Hybrid and Parallel Domain Approach)”, *Crypto and Information Security*, vol. 2, 2023 (<https://doi.org/10.23977/crypis.2023.020101>).
- [11] M.U. Shaikh, W.A.W. Adnan, and S.A. Ahmad, “Secured Electrocardiograph (ECG) Signal Using Partially Homomorphic Encryption Technique-RSA Algorithm”, *Pertanika Journal of Science and Technology*, vol. 28, pp. 231–242, 2020 (<https://doi.org/10.47836/pjst.28.s2.18>).
- [12] M.U. Shaikh, W.A.W. Adnan, and S.A. Ahmad, “Sensitivity and Positive Prediction of Secured Electrocardiograph (ECG) Transmission using Fully Homomorphic Encryption Technique (FHE)”, *2020 IEEE-EMBS Conference on Biomedical Engineering and Sciences (IECBES)*, Langkawi Island, Malaysia, 2021 (<https://doi.org/10.1109/IECBES48179.2021.9398792>).
- [13] A.A. Ahmed, M.M. Madboly, and S.K. Guirguis, “Securing Data Transmission and Privacy Preserving Using Fully Homomorphic Encryption”, *International Journal of Intelligent Engineering and Systems*, vol. 16, pp. 277–289, 2023 (<https://doi.org/10.22266/ijies2023.0228.25>).
- [14] P. Vithya KP, “Secured ECG Distribution Using Compression and RSA Algorithm for Telemedicine Application”, *International Journal of Recent Technology and Engineering (IJRTE)*, vol. 3, pp. 2277–3878, 2014 [Online]. Available: <https://api.semanticscholar.org/CorpusID:212557890>.
- [15] O. Kocabas and T. Soyata, “Medical Data Analytics in the Cloud Using Homomorphic Encryption”, *E-Health and Telemedicine*, pp. 751–768, 2016 (<https://doi.org/10.4018/978-1-4666-8756-1.ch038>).
- [16] G.B. Moody and R.G. Mark, “The Impact of the MIT-BIH Arrhythmia Database”, *IEEE Engineering in Medicine and Biology Magazine*, vol. 20, pp. 45–50, 2001 (<https://doi.org/10.1109/51.932724>).
- [17] W.L. Caldas, J.P.V. Madeiro, C.L.C. Mattos, and J.P.P. Gomes, “A New Methodology for Classifying QRS Morphology in ECG Signals”, *International Joint Conference on Neural Networks*, Glasgow, UK, 2020 (<https://doi.org/10.1109/IJCNN48605.2020.9206707>).

- [18] D. Menard, D. Chillet, and O. Sentieys, "Floating-to-fixed-point Conversion for Digital Signal Processors", *EURASIP Journal on Advances on Signal Processing*, vol. 2006, art. no. 096421, 2006 (<https://doi.org/10.1155/ASP/2006/96421>).
- [19] K. Sayood, *Introduction to Data Compression*, 5th ed., Elsevier, 765 p., 2018 (<https://doi.org/10.1016/C2015-0-06248-7>).
- [20] S. Banerjee and G.K. Singh, "A New Real-time Lossless Data Compression Algorithm for ECG and PPG Signals", *Biomedical Signal Processing and Control*, vol. 79, art. no. 104127, 2023 (<https://doi.org/10.1016/j.bspc.2022.104127>).
- [21] B. Yogapriya *et al.*, "Accelerating Linear Congruential Generators with Carbon Nanotube Field-effect Transistors", *IOP Conference Series: Materials Science and Engineering*, vol. 1316, art. no. 012005, 2024 (<https://doi.org/10.1088/1757-899X/1316/1/012005>).
- [22] S. Nisha and M. Farik, "RSA Public Key Cryptography Algorithm – A Review", *International Journal of Scientific & Technology Research*, vol. 6, pp. 187–191, 2017.
- [23] A.A. Ahmed, M.M. Madboly, and S.K. Guirguis, "Securing Signal Encryption Based on Reduced Round Homomorphic AES", *International Journal of Intelligent Engineering and Systems*, vol. 16, pp. 440–454, 2023 (<https://doi.org/10.22266/ijies2023.0630.35>).
- [24] V. Mondelo *et al.*, "Detection of Heart Beat Positions in ECG Recordings: A Lead-dependent Algorithm", *Journal of Information Systems Engineering & Management*, vol. 2, pp. 1–8, 2017 (<https://www.jisem-journal.com/download/64NV7FMF.pdf>).
- [25] H. De Melo Ribeiro *et al.*, "ECG-based Real-time Arrhythmia Monitoring Using Quantized Deep Neural Networks: A Feasibility Study", *Computers in Biology and Medicine*, vol. 143, art. no. 105249, 2022 (<https://doi.org/10.1016/j.compbiomed.2022.105249>).
- [26] B. Adithya and G. Santhi, "A DNA Sequencing Medical Image Encryption System (DMIES) Using Chaos Map and Knight's Travel

Map", *International Journal of Reliable and Quality E-Healthcare*, vol. 11, pp. 1–22, 2022 (<https://doi.org/10.4018/IJRQEH.308803>).

#### Fatma Zohra Besmi, Ph.D. Student

LIST Laboratory, Department of Electrical Systems Engineering

 <https://orcid.org/0009-0007-8216-0831>


E-mail: [f.besmi@univ-boumerdes.dz](mailto:f.besmi@univ-boumerdes.dz)

University of Boumerdes, Boumerdes, Algeria

<https://www.univ-boumerdes.dz>

#### Samia Belkacem, Ph.D., Associate Professor

Department of Electrical Systems Engineering

 <https://orcid.org/0000-0003-0912-3392>

E-mail: [s.belkacem@univ-boumerdes.dz](mailto:s.belkacem@univ-boumerdes.dz)

University of Boumerdes, Boumerdes, Algeria

<https://www.univ-boumerdes.dz>

#### Noureddine Messaoudi, Professor

LIST Laboratory, Department of Electrical Systems Engineering

 <https://orcid.org/0000-0002-7228-5784>

E-mail: [n.messaoudi@univ-boumerdes.dz](mailto:n.messaoudi@univ-boumerdes.dz)

University of Boumerdes, Boumerdes, Algeria

<https://www.univ-boumerdes.dz>

# ILP Optimized LSTM-based Autoscaling and Scheduling of Containers in Edge-cloud Environment

Shivan Singh, Narayan D.G., Sadaf Mujawar, G.S. Hanchinamani, and P.S. Hiremath

*KLE Technological University, Hubballi, Karnataka, India*

<https://doi.org/10.26636/jtit.2025.2.2088>

**Abstract** — Edge computing is a decentralized computing paradigm that brings computation and data storage closer to data sources, enabling faster processing and reduced latency. This approach is critical for real-time applications, but it introduces significant challenges in managing resources efficiently in edge-cloud environments. Issues such as increased response times, inefficient autoscaling, and suboptimal task scheduling arise due to the dynamic and resource-constrained nature of edge nodes. Kubernetes, a widely used container orchestration platform, provides basic autoscaling and scheduling mechanisms, but its default configurations often fail to meet the stringent performance requirements of edge environments, especially in lightweight implementations like KubeEdge. This work presents an ILP-optimized, LSTM-based approach for autoscaling and scheduling in edge-cloud environments. The LSTM model forecasts resource demands using both real-time and historical data, enabling proactive resource allocation, while the integer linear programming (ILP) framework optimally assigns workloads and scales containers to meet predicted demands. By jointly addressing auto-scaling and scheduling challenges, the proposed method improves response time and resource utilization. The experimental setup is built on a KubeEdge testbed deployed across 11 nodes (1 cloud node and 10 edge nodes). Experimental results show that the ILP-enhanced framework achieves a 12.34% reduction in response time and a 7.85% increase in throughput compared to the LSTM-only approach.

**Keywords** — *autoscaling, edge computing, ILP optimization, Kubernetes, LSTM, resource efficiency, scheduling, throughput*

## 1. Introduction

Edge computing represents an approach to data processing that enables computation and storage closer to the data source than using centralized cloud servers. This decentralized model improves real-time data analysis, reduces latency, and improves resource utilization, making it suitable for applications like autonomous systems, smart cities, and industrial automation.

The unique requirements of edge computing environments, including low latency responses and efficient resource utilization, create significant challenges in workload management and resource optimization.

KubeEdge is an open-source framework designed to extend Kubernetes functionality to edge computing environments,

enabling efficient management of containerized applications across distributed edge nodes. Bridges the gap between cloud infrastructure and edge devices, facilitating seamless deployment and orchestration of workloads in resource-constrained and geographically dispersed locations.

The architecture of KubeEdge includes components optimized for edge scenarios, such as the CloudCore module, which manages edge node control, configuration, and communication with the Kubernetes API server at the cloud level, and the EdgeCore module, which handles application deployment, resource monitoring, and local decision-making at the edge, reducing dependency on continuous cloud connectivity.

Additionally, KubeEdge incorporates an edge message bus for real-time communication between devices and applications, requiring low latency and high responsiveness. By enhancing Kubernetes with edge-specific capabilities, KubeEdge provides a robust platform for deploying scalable and reliable applications in distributed environments.

Autoscaling completes scheduling by dynamically adjusting the number of container replicas to match workload demands. The Kubernetes HorizontalPod Autoscaler (HPA) primarily relies on metrics such as CPU and memory usage to scale resources. For edge computing, network traffic information can play an important role in reducing response time [1]. Although effective in static or predictable workloads, this approach struggles in dynamic edge environments characterized by unpredictable workload patterns.

Proactive scaling mechanisms, using predictive models such as long-short-term memory (LSTM) networks, offer a promising solution [2]. By analyzing historical and real-time metrics, LSTM models can predict future resource demands, enabling preemptive scaling decisions. Reduce resource underutilization and overprovisioning and also ensure timely responses to workload fluctuations.

Scheduling in edge computing plays an important role in efficiently assigning tasks to nodes while minimizing latency and balancing workloads across distributed resources. Unlike traditional cloud environments, where computational resources are large, edge nodes operate under strict resource limitations. Effective scheduling requires consideration of factors such as network conditions, task relationship, and node heterogeneity. For example, tasks that require real-time processing must

be assigned to nodes closer to the data source to ensure low latency, while less critical tasks can be offloaded to distant nodes or cloud servers [3].

Recent advances have explored machine learning and integer linear programming (ILP) to enhance scheduling efficiency, but their integration with existing Kubernetes architectures remains a challenge. Integrating autoscaling with scheduling in Kubernetes enhances resource management by dynamically adjusting workloads to real-time demands. Autoscaling mechanisms such as HPA scale pod replicas based on resource utilization, while intelligent scheduling ensures optimal workload placement across clusters.

A survey on Kubernetes scheduling algorithms emphasizes the importance of autoscaling-enabled scheduling, advocating for algorithms that adapt to dynamic workloads by integrating autoscaling into the scheduling process [4]. This combined approach improves resource allocation and reduces latency under fluctuating workloads.

This work integrates AI-based autoscaling and scheduling mechanisms customized for Kubernetes-based edge-cloud environments. We propose an ILP-optimized LSTM-based approach that addresses the limitations of existing solutions. Using CPU and memory usage as well as RTT metrics, the LSTM model predicts workloads, enabling proactive scaling decisions. In addition, an ILP-based scheduling algorithm assigns tasks to nodes based on real-time and predicted resource availability, optimizing response time and resource utilization. By integrating these processes, the proposed framework ensures efficient operation in dynamic and resource-constrained environments.

The contributions of this work are as follows:

- an LSTM-based prediction model is proposed to forecast resource utilization such as CPU, memory, and RTT, enabling accurate workload predictions,
- a combined autoscaling and scheduling framework that integrates LSTM-based predictions with ILP-based optimization is proposed,
- the proposed approach is evaluated in a KubeEdge testbed environment to determine improvements in resource utilization, response time, and workload prediction accuracy and then compared to traditional methods.

The rest of the article is as follows. Section 2 reviews related research and background study. Sections 3 and 4 outline the mathematical model, the proposed model, and the algorithms. The results are presented in Section 5, and conclusions are given in Section 6.

## 2. Background Study

The KubeEdge architecture, as shown in Fig. 1, is designed to seamlessly integrate cloud and edge computing environments. It consists of two main layers: the cloud layer and the edge layer. The cloud layer hosts the master node, which contains critical components like the scheduler, metric server, deployments, and autoscaler. These components ensure that tasks are managed efficiently, that resources are optimally

allocated, and that the overall system remains scalable and reliable.

The edge layer, on the other hand, includes multiple worker nodes that are connected through EdgeHub. These worker nodes host applications and devices, providing compute power at the edge of the network closer to end users.

The CloudCore, which operates in the cloud, is a key component of the KubeEdge architecture. It consists of various modules that handle critical communication, synchronization, and management tasks. One of these modules is CloudHub, which acts as the primary gateway for communication between the cloud and the edge nodes. It is responsible for maintaining secure web socket connections and routing messages efficiently. Another important module is the edge controller, which oversees the management of edge nodes and synchronizes pod metadata between the edge and the Kubernetes API server. Additionally, the device controller plays a crucial role in ensuring that device metadata is accurately synchronized between the cloud and edge environments.

The KubeEdge architecture incorporates sophisticated mechanisms to handle varying workloads and dynamic resource requirements. One of the standout features is its autoscaling capability, which ensures that the system can adapt to changes in computational demands. This is particularly important for edge computing scenarios where workloads can fluctuate based on real-time events or user interactions. By dynamically adjusting resources, KubeEdge maintains optimal performance without overprovisioning or underutilizing resources.

Scheduling in KubeEdge is another critical aspect of its architecture. The scheduling framework operates both at the cloud and the edge levels, ensuring that tasks are effectively distributed across available resources. At the cloud level, the master scheduler is responsible for global resource allocation. It evaluates various factors such as resource availability, network conditions, and task priorities to make informed scheduling decisions.

At the edge level, the EdgeCore components handle local scheduling. These components are designed to optimize resource utilization while considering factors such as net-

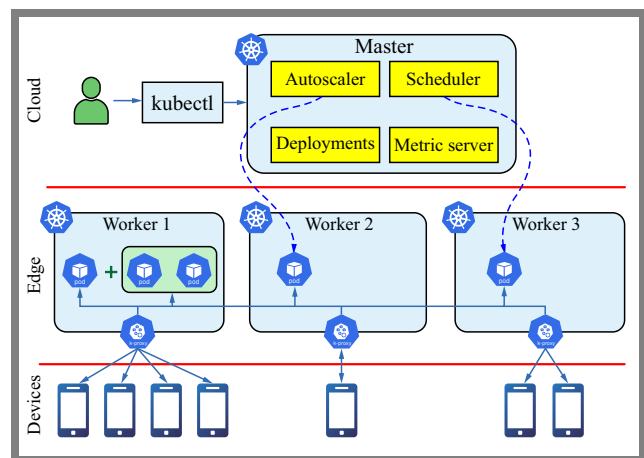


Fig. 1. KubeEdge architecture.

work latency and data locality. The distributed nature of this scheduling framework allows KubeEdge to achieve a balance between centralized control and autonomous operations at the edge.

One of the most important aspects of KubeEdge is its ability to maintain operational consistency even in the face of discontinuous connectivity. Edge environments often operate in challenging conditions where stable network connections cannot be guaranteed. KubeEdge addresses this challenge by ensuring that edge nodes can continue to function autonomously even when disconnected from the cloud. This resilience makes it particularly suitable for scenarios like industrial automation, smart cities, and remote monitoring, where edge devices must continue to operate independently.

### 2.1. Related Work

In the realm of autoscaling, a hybrid proactive autoscaler optimized for edge computing scenarios has been proposed in [5]. Utilizing a bidirectional long- and short-term memory (Bi-LSTM) based load prediction model, this autoscaler predicts future workloads and performs scaling operations preemptively. Furthermore, an overload compensation algorithm is implemented to ensure quality of service (QoS) degradation due to underprediction, and a hybrid scaling method is applied to simultaneously adjust the number of pods and their resource quota without restarting [5].

In [6], the authors introduce an integrated framework that combines autoscaling and scheduling for edge-cloud environments. This framework dynamically adjusts resources and schedules tasks based on real-time workload analysis, enhancing system performance and resource utilization. Similarly, article [7] presents a proactive autoscaling approach for edge computing systems managed with Kubernetes. By forecasting workloads using multiple user-defined metrics, the proposed proactive pod autoscaler (PPA) scales applications accordingly, outperforming the default autoscaler in both resource utilization efficiency and application performance. These studies highlight the importance of integrating predictive models with resource management strategies to effectively handle dynamic workloads in edge-cloud environments.

Node-aware autoscaling has been extensively explored to address dynamic workloads in edge computing environments. Paper [8] introduces a node-based horizontal pod autoscaler (NHPA) designed for KubeEdge environments. Focuses on the use of resource at the individual node level, allowing dynamic adjustment of pod numbers independently for each node. This ensures optimized pod allocation and seamless scaling, particularly in scenarios where traffic volume fluctuates over time and location, and communication links between edge nodes may be unstable.

Integrating machine learning models with autoscaling has been investigated in [9], while the authors present FedAvg-BiGRU, a proactive autoscaling method in edge computing that combines federated averaging (FedAvg) and multistep prediction using a bidirectional gated recurrent unit (BiGRU). This approach reduces network traffic by exchanging model

updates instead of raw data, thereby enhancing resource allocation efficiency. Similarly, [10] proposes a hybrid proactive autoscaler that combines horizontal and vertical scaling based on future workload predictions. This method aims to improve QoS and resource utilization efficiency in Kubernetes clusters.

Predictive workload modeling is crucial to improve autoscaling and scheduling in dynamic environments. In [11], the authors address container job scheduling as a multiobjective optimization problem, proposing a linear programming model to address this issue. They highlight the limitations of traditional approaches in capturing the non-linearities associated with resource usage patterns, suggesting that deep neural networks could offer a more effective solution.

Dynamic resource allocation frameworks have been developed to enhance Kubernetes cluster management. For example, a study [12] proposes a dynamic task offloading framework in KubeEdge-based edge computing environments, utilizing machine learning to optimize resource allocation and ensure data privacy. This approach addresses the challenges of resource limitations and privacy concerns in edge computing scenarios. Similarly, work [13] introduces a proactive hybrid autoscaler designed for edge applications in Kubernetes. By employing a Bi-LSTM based load prediction model, this autoscaler anticipates future workloads, enabling preemptive scaling actions that improve resource utilization and maintain QoS. These frameworks demonstrate the effectiveness of integrating predictive models and machine learning techniques for dynamic resource management in Kubernetes clusters.

Efficient task scheduling is a critical challenge in edge computing environments. In [3], the authors propose a network-based container scheduling approach that considers the various edge node network performance, such as geographical location and network topology, to optimize resource allocation and application performance. Furthermore, [14] investigates the performance of KubeEdge in terms of computational resource distribution and latency between edge nodes. The study reveals that forwarding traffic between edge nodes leads to a degraded throughput and an increased service delay in an edge computing environment. To mitigate this problem, the authors propose a local scheduling scheme that processes user traffic locally at each edge node, enhancing the performance of edge devices.

In edge computing, efficient resource scheduling is crucial to manage the limited computational resources and dynamic workloads characteristic of edge environments. A comprehensive taxonomy of resource scheduling techniques is presented in [15], categorizing approaches based on application scenarios, computational platforms, algorithm paradigms, and optimization objectives. This taxonomy addresses challenges such as heterogeneity, workload dynamics, and the need for real-time processing, providing a structured framework for developing effective scheduling strategies. In addition, a multi-objective optimization algorithm is introduced in [16], with the aim of minimizing latency and maximizing resource utilization in edge systems. This algorithm considers factors such as task allocation, resource availability, and

**Tab. 1.** Summary of related work and research gaps.

Ref.	Focus area	Methodology	Research gap
[1]	Traffic-aware autoscaling	Incorporates traffic patterns into HPA for Kubernetes	Lacks scalability for multi-cluster deployments in edge computing
[2]	LSTM-based autoscaling	Analyzes real-time and historical metrics to forecast resource demands, enabling proactive scaling	Limited integration with scheduling frameworks for task placement
[5]	Proactive autoscaling	Predicts workload demands using Bi-LSTM models for preemptive scaling	Limited focus on heterogeneous resource capacities in edge environments
[9]	Federated autoscaling	Leverages FL models to predict workload variations across distributed edge nodes	High communication overhead and synchronization challenges
[13]	ILP-based scheduling	Optimizes container task placement using ILP models	Limited adaptability to dynamic workloads in real-time systems
[14]	Real-time scheduling	Allocates tasks in KubeEdge environments, focusing on real-time edge task execution	Does not incorporate latency-awareness in multi-node edge deployments
[15]	Resource-aware scheduling	Proposes a taxonomy for scheduling algorithms with a focus on adaptive mechanisms	Absence of integration with predictive autoscaling mechanisms
[16]	Multi-objective scheduling	Employs multi-objective optimization to balance latency and resource usage in task scheduling	Requires enhanced scalability for large-scale edge systems

network conditions to ensure efficient processing and timely responses to user demands. By integrating these heuristic and taxonomy-based approaches, edge computing systems can achieve improved performance, adaptability, and resource management.

Energy efficiency is a critical concern in edge computing scheduling algorithms. In [18] a workload scheduling approach based on deep reinforcement learning (DRL) has been proposed to balance workloads, reduce service time, and decrease task failure rates in edge environments. This method utilizes deep Q-network (DQN) algorithms to address the complexities of workload scheduling, with the aim of improving virtual machine utilization and overall system performance. Article [19] introduces energy-efficient scheduling algorithms to minimize computational overhead while maintaining performance in edge environments.

In the realm of edge computing, DRL has been used to improve task scheduling efficiency. In [20] the DRL-based task scheduling algorithm has been introduced to intelligently manage tasks in edge computing settings, focusing on reducing service delay and traffic load. This approach uses reinforcement learning to optimize task assignments, thereby enhancing the efficiency of edge computing systems.

ILP-based approaches have also been used to optimize resource allocation. The authors of [21] use ILP models to optimize the placement of service function chains (SFC) in edge cloud environments, integrating workload predictions from LSTM models to improve efficiency. Similarly, paper [22] proposes the combined predictive autoscaler (COPA), which combines horizontal and vertical scaling to optimize resource usage in Kubernetes clusters.

Research is summarized in the Tab. 1 highlights significant advances in the fields of autoscaling and scheduling techniques.

### 3. ILP Formulation for Joint Autoscaling and Scheduling

To optimize resource allocation, task placement, and replica scaling simultaneously in the Kubernetes-based edge-cloud environment using a linear framework, we propose an ILP model. This model uses predictions from the LSTM model to make proactive decisions, with the aim of minimizing a composite cost function reflecting task priorities, network latency, and a linear cost for scaling, while adhering to node resource constraints. The decision variables are the following.

- $x_{p,n}$ : binary decision variable, where  $x_{p,n} = 1$  if the task  $p$  is assigned to node  $n$  and  $x_{p,n} = 0$  otherwise. This represents the *scheduling* decision.
- $R$ : non-negative integer decision variable representing the total target number of container replicas to be maintained across the cluster, as determined by the integrated optimization. This represents the *autoscaling* decision.

Parameters used in the proposed model:

- $c_p$ : cost or inverse priority associated with placing task  $p$ . A higher  $c_p$  might represent a lower priority task or a higher intrinsic cost of running it.
- $RTT_n$ : predicted round-trip time parameter for node  $n$ . In this work,  $RTT_n$  is defined as the predicted latency between edge node  $n$  and a designated central point within

the cluster, the Kubernetes API server. This serves as a proxy for the general responsiveness and accessibility of node  $n$  from a control or coordination perspective. This value is forecasted using a dedicated LSTM model.

The inputs to this LSTM model to predict  $\text{RTT}_n$  for a specific node  $n$  include its own historical  $\text{RTT}_n$  values (to the central point) from previous time intervals, the current CPU and memory utilization of node  $n$ , and the current network interface traffic statistics for node  $n$  (e.g., bytes in/out, packets in/out).

These input metrics are collected through Prometheus. The LSTM is trained offline on historical data to learn patterns and predict  $\text{RTT}_n$  for the subsequent operational interval. This predicted  $\text{RTT}_n$  is then fed as a constant parameter into each instance of the ILP optimization problem.

- $\gamma$ : non-negative penalty factor for network latency. Controls the importance of minimizing latency during task placement.
- $\beta$ : non-negative *linear* scaling cost factor. This weight represents the cost associated with the deployment and maintenance of each replica. Penalize solutions linearly on the basis of the total number of replicas  $R$ .
- $\text{CPU}_p$ ,  $\text{memory}_p$ : resource requirements (CPU cores, memory units) of a single instance/replica of task  $p$ .
- $\text{CPU}_n$ ,  $\text{memory}_n$ : resource capacities (available CPU cores, memory units) of edge node  $n$ .
- $N$ : total number of edge nodes.
- $P$ : total number of tasks to be scheduled.

The following constraints ensure valid scheduling and resource allocation. Each task must be assigned to exactly one node. To achieve this, the following formula is applied:

$$\sum_{n=1}^N x_{p,n} = 1, \quad \forall p \in \{1, \dots, P\}. \quad (1)$$

The total CPU usage of tasks assigned to a node must not exceed the CPU capacity. To prevent over-commitment of CPU resources on any node, we use:

$$\sum_{p=1}^P x_{p,n} \cdot \text{CPU}_p \leq \text{CPU}_n, \quad \forall n \in \{1, \dots, N\}. \quad (2)$$

Total memory usage of tasks assigned to a node must not exceed the memory capacity:

$$\sum_{p=1}^P x_{p,n} \cdot \text{Memory}_p \leq \text{Memory}_n, \quad \forall n \in \{1, \dots, N\}. \quad (3)$$

The task-to-node assignment must be binary:

$$x_{p,n} \in \{0, 1\}, \quad \forall p \in \{1, \dots, P\}, \quad n \in \{1, \dots, N\}. \quad (4)$$

The total number of replicas must be non-negative:

$$R \geq 0, \quad R \in \mathbb{Z}. \quad (5)$$

The goal is to minimize the combined linear cost of task placement, network latency, and replica scaling:

$$\min \sum_{p=1}^P \sum_{n=1}^N (c_p \cdot x_{p,n} + \gamma \cdot \text{RTT}_n \cdot x_{p,n}) + \beta \cdot R. \quad (6)$$

The objective function (6) aims to find the optimal balance between scheduling efficiency and resource scaling costs using a purely linear formulation:

- The scheduling cost term ( $c_p \cdot x_{p,n}$ ) accumulates the intrinsic cost  $c_p$  of placing the task  $p$  on node  $n$ . Favors placing low-cost (high-priority) tasks.
- The latency penalty term ( $\gamma \cdot \text{RTT}_n \cdot x_{p,n}$ ) adds a penalty proportional to the predicted latency ( $\text{RTT}_n$ ) of node  $n$  where the task  $p$  is placed. The factor  $\gamma$  weights the importance of latency. Minimizing this drives tasks towards low-latency nodes.
- The linear scaling cost term ( $\beta \cdot R$ ) adds a cost that increases *linearly* with the total number of replicas  $R$  deployed. The factor  $\beta$  represents the cost per replica. This encourages resource efficiency by penalizing unnecessarily high replica counts.

The ILP solver finds the optimal integer values for  $x_{p,n}$  (task placement) and  $R$  (total replicas) that minimize this linear objective function while satisfying all constraints (1)–(5).

The integration with LSTM remains the same as for ILP. The LSTM model provides predictive inputs:

- The predicted round-trip time  $\text{RTT}_n$  is used directly in the latency penalty term.
- Other LSTM predictions can inform the setting of parameters like  $c_p$ ,  $\text{CPU}_p$ , or  $\text{memory}_p$  before solving the ILP.

By integrating autoscaling and scheduling equations into the ILP framework, the proposed model addresses the challenges of dynamic workload management in edge-cloud environments. The Gurobi solver is used to solve the ILP, ensuring efficient optimization of resources while maintaining low response times.

## 4. Proposed Methodology

This section details the methodology of the proposed system, covering system workflow, dataset preparation, scheduling algorithm, auto-scaling mechanism, and mathematical modeling using ILP for combined scheduling and autoscaling.

### 4.1. System Model

The system model, depicted in Fig. 1, integrates auto-scaling, scheduling, and predictive modeling to achieve efficient resource management in Kubernetes-based edge-cloud environments. At the core of the design depicted in Fig. 2 is a seamless workflow that begins with collecting real-time resource metrics, such as CPU usage, memory consumption, and RTT from all nodes within the cluster.

Metric server as the monitoring agent, providing an uninterrupted stream of data essential for decision-making processes.

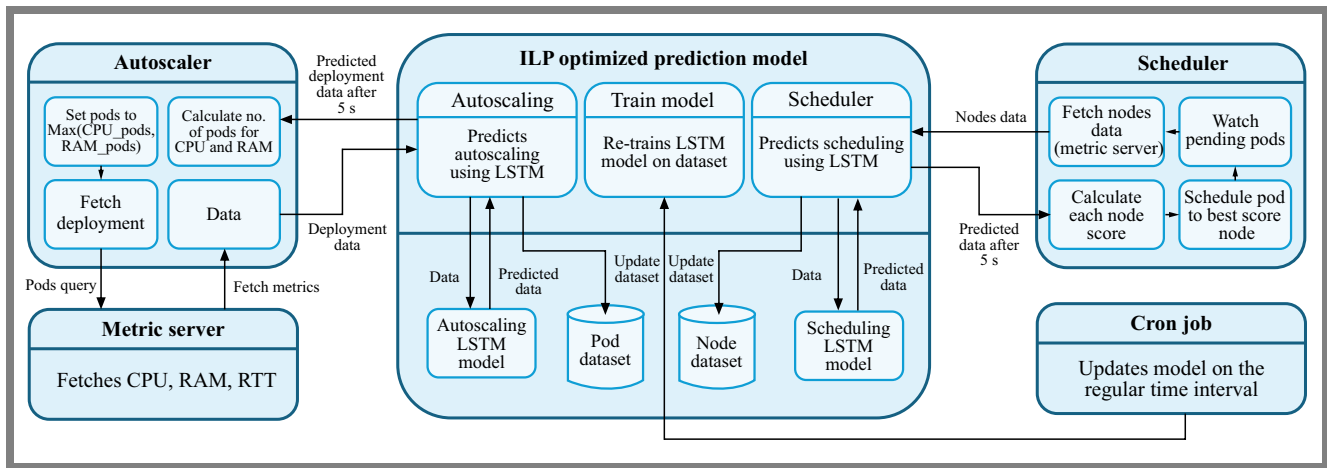


Fig. 2. Detailed system model.

The autoscaler processes these metrics and prepares them to forecast future workloads for 5 s using an LSTM-based prediction model. This predictive model utilizes historical and real-time data to estimate future resource requirements, including CPU and memory usage, network delays, and application request rates. Predicted values are stored alongside the actual metrics, creating a continuously updated dataset to retrain the LSTM model. This approach ensures that the accuracy of the prediction evolves with changes in workload patterns, making the system highly adaptive.

When resource demands exceed predefined thresholds, the autoscaler proactively scales resources by either increasing replicas or reallocating workloads to mitigate performance bottlenecks. Similarly, the scheduler evaluates all nodes within the cluster to determine the most suitable deployment location for new or pending pods. This evaluation is based on a scoring mechanism that incorporates predicted CPU and memory usage, task priority, pod affinity, and network parameters such as RTT. The node with the highest score is selected for pod deployment, ensuring efficient and balanced resource allocation.

To further enhance system performance, the model is designed to retrain the LSTM predictor at regular intervals using a cron job. The updated data set collected during operations allows the retraining process to adapt to evolving workload behaviors, ensuring the system remains capable of handling unpredictable and dynamic resource demands.

4.2. Dataset

The data set used in this investigation was created using Prometheus, a powerful open-source monitoring tool. Prometheus collects real-time metrics related to resource usage from Kubernetes-based edge-cloud environments. The dataset consists of two main components: pod-level data for autoscaling and node-level data for scheduling. Each data set contains features that help predict future resource utilization and optimize resource allocation strategies.

The pod-level dataset focuses on metrics crucial to making auto-scaling decisions. The data set includes information

Tab. 2. Description of the dataset at the pod level.

Feature	Description
timestamp	The timestamp indicating the time of the measurement
cpu	CPU usage of the pod
memory	Memory usage of the pod
rtt	Round-trip time or latency for the pod
next_rtt	Round-trip time or latency 5 s after the current time
next_memory	Memory usage 5 s after the current time
next_cpu	CPU usage 5 s after the current time

Tab. 3. Description of the node-level dataset.

Feature	Description
timestamp	The timestamp indicating the time of the measurement
nodename	The name of the node
cpu	CPU usage of the node
memory	Memory usage of the node
rtt	Round-trip time or latency for the node
next_rtt	Round-trip time or latency 5 s after the current time
next_memory	Memory usage 5 s after the current time
next_cpu	CPU usage 5 s after the current time

on CPU usage, memory utilization, RTT, and future predictions. The data is continuously monitored and collected by Prometheus. A detailed description of the pod-level dataset features is provided in the Tab. 2.

The node-level data set is designed to help schedule pods to the most optimal nodes. It includes both current and predicted metrics for nodes such as CPU usage, memory usage, and RTT. This data set is essential to implement an efficient

scheduling algorithm. Table 3 provides a detailed description of the features of the node-level dataset.

### 4.3. The Algorithm Description

The proposed methodology, shown in Algorithm 1, adopts a cyclical, multiphase strategy to integrate predictive insights with optimization for resource management in Kubernetes-based edge environments. This structure aims at both proactive adjustments and refined decision-making. The process involves the following phases:

- 1) Proactive autoscaling (lines 11–23). This phase leverages LSTM predictions (PredictMetricValues), derived from current Prometheus metrics, to forecast future resource demands. Based on these predictions versus operational thresholds, a preliminary target replica count (NoofReplicasHeuristic) is calculated heuristically (lines 20–21). An initial scaling action (ScaleDeployment) might optionally be performed at this stage (line 22) if the heuristic count differs significantly from the current state, allowing a rapid response to predicted load shifts.
- 2) Heuristic scheduling (lines 24–29). Following autoscaling considerations, this phase performs a quick heuristic placement for pending tasks. Calculate a score ( $Score_n$ ) for each node using the function  $f$ , based on current or predicted resources and latency  $RTT_n$  (line 26). Each task is then preliminarily assigned to the node identified as having the best score (BestNode) (lines 28–29).
- 3) ILP optimization and refinement (lines 30–39): The final phase employs ILP for comprehensive optimization. An ILP problem is formulated using the revised model, i.e. constraints (1)–(4) and objective (6) (lines 32–33). The inputs include node states, task requirements, predicted  $RTT_n$  by LSTM, and tuning parameters  $\gamma, \beta, c_p$ . The ILP is solved using Gurobi (line 34) to determine both the final optimal task placement  $x_{p,n}$  and the final optimal total replica count  $R_{optimal}$ . This ILP solution refines the decisions of the previous phases. The resulting  $R_{optimal}$  dictates the definitive scaling action (lines 36–37), and the optimized  $x_{p,n}$  determines the final task deployment (line 39).

This multiphase design allows the system to potentially react quickly using predictive heuristics (phases 1 and 2) while leveraging the comprehensive optimization power of ILP (phase 3) for refinement and determining the definitive scaling and placement actions. Practical considerations such as prediction accuracy, interphase delays, and ILP solve time remain relevant for real-world performance.

### 4.4. Cron Job for LSTM Retraining

To maintain prediction accuracy, the LSTM model is periodically re-trained using updated datasets. A cron job is implemented to automate this process. The steps involved are as follows:

- Fetch updated metrics (CPU, memory, RTT) from Prometheus,

**Tab. 4.** Hardware and software configuration.

Component	Specification
Processor	Intel core i5 7th Gen
RAM	128 GB DDR4
Storage	30 GB SSD per node
Operating system	Ubuntu 20.04 LTS
Kubernetes version	1.29
KubeEdge version	1.19
Nodes	11: 1 cloud node, 10 edge nodes
RAM per node	4 GB
Virtualization	VMware workstation
Orchestration tool	Kubectrl
Programming framework	TensorFlow, Python
Deployed application	Web App using Ngnix

**Tab. 5.** Configuration of the LSTM model.

Parameter	Value
Number of LSTM layers	16
Number of epochs	50
Batch size	64
Optimizer	Adam
Loss function	Categorical cross entropy

- Update the training dataset with the latest metrics,
- Retrain the LSTM model with the updated dataset to improve prediction accuracy,
- Deploy the updated model into the system for future predictions.

## 5. Results and Discussion

The data set was generated in a multinode KubeEdge setup, consisting of one cloud node and ten edge nodes, and metrics were collected at both node and pod levels using Prometheus and bash scripts over a 10-hour period, capturing diverse workloads and resource utilization patterns. Table 4 summarizes the configuration used, while the LSTM model configuration is summarized in Table 5. Workloads to evaluate autoscaling mechanisms were generated using the Apache Benchmark (ab) tool, which aims at the deployed web server application.

For the evaluations focusing on fixed sustained load for response time, CPU/memory utilization), a total of 100 000 HTTPS GET requests ab Apache tool with concurrency of 100. This configuration creates a *closed-loop workload model*. In this model, ab attempts to maintain 100 concurrent active connections to the server. As soon as a request receives a response, ab immediately issues a new request on

---

**Algorithm 1** ILP optimized LSTM-based autoscaling and scheduling.
 

---

```

1: Input:
2:   Historical metrics  $\{CPU_t, MEM_t, RTT_t\}_{t=1}^T$  from metric server
3:   Current replica count  $R_{current}$  from Kubernetes
4:   Threshold values for CPU, memory, and RTT:  $Threshold_{CPU}, Threshold_{memory}, Threshold_{RTT}$ 
5:   Resource capacities limit of nodes:  $CPU_n, Memory_n,$  and  $RTT_n$ 
6:   Tasks and their resource requirements:  $\{CPU_p, MEM_p\}_{p=1}^P$ 
7: Output:
8:   Final optimized replica count  $R_{optimal}$  deployed
9:   Task-to-node mapping  $x_{p,n}$  used for deployment
10: while true do
11:   Phase 1. Autoscaling
12:   DesiredReplicaHeuristic  $\leftarrow 1$ 
13:   FetchedMetrics  $\leftarrow$  FetchMetricsFromPrometheus()
14:   CurrentReplicas  $\leftarrow$  GetCurrentReplicas(Deployment)
15:   for metric in FetchedMetrics do
16:     PredictedMetricValue  $\leftarrow$  PredictMetricValues(metric)
17:     DesiredMetricValue  $\leftarrow$  GetDesiredMetricValue(metric)
18:     MetricDesiredReplicas  $\leftarrow \lceil CurrentReplicas \cdot (PredictedMetricValue / DesiredMetricValue) \rceil$ 
19:     DesiredReplicaHeuristic  $\leftarrow \max(DesiredReplicaHeuristic, MetricDesiredReplicas)$ 
20:   end for
21:   NoofReplicasHeuristic  $\leftarrow$  DesiredReplicaHeuristic
22:   if NoofReplicasHeuristic  $\neq$  CurrentReplicas then
23:     end if
24:   Phase 2. Scheduling
25:   For each task  $p$  calculate node scores using resource metrics:
26:      $Score_n \leftarrow f(CPU_n, Memory_n, RTT_n)$ 
27:   Identify the most suitable node for each task:
28:      $BestNode \leftarrow \arg \max_n (Score_n)$ 
29:   Schedule tasks to nodes based on the highest scores
30:   Phase 3. ILP Optimization
31:   Define decision variables:  $x_{p,n}$  (binary),  $R_{optimal}$  (integer)
32:   Formulate constraints using Eqs. (1)–(4)
33:   Reference objective function from Eq. (6)
34:   Solve ILP using Gurobi solver to determine optimal  $x_{p,n}$  and optimal replica count  $R_{optimal}$ 
35:      $\triangleright$  The solved  $R_{optimal}$  refines/overrides NoofReplicasHeuristic from phase 1
36:   if  $R_{optimal} \neq$  CurrentReplicas then
37:     ScaleDeployment( $R_{optimal}$ )
38:   end if
39:   Deploy tasks to nodes based on the optimized  $x_{p,n}$ 
40: end while
    
```

---

that connection, continuing until 100 000 total requests are completed. The request arrival process throughput measurements obtained using ab tool are based on a deterministic arrival process, where new requests are initiated as soon as existing ones complete, constrained by the specified concurrency level. This approach subjects the system to continuous high stress.

The Apache Benchmark (ab) tool was used with different total request counts (-n parameter), specifically 10 000, 30 000, 50 000, 75 000, and 100 000 requests. The concurrency level was kept constant at 100 (-c 100) for all of these runs. This variation in the total number of requests, while maintaining the same concurrency, effectively changes the duration of the sustained load test. The nature of the request generation

remained a closed-loop model (100 concurrent connections sending requests as fast as the server responds), allowing for the measurement of sustained throughput under different total work performed.

For the context of these experiments, throughput is defined as the average rate at which the system successfully processes requests over the entire duration of a given test run. It is calculated as follows:

$$\text{Throughput} = \frac{\text{Total successfully completed requests}}{\text{Total time taken for the test run [s]}}. \quad (7)$$

This value is directly reported by the ab tool upon completion of each test. During the experiments, it was ensured that the server did not explicitly reject requests due to overload; thus, the measured throughput primarily reflects the sustained

**Tab. 6.** Evaluation metrics for scheduling predictions.

Model	Evaluation metric	CPU	Memory
LSTM	MSE	0.28	0.45
	MAE	0.22	0.39

**Tab. 7.** Optimized LSTM predictions with ILP.

Evaluation metric	CPU	Memory
MSE	0.166	0.304
MAE	0.210	0.445

**Tab. 8.** Evaluation metrics for autoscaling predictions.

Model	Evaluation metric	CPU	Memory
LSTM	MSE	0.156	0.104
	MAE	0.247	0.195

service capacity of the autoscaled application deployment rather than being affected by significant request loss.

**5.1. Scheduling Predictions (Nodes)**

Evaluation of models to predict CPU and memory usage revealed that the LSTM model provides highly accurate predictions. The evaluation metrics for the CPU and memory usage predictions are presented in Table 6. One may notice that the LSTM model achieved an MSE of 0.28 (CPU) and 0.45 (memory), with MAE values of 0.22 (CPU) and 0.39 (memory).

Table 7 demonstrates the impact of optimization using ILP. The MSE and MAE values for CPU and memory predictions are reduced. This improvement is achieved because ILP minimizes the prediction error by systematically adjusting the task allocations, ensuring greater accuracy.

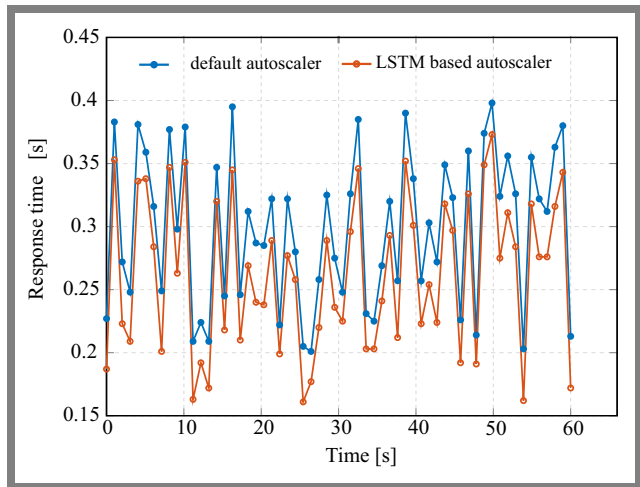
**5.2. Autoscaling Prediction (Pods)**

The LSTM model is also employed to predict CPU and memory usage for the auto-scaling mechanism. The evaluation metrics for autoscaling predictions are detailed in the Tab. 8. The values obtained show the performance of the LSTM model for auto-scaling predictions across CPU and memory usage. The model achieved a mean squared error (MSE) of 0.156 for the CPU and 0.104 for memory, indicating high precision in predicting resource requirements. Similarly, the mean absolute error (MAE) values were 0.247 for the CPU and 0.195 for memory, demonstrating the reliability in minimizing prediction deviations.

Table 9 highlights the improved prediction metrics after optimizing LSTM predictions with ILP. By reducing errors in resource estimation, the system achieves a better alignment between predicted and actual usage, leading to improved resource allocation efficiency.

**Tab. 9.** Optimized LSTM predictions with ILP for autoscaling.

Evaluation metric	CPU	Memory
MSE	0.135	0.089
MAE	0.356	0.202



**Fig. 3.** Response time comparison: HPA vs. LSTM autoscaling.

**5.3. Autoscaling Results**

The evaluation of CPU and memory utilization, as well as response time, was performed under a fixed workload of 50 000 HTTPS requests using the ab tool with 100 concurrent requests at a time. This workload level was selected to simulate medium to high resource utilization scenarios, providing insight into system behavior under significant demand. The metrics obtained offer a detailed view of resource consumption and autoscaling efficiency under real conditions.

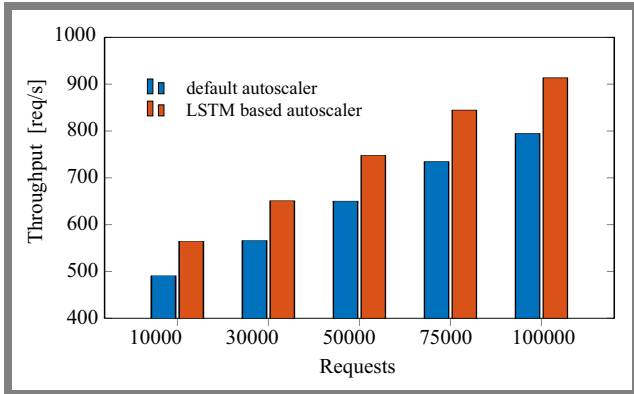
Response time refers to the end-to-end duration measured by the Apache Benchmark client for each individual HTTPS request, capturing the total time from request initiation to the reception of the complete response from the application server. This metric reflects the perceived latency under the applied load. This measured application response time should be distinguished from the predicted network latency parameter  $RTT_n$  used within the ILP objective function, which serves as an internal factor optimized to influence this overall externally measured response time.

Figure 3 presents a comparative analysis of the response times between the default autoscaling mechanisms over time. The graph reveals that the LSTM-based autoscaler consistently maintains lower response times throughout the observation period, with values ranging between 0.15 and 0.5 s. The LSTM model demonstrates superior performance by achieving a 12.5% reduction in response time compared to the default autoscaler.

The default autoscaler exhibits more pronounced fluctuations and generally higher response times, with peaks reaching approximately 0.4 s. In contrast, the LSTM-based autoscaler maintains more stable performance with an average response time of 0.262 s, compared to 0.298 s. This improvement can be attributed to the ability of the LSTM model to predict

**Tab. 10.** Response time metrics: HPA vs. LSTM autoscaling.

Metric	HPA RT	LSTM autoscaling RT
Minimum	0.201 s	0.161 s
Average	0.298 s	0.262 s
Maximum	0.398 s	0.373 s



**Fig. 4.** Throughput comparison: HPA vs. LSTM autoscaling.

resource requirements and proactively adjust allocations, resulting in more efficient resource utilization and reduced latency.

The temporal pattern shows that while both autoscalers experience periodic fluctuations in response times, the LSTM-based approach maintains better consistency and lower overall latency. This improved stability and reduced response time demonstrate the effectiveness of the LSTM model in dynamic resource allocation, particularly in handling varying workload conditions in edge cloud environments.

Table 10 highlights the statistical metrics for response time. LSTM autoscaling consistently achieves better minimum, average, and maximum response times compared to HPA, ensuring faster response to workload fluctuations.

Figure 4 presents a comparison of throughput performance between the LSTM-based autoscaler and the default autoscaler (HPA) across varying request loads from 10 000 to 100 000 requests. The LSTM autoscaler consistently demonstrates superior performance, achieving a 29.2% improvement in throughput compared to HPA. The graph illustrates the progressive increase in throughput as the request volume increases. The LSTM autoscaler maintains a steeper growth trajectory, starting at approximately 564.52 req/sec at 10 000 requests and reaching 913.54 req/s at 100 000 requests. On the contrary, the default autoscaler shows a more modest progression from 490.89 req/s to 794.39 req/s over the same range.

Figure 4 also provides a visualization of the performance gap between the two approaches at specific request intervals. The LSTM autoscaler achieves an average throughput of 744.36 req/s, outperforming the default autoscaler’s 647.18 req/s. This performance differential becomes more pronounced with higher request volumes, demonstrating superior capability in handling increased workload demands.

**Tab. 11.** Summary of performance metrics.

Metric	HPA	LSTM autoscaling
Minimum	490.888 req/s	564.521 req/s
Average	647.169 req/s	744.644 req/s
Maximum	794.385 req/s	913.543 req/s

This improvement in performance can be attributed to the ability of the LSTM model to predict resource requirements and proactively adjust scaling decisions, resulting in more efficient resource utilization and better handling of varying workload patterns in edge cloud environments.

Table 11 presents a summary of the performance for HPA and LSTM autoscaling. The LSTM autoscaling model consistently delivers higher throughput, with a minimum of 564.52 req/s, exceeding HPA’s 490.89 req/s.

The average throughput for LSTM autoscaling is 744.644 req/s, while HPA achieves 647.169 req/s, indicating an improvement. At peak load, LSTM autoscaling reaches a maximum throughput of 913.543 req/s compared to HPA’s 794.385 req/s. These results highlight the scalability and efficiency of LSTM-based autoscaling over HPA.

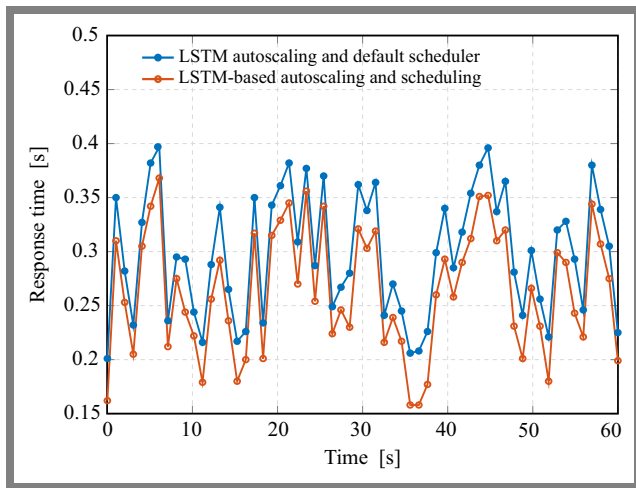
#### 5.4. Combined Autoscaling and Scheduling Results

Integrating LSTM-based autoscaling and scheduling ensures a coordinated approach, improving workload distribution and resource management. Evaluation of CPU and memory utilization, as well as response time, was performed using a fixed workload of 50 000 requests generated with the Apache Benchmark (ab) tool. This workload level was selected to represent a realistic medium-load scenario, providing a comprehensive view of system performance under consistent demand.

Figure 5 illustrates the response time comparison between LSTM autoscaling with default scheduler and LSTM-based autoscaling and scheduling over time. The LSTM autoscaling with default scheduler exhibits higher response times throughout the observation period, fluctuating between 0.201 and 0.397 s. The response time pattern shows notable variations, particularly between the 2040 s interval, indicating less stable performance. On the contrary, the LSTM-based autoscaling and scheduling approach demonstrates consistently lower response times across the entire timeline.

This combined approach maintains response times between 0.158 and 0.368 s, achieving a 12% improvement in the median response time (0.259 s vs. 0.294 s). The graph shows more stable performance with fewer fluctuations, particularly evident in the 30–50 s range, where the response time variations are notably smaller than the default approach.

The improved stability and lower response times can be attributed to the combined effect of auto-scaling with scheduling. The integrated approach demonstrates superior resource allocation efficiency, with the LSTM models working in tandem to predict resource requirements and optimize pod placement. This coordinated decision-making results in more predictable performance patterns and reduced latency, as ev-



**Fig. 5.** Response time vs. time for combined scheduling and autoscaling.

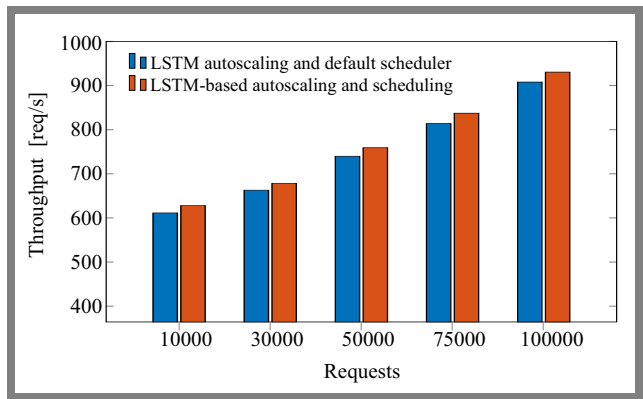
**Tab. 12.** Comparison of response time for combined scheduling and autoscaling.

Metric	LSTM based autoscaling with default scheduler	LSTM autoscaling and scheduling
Minimum	0.201 s	0.158 s
Average	0.294 s	0.259 s
Maximum	0.397 s	0.368 s

identified by the 21.4% improvement in minimum response time (from 0.201 to 0.158 s) and the 7.3% reduction in maximum response time (from 0.397 to 0.368 s). Performance improvements in all metrics underscore the effectiveness of the combined LSTM-based approach in maintaining optimal system responsiveness under varying workload conditions.

Figure 6 presents a comparison of throughput performance between LSTM autoscaling with the default Kubernetes scheduler and the combined LSTM-based autoscaling and scheduling approach. The analysis reveals several significant performance patterns across varying request loads. At the lower end of the request spectrum (10 000 requests), the combined LSTM-based approach demonstrates an initial throughput advantage of 628.14 req/s compared to 611.33 req/s for the default scheduler, that is, a 2.75% improvement. This performance gap widens with load level.

The throughput enhancement becomes pronounced at higher request volumes, reaching 930.42 req/s vs. 907.78 req/s at 100 000 requests, maintaining a 2.49% performance gain. The system shows excellent scalability, with both approaches maintaining near-linear throughput growth from 10 000 to 100 000 requests. The LSTM-based combined approach maintains a consistent performance improvement of average 2.61% at all test points. At medium load (50 000 requests), the LSTM-based system processes 758.99 req/s compared to 739.73 req/s for the default scheduler, demonstrating robust performance under typical operating conditions. The highest absolute performance gain is observed at 100 000 requests, where the LSTM-based system processes an additional 22.64



**Fig. 6.** Comparison of performance for combined scheduling and autoscaling.

**Tab. 13.** Summary of performance metrics: HPA vs. LSTM autoscaling.

Metric	HPA	LSTM autoscaling
Minimum	611.331 req/s	628.14 req/s
Average	747.0824 req/s	766.217 req/s
Maximum	907.778 req/s	930.416 req/s

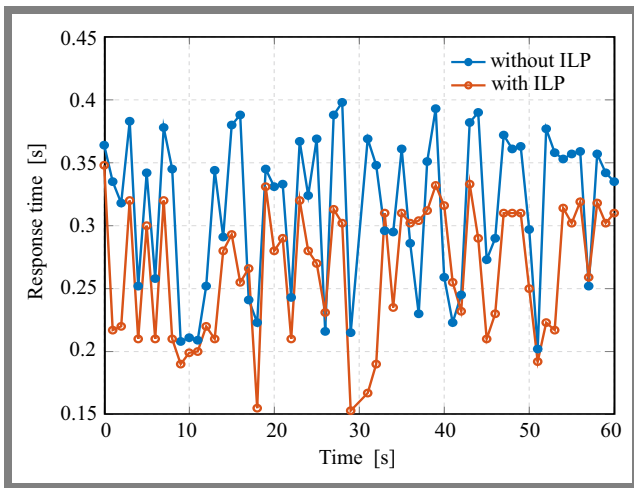
req/s. This sustained performance improvement across all request levels demonstrates the effectiveness of integrating LSTM-based decision-making in both the scheduling and autoscaling components. Table 13 summarizes the throughput performance of HPA and LSTM autoscaling on different request loads. LSTM-based autoscaling consistently outperforms HPA in all scenarios.

### 5.5. Combined Autoscaling and Scheduling with and without ILP

Here a comparative analysis of the combined autoscaling and scheduling approach with and without ILP optimization is provided. The evaluation focuses on response time as the key performance metric and data collected under a fixed workload of 50 000 HTTPS requests using the ab tool. This workload level provides a realistic scenario to assess the efficiency of ILP optimization in minimizing latency.

The response time analysis reveals distinct performance characteristics between the two approaches (Fig. 7). The combined autoscaling and scheduling without ILP exhibits response times fluctuating between 0.21 and 0.40 s, with variations particularly in the 20–30 s interval. In contrast, the ILP-enhanced approach demonstrates superior stability, maintaining response times between 0.15 and 0.33 s with reduced variance. The throughput comparison across varying request loads shows consistent performance advantages for the LSTM-based approach. Starting at 10 000 requests, it achieves 628.14 req/s compared to 611.33 req/s for the default scheduler, i.e. a 2.75% improvement. This performance differential persists through higher loads, reaching 930.42 req/s versus 907.78 req/s at 100 000 requests.

The system demonstrates excellent scalability with near-linear throughput growth throughout the test range, maintaining



**Fig. 7.** Comparison of response times for combined autoscaling and scheduling with and without ILP.

**Tab. 14.** Comparison of response time for combined scheduling and autoscaling.

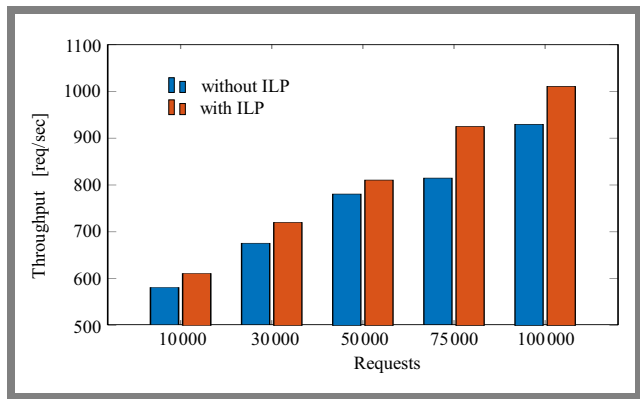
Metric	LSTM autoscaling and scheduling	ILP-enhanced LSTM autoscaling and scheduling
Minimum	0.210 s	0.152 s
Median	0.315 s	0.262 s
Maximum	0.397 s	0.348 s

an average improvement of 2.61% across all test points. At medium load (50 000 requests), the LSTM-based system processes 758.99 req/s compared to 739.73 req/s for the default scheduler, while the highest absolute performance gain is observed at 100 000 requests with an additional 22.64 req/s. The impact of ILP optimization is particularly evident in maintaining more consistent performance levels, especially during the 40–60 s period, where it stabilizes around 0.30 s, demonstrating enhanced efficiency in resource allocation and workload distribution.

The metrics in the Tab. 14 highlight the improvements achieved by the ILP-enhanced LSTM autoscaling and scheduling approach. Compared to the non-ILP method, the median response time was reduced from 0.315 to 0.263 s, demonstrating enhanced efficiency, particularly in handling dynamic workloads with lower latency.

Figure 8 illustrates an analysis of throughput performance between the ILP-enhanced and standard LSTM approaches for combined autoscaling and scheduling. The ILP-enhanced solution demonstrates superior performance across all request volumes, with the average throughput increasing from 756.34 req/s to 815.50 req/s, representing a 7.8% improvement.

The performance advantage becomes more pronounced under higher workloads, with maximum throughput reaching 1 010.75 req/s compared to 929.68 req/s in the non-ILP approach, showing an 8.7% increase. Even at lower request volumes, the ILP-enhanced method maintains better efficiency, with minimum throughput improving from 581.02 req/s



**Fig. 8.** Performance comparison for combined LSTM autoscaling and scheduling with and without ILP.

**Tab. 15.** Performance metrics: combined LSTM autoscaling and scheduling vs. ILP-enhanced approach.

Metric	LSTM autoscaling and scheduling	ILP-enhanced LSTM autoscaling and scheduling
Minimum	581.02 req/s	610.88 req/s
Average	756.34 req/s	815.50 req/s
Maximum	929.68 req/s	1010.75 req/s

to 610.88 req/s. The graph highlights consistent performance gains across all request volumes, particularly in the 75 000–100 000 request range, where the system demonstrates optimal resource utilization and workload management capabilities. Table 15 provides detailed comparison of throughput metrics.

## 6. Conclusions

This study addresses the limitations of default Kubernetes resource management by proposing an integrated framework that combines LSTM-based autoscaling and scheduling with ILP-based optimization. Using predictive modeling in conjunction with intelligent resource allocation, the ILP and LSTM-based system improves overall efficiency. A comparative evaluation between the LSTM-based autoscaling and scheduling system and the ILP- and LSTM-based system demonstrated a 12.34% reduction in response time and a 7.85% increase in throughput.

Despite the improvement in results, several limitations must be considered. The experiments were conducted in a virtualized environment using a stateless Web application. While efforts were made to approximate real-world conditions, the controlled nature of the testbed may not fully reflect the complexities encountered in practical deployments. Furthermore, the LSTM model was specifically trained and fine-tuned for the given use case. Applying the framework to other applications would likely require retraining the model with domain-specific historical data and adjusting parameters to suit different system dynamics.

## References

- [1] L.H. Phuc, L.-A. Phan, and T. Kim, "Traffic-aware Horizontal Pod Autoscaler in Kubernetes-based Edge Computing Infrastructure", *IEEE Access*, vol. 10, pp. 18966–18977, 2022 (<https://doi.org/10.1109/ACCESS.2022.3150867>).
- [2] S.T. Singh, M. Tiwari, and A.S. Dhar, "Machine Learning based Workload Prediction for Auto-scaling Cloud Applications", 2022 *OPJU International Technology Conference on Emerging Technologies for Sustainable Development (OTCON)*, Raigarh, India, 2023 (<https://doi.org/10.1109/OTCON56053.2023.10114033>).
- [3] I. Ahmad, M.G. AlFailakawi, A. AlMutawa, and L. Alsaman, "Container Scheduling Techniques: A Survey and Assessment", *Journal of King Saud University – Computer and Information Sciences*, vol. 34, pp. 3934–3947, 2022 (<https://doi.org/10.1016/j.jksuci.2021.03.002>).
- [4] K. Senjab, S. Abbas, N. Ahmed, and A.R. Khan, "A Survey of Kubernetes Scheduling Algorithms", *Journal of Cloud Computing*, vol. 12, art. no. 87, 2023 (<https://doi.org/10.1186/s13677-023-00471-1>).
- [5] K. Zhu *et al.*, "Proactive Hybrid Autoscaling for Container-Based Edge Applications in Kubernetes", *Lecture Notes of the Institute for Computer Sciences*, vol. 574, pp. 330–345, 2024 ([https://doi.org/10.1007/978-3-031-65123-6\\_24](https://doi.org/10.1007/978-3-031-65123-6_24)).
- [6] Z. Wang, Q. Zhu, and Y. Hou, "Multiworkflow Scheduling in Edge-cloud Computing by African Vulture Optimization Algorithm", 2024 *11th International Forum on Electrical Engineering and Automation (IFEAA)*, Shenzhen, China, 2024 (<https://doi.org/10.1109/IFEAA64237.2024.10878706>).
- [7] J. Li, P. Singh, and S. Toor, "Proactive Autoscaling for Edge Computing Systems with Kubernetes", *Proc. of the 14th IEEE/ACM International Conference on Utility and Cloud Computing Companion (UCC'21)*, art. no. 22, pp. 1–8, 2022 (<https://doi.org/10.1145/3492323.3495588>).
- [8] L.H. Phuc *et al.*, "Node-based Horizontal Pod Autoscaler in KubeEdge-based Edge Computing Infrastructure", *IEEE Access*, vol. 10, pp. 134417–134426, 2022 (<https://doi.org/10.1109/ACCESS.2022.3232131>).
- [9] J. Dogani and F. Khunjush, "Proactive Auto-scaling Technique for Web Applications in Container-based Edge Computing Using Federated Learning Model", *Journal of Parallel and Distributed Computing*, vol. 187, art. no. 104837, 2024 (<https://doi.org/10.1016/j.jpdc.2024.104837>).
- [10] T.-X. Do and V.K.N. Tan "Hybrid Autoscaling Strategy on Container-Based Cloud Platform", *International Journal of Software Innovation*, vol. 10, 2022, pp. 1–12. (<https://doi.org/10.4018/IJSI.292019>).
- [11] X. Feng *et al.*, "Adaptive Container Auto-scaling for Fluctuating Workloads in Cloud", *Future Generation Computer Systems*, vol. 172, art. no. 107872, 2025 (<https://doi.org/10.1016/j.future.2025.107872>).
- [12] S.D. Konidena, "Efficient Resource Allocation in Kubernetes Using Machine Learning", *International Journal of Innovative Science and Research Technology*, vol. 9, 2024 (<https://doi.org/10.38124/ijisrt/IJISRT24JUL607>).
- [13] J. Zhou, S. Pal, C. Dong, and K. Wang, "Enhancing Quality of Service Through Federated Learning in Edge-cloud Architecture", *Ad Hoc Networks*, vol. 156, art. no. 103430, 2024 (<https://doi.org/10.1016/j.adhoc.2024.103430>).
- [14] S.-H. Kim and T. Kim, "Local Scheduling in KubeEdge-based Edge Computing Environment", *Sensors*, vol. 23, art. no. 1522, 2023 (<https://doi.org/10.3390/s23031522>).
- [15] M. Raesi-Varzaneh, O. Dakkak, A. Habbal, and B.-S. Kim, "Resource Scheduling in Edge Computing: Architecture, Taxonomy, Open Issues and Future Research Directions", *IEEE Access*, vol. 11, pp. 25329–25350, 2023 (<https://doi.org/10.1109/ACCESS.2023.3256522>).
- [16] Z. Shi and Z. Shi, "Multi-node Task Scheduling Algorithm for Edge Computing Based on Multi-Objective Optimization", *Journal of Physics: Conference Series*, vol. 1607, art. no. 012017, 2020 (<https://doi.org/10.1088/1742-6596/1607/1/012017>).
- [17] K. Wang *et al.*, "Computing Aware Scheduling in Mobile Edge Computing System", *Wireless Networks*, vol. 27, pp. 4229–4245, 2021 (<https://doi.org/10.1007/s11276-018-1892-z>).
- [18] G. Vijayasekaran and M. Durairandian, "Resource Scheduling in Edge Computing IoT networks Using Hybrid Deep Learning Algorithm", *System Research and Information Technologies*, pp. 86–101, 2022 (<https://doi.org/10.20535/SRIT.2308-8893.2022.3.06>).
- [19] Y. Lu *et al.*, "EA-DFPSO: An Intelligent Energy-efficient Scheduling Algorithm for Mobile Edge Networks", *Digital Communications and Networks*, vol. 8, pp. 237–246, 2022 (<https://doi.org/10.1016/j.dcan.2021.09.011>).
- [20] P. Khoshvaght *et al.*, "A Multi-objective Deep Reinforcement Learning Algorithm for Spatio-temporal Latency Optimization in Mobile IoT-enabled Edge Computing Networks", *Simulation Modelling Practice and Theory*, vol. 143, art. no. 103161, 2025 (<https://doi.org/10.1016/j.simpat.2025.103161>).
- [21] P. Vishesh *et al.*, "Optimized Placement of Service Function Chains in Edge Cloud with LSTM and ILP", *SN Computer Science*, vol. 6, art. no. 44, 2024 (<https://doi.org/10.1007/s42979-024-03539-0>).
- [22] Z. Ding and Q. Huang, "COPA: A Combined Autoscaling Method for Kubernetes", 2021 *IEEE International Conference on Web Services (ICWS)*, Chicago, USA, 2021 (<https://doi.org/10.1109/ICWS53863.2021.00061>).

---

### Shivan Singh, B.Eng.

School of Computer Science and Engineering

 <https://orcid.org/0009-0004-7894-3858>

E-mail: 01fe21bcs246@kletech.ac.in

KLE Technological University, Hubballi, Karnataka, India

<https://www.kletech.ac.in>

### Narayan D.G., Ph.D.

School of Computer Science and Engineering

 <https://orcid.org/0000-0002-2843-8931>

E-mail: narayan\_dg@kletech.ac.in

KLE Technological University, Hubballi, Karnataka, India

<https://www.kletech.ac.in>

### Sadaf Mujawar, M.Tech.

Department of Computer Science and Engineering

 <https://orcid.org/0009-0007-5434-0114>

E-mail: sadaf.savanur@kletech.ac.in

KLE Technological University, Hubballi, Karnataka, India

<https://www.kletech.ac.in>

### G.S. Hanchinamani, Ph.D.

Department of Computer Science and Engineering

 <https://orcid.org/0000-0002-8791-0351>

E-mail: gs\_hanchinamani@kletech.ac.in

KLE Technological University, Hubballi, Karnataka, India

<https://www.kletech.ac.in>

### P.S. Hiremath, Ph.D.

Department of MCA

 <https://orcid.org/0000-0001-7640-6937>

E-mail: pshiremath@kletech.ac.in

KLE Technological University, Hubballi, Karnataka, India

<https://www.kletech.ac.in>

# TinyML-driven Sensor Nodes for Energy-efficient Acoustic Event Detection in Pervasive Acoustic WSNs

Bibek B. Roy<sup>1</sup>, Sushovan Das<sup>2</sup>, and Uttam Kr. Mondal<sup>1</sup>

<sup>1</sup>Vidyasagar University, Midnapore, WB, India,

<sup>2</sup>College of Engineering & Management, Kolaghat, WB, India

<https://doi.org/10.26636/jtit.2025.2.2084>

**Abstract** — The process of sensing and transmitting acoustic signals by pervasive acoustic wireless sensor networks (PAWSNs) poses considerable energy challenges. These problems may be mitigated by filtering only relevant acoustic events from the sensor network. By reducing the number of acoustic events, the frequency of communication may be decreased, thereby enhancing energy efficiency. Although traditional machine learning models are capable of predicting relevant acoustic events by being trained on suitable data sets, they are impractical for direct implementation on resource-limited acoustic sensor nodes. To address this issue, this research introduces TinyML-based acoustic event detection (AED) models which facilitate efficient real-time processing on microcontrollers with scarce hardware resources. The study develops several TinyML models using an environmental dataset and evaluates their accuracy. These models are then deployed in hardware to assess their performance in terms of AED. Thanks to such an approach, only predicted events that exceed a certain threshold are transmitted to the base station via router nodes, which reduces the transmission burden, thus improving energy efficiency of PAWSNs. Real-time experiments confirm that the proposed method significantly improves energy efficiency and boosts node lifetime.

**Keywords** — *acoustic event detection, energy efficiency, pervasive acoustic WSN, TinyML*

## 1. Introduction

Pervasive acoustic wireless sensor networks (PAWSNs) [1], [2] are composed of ubiquitous sensors dedicated to monitoring various environmental settings. Such networks are essential in real-time applications, such as underwater monitoring, detection of anomalies in industrial devices, smart city infrastructure, and observation of wildlife. However, spatially distributed battery-operated acoustic sensors continuously capture, process and transmit acoustic signals, raising concerns about excessive energy consumption and reducing the overall lifespan of the network. The conventional approach relying on central processing of the sensed data is inefficient for large volumes of data produced by PAWSNs. Typically, each detected acoustic event is sent to a base station (BS) through router nodes using traditional routing algorithms. This results in significant amounts of energy being consumed to conduct the transmissions and causes network congestion.

It also needs to be borne in mind that PAWSNs, which operate in various environments, are faced with distinctive acoustic conditions.

Although traditional ML models could address this diversity by predicting or filtering acoustic events to lower transmission costs, they require large amounts of computing resources which are often unavailable in battery-powered leaf nodes that offer limited processing power. Therefore, this study introduces a lightweight, decentralized method that reduces data transmission while preserving high AED prediction accuracy.

This work proposes a TinyML-powered AED model that enables lightweight ML models to run on battery-operated microcontrollers with limited hardware, as a low-power real-time solution. Due to this, such devices as the Raspberry Pi platform are capable of predicting specific acoustic events at leaf nodes and only send, to the base station, the detected filtered events, thus significantly lowering energy use by reducing the amount of data exchanged between leaf and router nodes. This decentralized method addresses the challenge posed by the pervasive acoustic environment, as it decreases transmission loads, extends the useful life of the network, and improves the overall effectiveness of the system.

This work also employs a context-aware TinyML model selection strategy, in which models are deployed and tested for optimal performance on leaf nodes, based on environmental conditions. By selecting models with the highest level of accuracy for AED, each leaf node effectively detects relevant acoustic events in its specific acoustic zone. Moreover, due to security concerns, such as intruders mimicking environmental sounds, TinyML models must be updated regularly based on newly collected data, reaching beyond the initial dataset. The base station periodically updates the models and re-sends them to specific leaf nodes to maintain long-term event detection accuracy.

The primary advantages of the proposed work include incorporating TinyML-driven AED models into leaf nodes, optimization of PAWSN efficiency by minimizing the transmission of unnecessary acoustic data and decreasing storage needs, all while preserving high accuracy and ensuring secure environmental monitoring. This study helps improve AI-driven lightweight TinyML models, with the aim of en-

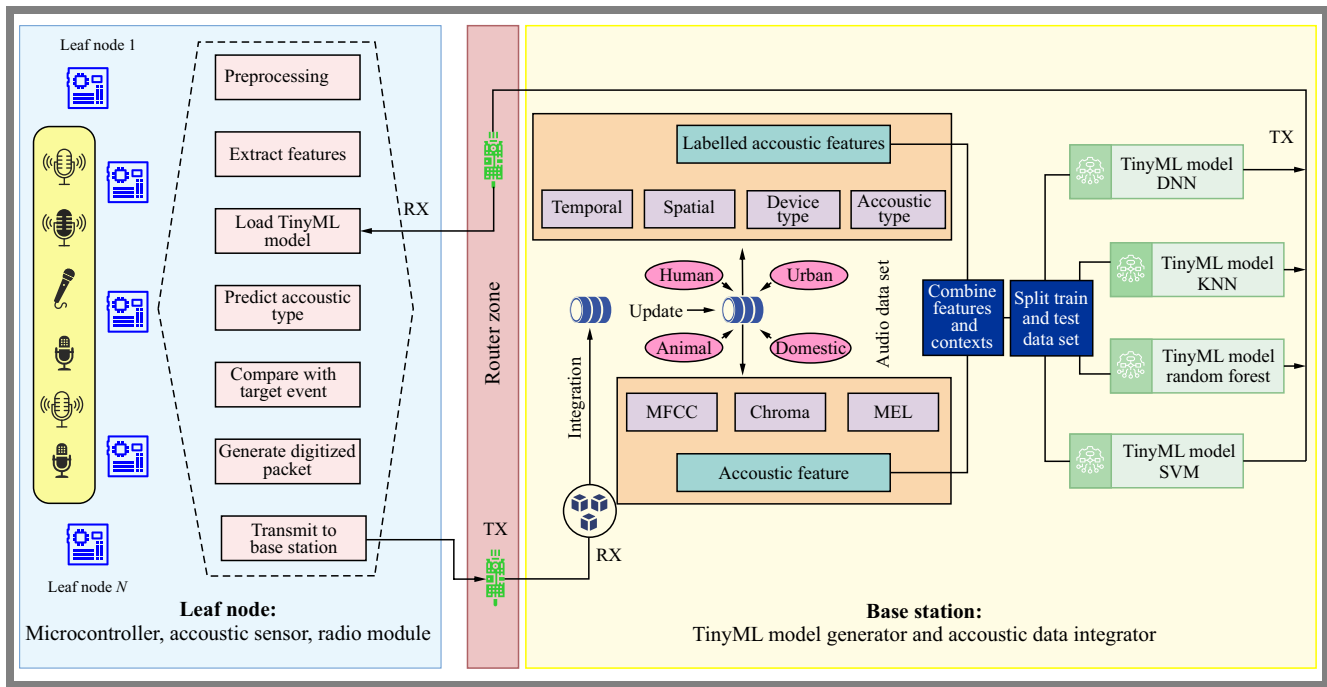


Fig. 1. PAWSN system architecture.

uring energy-efficient real-time detection of acoustic events with the use of resource-limited battery-powered acoustic sensors within PAWSN settings. The framework facilitates the prolonged deployment of PAWSNs across diverse monitoring applications by enabling nodes to dynamically adjust and update their models in response to changing acoustic environments, while simultaneously minimizing energy consumption.

## 2. Literature Review

Incorporating TinyML-based models into PAWSNs is crucial, as it allows real-time on-device processing in battery-operated, resource-limited acoustic sensor nodes [3]–[10]. Since the process of monitoring environmental conditions in real time requires acoustic data to be transmitted on a continuous basis, an activity that results in significant energy consumption [11]–[13], viable solutions must incorporate mechanisms to filter irrelevant data. Traditional ML models, which require considerable computational capabilities, are incompatible with resource-limited sensor nodes. In contrast, TinyML is designed to operate efficiently on such sensor nodes, using minimal computational resources and energy, especially to detect acoustic events with a high accuracy rate at the sensor node level [11].

Optimizing and selecting models, such as quantifying and deploying neural networks, is crucial for sensor nodes with limited resources. These models aim to balance computational efficiency with improved precision [13]. The challenges posed by the dynamic nature of acoustic environments require periodic updating of ML models to ensure long-term viability. Various learning methodologies are used to update TinyML

models at the base station, improving PAWSN’s ability to adapt effectively to changing environmental conditions [14]. Furthermore, the choice of low-power hardware is crucial for implementation in real-time environmental or industrial settings, such as detecting accidents involving workers through audio classification [15].

Although numerous studies [11], [13]–[15] focus on resolving issues such as achieving high accuracy and energy efficiency, many problems caused by challenging environmental conditions and noise still exist in practical applications. To address these challenges, future research should emphasize the development of TinyML model architectures by investigating supervised and unsupervised learning approaches.

Multiple studies have explored energy efficiency in PAWSNs by optimizing data collection, improving routing protocols, and incorporating ML-based techniques for data integration in BS. In [16], the authors introduced an energy efficient data aggregation algorithm that effectively collected data for unmanned aerial vehicles (UAVs). The authors of [17] presented an energy efficient data collection approach using autonomous underwater vehicles (AUVs) for underwater acoustic sensor networks, highlighting the progress in optimizing data transmission strategies.

In [18], the authors introduced a single-relay selective method for WSNs aimed at reducing energy use while maintaining reliable communication. Similarly, [19] explored energy efficiency in the industrial monitoring of WSNs with limited energy, highlighting the advantages of a cooperative communication protocol. The strategies implemented by these researchers promote more sustainable and durable sensor networks by reducing transmission overhead and optimizing resource use.

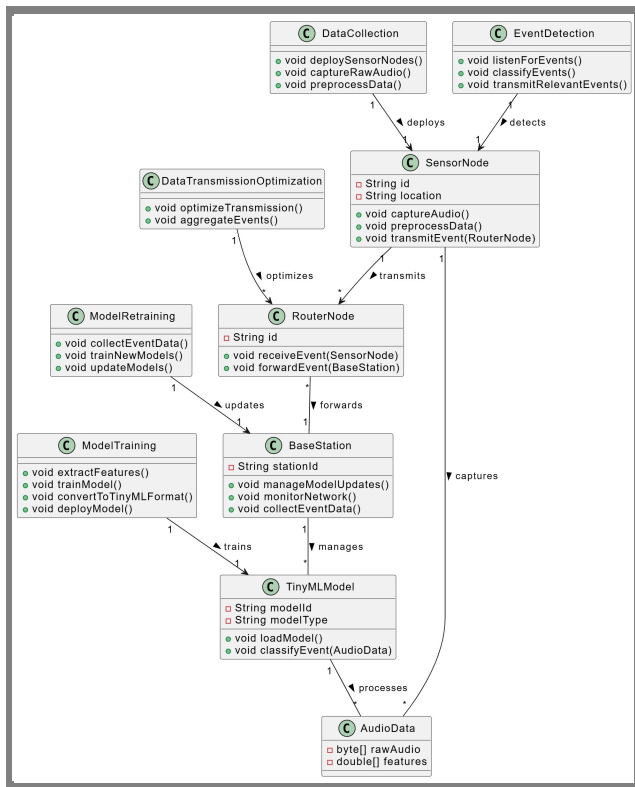


Fig. 2. PAWSN workflow with TinyML-driven sensor.

Incorporation of TinyML-based AED models into sensor nodes is becoming a trend for energy-efficient PAWSNs. Study [11] explored the use of RNN-based ML models in low-power sensor devices for AED, achieving real-time processing. In [14], the researchers presented a versatile and unsupervised TinyML framework designed to identify anomalies in industrial settings, incorporating on-device learning and training capabilities.

Although technological advancements have been introduced to sensors, communications, data processing, integration, and ML-based prediction, PAWSNs still face challenges related to energy efficiency and security surveillance during continuous environmental monitoring. Furthermore, complex processing of acoustic signals and handling of large amounts of acoustic data can negatively impact the accuracy of event detection processes.

### 3. Proposed Technique

Figure 1 illustrates the proposed TinyML-based PAWSN architecture which consists of three primary components: leaf nodes, a router zone, and a base station. The leaf nodes, located at the sensing edge, contain microcontrollers integrated with acoustic sensors and a radio module. Each leaf node performs local acoustic processing, including feature extraction (e.g. MFCC, chroma, MEL), TinyML model inference, and acoustic-type prediction. If the prediction aligns with a pre-defined target event, the node generates a digitized packet and transmits it to the base station through the router zone. These nodes are battery operated, and due to the use of compact

TinyML models (with a code size of 200 KB, their power consumption remains low at 0.25 ... 0.45 W), allowing extended operational lifetimes depending on the duty cycle.

The router zone enables multihop or range-extended wireless communication using such platforms as ESP-Now or XBee. The BS acts as the central intelligence hub, typically a high-end GPU-based computer, as it integrates and labels the acoustic dataset, extracts contextual features (temporal, spatial, device type), splits the data for training and testing, and generates optimized TinyML models.

Although the PC is essential during model development and training, it can also be employed in field deployments for real-time integration and reconfiguration, although it can be replaced with lightweight embedded computing platforms in resource-constrained environments.

TinyML model development [20], [21] frequently favors the use of SVM, KNN, random forest, and dense neural networks (DNN) due to their empirical performance in acoustic classification tasks and their compatibility with microcontroller-class hardware constraints. These models are effective in classifying structured, low-dimensional acoustic features such as MFCCs, chroma, and MEL spectrograms, which are essential for identifying diverse environmental sound categories, as found in the ESC-50 dataset. Each algorithm was configured through preliminary tuning to balance inference accuracy and computational efficiency on resource-constrained hardware.

Specifically, we used an RBF kernel for SVM, KNN with  $k = 5$ , a random forest with 40 trees, and a DNN with two hidden layers (64 and 32 neurons, respectively). Classical models, such as decision trees and naive Bayes classifiers, although computationally lightweight, underperformed in generalization during initial evaluations. Random forest was selected over single decision trees for its ensemble robustness, while naive Bayes was excluded due to its strong independence assumptions, which are not well suited to the correlated nature of time-frequency audio features. The selected models strike an optimal trade-off between expressiveness and resource use, making them suitable for quantization and deployment on devices like the Seed XIAO ESP32S3 for energy-efficient, real-time inference in edge-based acoustic wireless sensor networks.

The data analysis begins by categorically separating various acoustic classes from 40 different environmental sounds, such as animal, urban, and human noises, using the ESC-50 dataset, complete with the appropriate labeling and target naming. Then, several TinyML models are trained using different classification algorithms such as K-nearest neighbors (KNN), support vector machine (SVM), random forest (RF), and DNN to handle the specific acoustic data present in the dataset. These models employ Mel-frequency cepstral coefficients (MFCCs) and chroma features to extract information from each category of environmental data. Next, the different TinyML models tailored for various environments are deployed to detect specific events. Upon sensing an acoustic event, the model forecasts its type. If the prediction score exceeds a predetermined threshold, only then is the detected

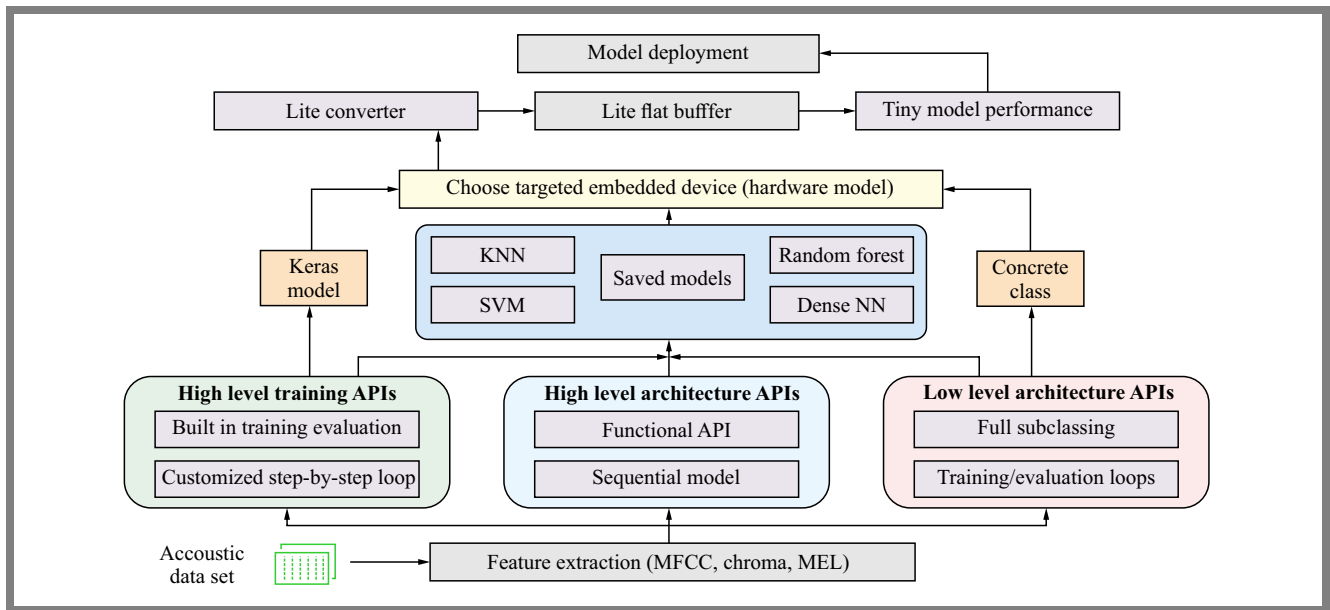


Fig. 3. Generation and deployment of the TinyML model for acoustic sensor nodes.

acoustic signal forwarded to router nodes for transmission to the BS.

Figure 2 illustrates the operational flow of the processes associated with the proposed system. Sensor nodes capture and preprocess raw audio data before transmitting filtered events to router nodes. The event detection module ensures accurate sound classification, while the data transmission optimization module reduces redundancy by aggregating and filtering events. The router nodes forward the filtered events to the BS which manages network operations, collects data for model updates, and optimizes the system’s performance. The TinyML model is responsible for audio classification, with additional modules handling feature extraction, model training, and retraining to improve accuracy and adaptability over time.

The process of creating a TinyML model, as depicted in Fig. 3, begins with an acoustic dataset, where the raw audio data undergoes feature extraction using techniques such as MFCC, chroma, and MEL spectrograms. These extracted features serve as numerical representations of the sound data, which are then processed to train machine learning (ML) models.

The architecture leverages TensorFlow and Keras, providing multiple abstraction levels for model development. In high-level architectures, models can be structured using either the functional API or sequential model, offering flexibility in defining deep learning networks. For more advanced customization, the low-level API allows for full subclassing, where users can define custom layers and training loops. Training is facilitated through high-level training APIs, which include built-in evaluation methods and customized step-by-step loops for model optimization.

Once models are trained, they are saved for evaluation, where multiple algorithms such as SVM, KNN, RF and DNN are tested for performance. After selecting the best-performing model, the next step is to prepare it for embedded deployment.

The model is optimized by choosing a targeted embedded device, with the next step consisting in converting it into a TensorFlow Lite (TFLite) model using the TFLite converter. To ensure that the model remains lightweight and efficient, the TFLite Flat buffer is used, facilitating model compression and quantization for tiny model’s performance.

The final model is then deployed on edge devices, enabling low-power, real-time inference for applications such as IoT, smart devices, and real-time acoustic analysis. This process ensures that TinyML models are optimized to run on microcontrollers and embedded hardware while maintaining accuracy and efficiency.

### 3.1. TinyML-based Acoustic Event Detection

Algorithm 1 defines TinyML-based acoustic event detection with network simulation for efficient real-time sound classification and data transmission. Audio signals are processed in a loop over a predefined sensing period, classifying each captured sound using a TinyML model adapted to the sensor’s zone. A confidence score is computed, and only high-confidence classifications are transmitted via XBee to the router nodes, which forward the data to the BS.

## 4. Mathematical Model

To quantify the energy efficiency of the proposed TinyML-driven acoustic event detection framework, this work develops a mathematical model considering event sensing, data transmission, and model retraining. Energy consumption is estimated using power specifications and operating durations derived from the Grove sensor, the XBee module, and the XIAO ESP32S3 microcontroller. Total energy is calculated as a weighted sum of the energy used during the sensing, pro-

---

**Algorithm 1** TinyML-based AED for PAWSN.
 

---

```

1: Input: Acoustic signals from pervasive sensor zone
2: Output: Transmission of detected signals to BS
3: Parameters:  $T_{total}$  = total sensing duration,  $T_{interval}$  =
   sensing interval,  $S_{threshold}$  = prediction score threshold
4: procedure NETWORKSIMULATION-TINYML()
5:   Deploy acoustic sensor nodes in pervasive zones
6:   Establish wireless communication using XBEE
7:   Sense acoustic signals and preprocess (noise
   reduction, normalization, feature extraction)
8:   Train ML models using suitable dataset with various
   ML algorithms at BS
9:   Convert trained models to TinyML-compatible
   formats (TensorFlow Lite or pickle)
10:  Deploy models onto leaf nodes
11:   $t = 0$ 
12:  while  $t < T_{total}$  do
13:    Sense acoustic signal from sensor zone
14:    Preprocess and extract MFCC features
15:    for each TinyML model do
16:      Perform AED using respective model
17:      Compute confidence score  $S$ 
18:      if  $S < S_{threshold}$  then
19:        Discard acoustic signal
20:      else
21:        Transmit detected event to routers
22:        Router node forwards data to the BS
23:      end if
24:    end for
25:    Wait for  $T_{interval}$  before next sensing cycle
26:     $t = t + T_{interval}$ 
27:  end while
28:  BS collects, stores and integrates acoustic data
29:  Retrain TinyML models periodically
30:  Redeploy updated models to sensor nodes
31: end procedure
    
```

---

cessing, and communication phases, based on realistic task durations.

The total energy consumption of the PAWSN can be expressed as:

$$E_{total} = E_{leaf} + E_{router}, \quad (1)$$

where  $E_{leaf}$  is the total energy consumed by all leaf nodes and  $E_{router}$  stands for the total energy consumed by all router nodes.

#### 4.1. Energy Model for Leaf Nodes

The energy consumed by a single leaf node is given by:

$$\begin{aligned}
 E_{leaf} = & \sum_{i=1}^{T_{obs}/T_{event}} (P_{sense} T_{sense,i} + P_{process} T_{process,i}) \\
 & + \sum_{i=1}^{T_{obs}/T_{event}} (P_{tx} T_{tx,i}), \quad (2)
 \end{aligned}$$

where:  $P_{sense}$ ,  $P_{process}$ , and  $P_{tx}$  are the power consumption rates for sensing, processing, and transmission.  $T_{sense,i}$ ,  $T_{process,i}$ ,  $T_{tx,i}$  are the corresponding time durations at interval  $i$ .

Since TinyML reduces transmission data volume, the effective energy consumption for a single-leaf node with TinyML can be expressed as:

$$\begin{aligned}
 E_{leaf}^{TinyML} = & \sum_{i=1}^{T_{obs}/T_{event}} (P_{sense} T_{sense,i} + P_{process} T_{process,i}) \\
 & + \sum_{i=1}^{T_{obs}/T_{event}} (P_{tx} T_{tx,i} \pi_{detect,i}), \quad (3)
 \end{aligned}$$

where  $\pi_{detect,i}$  accounts for the probability of a transmission-triggering event being detected.

The probability of detecting an acoustic event in time interval  $i$  is as follows:

$$\pi_{detect,i} = \frac{SNR_i}{SNR_i + \theta} \cdot \Phi(F_i, M), \quad (4)$$

where  $SNR_i$  is the signal-to-noise ratio in time interval  $i$ ,  $\theta$  is the detection threshold, a system-defined parameter that determines the sensitivity to acoustic events, and  $\Phi(F_i, M)$  is the TinyML model classification confidence function, which depends on:

- $F_i$  – feature vector extracted from the acoustic signal at interval  $i$ ,
- $M$  – TinyML model used for event classification.

#### 4.2. Energy Model for Router Nodes

The router nodes are responsible for receiving and forwarding data to the BS. The energy consumption per router node is given by:

$$E_{router} = \sum_{i=1}^{T_{obs}/T_{event}} (P_{rx} T_{rx,i} + P_{tx} T_{tx,i}), \quad (5)$$

where  $P_{rx}$ ,  $P_{tx}$  are the power consumption rates for receiving and transmitting data, while  $T_{rx,i}$ ,  $T_{tx,i}$  are the corresponding time durations at interval  $i$ .

With TinyML-based event detection, fewer events are transmitted from leaf nodes, reducing the forwarding burden on router nodes represented by:

$$E_{router}^{TinyML} = \sum_{i=1}^{T_{obs}/T_{event}} (P_{rx} T_{rx,i} P_{detect,i} + P_{tx} T_{tx,i} P_{detect,i}), \quad (6)$$

#### 4.3. Total Energy Consumption and Efficiency

The total energy consumption of the network is represented by:

$$E_{total} = \sum_{n \in N_{leaf}} E_{leaf,n} + \sum_{m \in N_{router}} E_{router,m}, \quad (7)$$

where  $N_{leaf}$  and  $N_{router}$  are the total numbers of leaf and router nodes, respectively.

With the proposed TinyML-based approach, the total energy consumption is:

$$E_{total}^{TinyML} = \sum_{n \in N_{leaf}} E_{leaf,n}^{TinyML} + \sum_{m \in N_{router}} E_{router,m}^{TinyML}, \quad (8)$$

Energy efficiency improvement can be quantified as:

$$\eta = \frac{E_{total} - E_{total}^{TinyML}}{E_{total}} \times 100\%, \quad (9)$$

where higher values of  $\eta$  indicate greater energy savings using the proposed model.

## 5. Experimental Results

To assess the effectiveness of the proposed system, this work sets up an acoustic wireless sensor network consisting of four sensor nodes, each equipped with a different training model for detecting various types of acoustic events. The network also includes four router nodes and a central BS. This study evaluates and compares different TinyML models with standard ML models in terms of hardware requirements and energy efficiency.

To establish a wireless acoustic sensor network, we use four Seeed XIAO ESP32S3 microcontrollers, each integrated with Grove sound sensors (LM358) and XBee radio modules as battery-operated leaf nodes. Four additional XBee modules serve as router nodes to enable multi-hop communication. A high-performance BS PC computer handles model training and evaluation.

The complete software stack is implemented in Python 3.10, using TensorFlow, Keras, scikit-learn, and Librosa for pre-processing, feature extraction, and model inference. Feature extraction (e.g., MFCCs) and model training are performed entirely on the BS, not on the sensor nodes. On the other hand, feature extraction for sensed audio signals is performed within the leaf nodes. The work deploys the trained TinyML models, quantized and memory-optimized (e.g., tflite, pkl) to the microcontrollers, which execute MicroPython-based scripts to perform real-time inference on acoustic inputs. Each node is configured to detect a specific class of sound events and transmits data only when a detection exceeds a confidence threshold, thus minimizing power consumption and communication overhead.

Table 1 summarizes hardware-related requirements for running TinyML and standard ML models on the proposed network.

### 5.1. System Output

Figures 4–5 describe the output for a single listening of the PAWSN for four different pervasive zones with different types of acoustic events at a specific time instant using TinyML models. The TinyML model is based on the RF algorithm, chosen for its performance, as shown in Figure 8.

Figure 4 illustrates that only high-confidence acoustic events are transmitted, reducing unnecessary transmissions and improving energy efficiency, although some lower-confidence

```
Sensor_Node_1 is capturing and processing acoustic event from pervasive Zone...1
Sensor_Node_1 detected 'dog' with confidence 0.71, forwarding to Router_2
Router_2 received 'dog' with confidence 0.71, forwarding to Base_Station
Base_Station received final event: 'dog' with confidence 0.71, logging data.
-----
Sensor_Node_2 is capturing and processing acoustic event from pervasive Zone...2
Sensor_Node_2 discarded low-confidence detection (0.64)
Sensor_Node_3 is capturing and processing acoustic event from pervasive Zone...3
Sensor_Node_3 detected 'clapping' with confidence 0.81, forwarding to Router_1
Router_1 received 'clapping' with confidence 0.81, forwarding to Base_Station
Base_Station received final event: 'clapping' with confidence 0.81, logging data.
-----
Sensor_Node_4 is capturing and processing acoustic event from pervasive Zone...4
Sensor_Node_4 discarded low-confidence detection (0.70)
```

**Fig. 4.** Sensor node output during filtering of acoustic events according to the accuracy level ( $> 0.75$ ).

```
Sensor_Node_1 is capturing and processing acoustic event from pervasive Zone...1
Sensor_Node_1 detected 'dog' with confidence 0.71 (All Events Mode), forwarding to Router_4
Router_4 received 'dog' with confidence 0.71, forwarding to Base_Station
Base_Station received final event: 'dog' with confidence 0.71, logging data.
-----
Sensor_Node_2 is capturing and processing acoustic event from pervasive Zone...1
Sensor_Node_2 detected 'frog' with confidence 0.64 (All Events Mode), forwarding to Router_4
Router_4 received 'frog' with confidence 0.64, forwarding to Base_Station
Base_Station received final event: 'frog' with confidence 0.64, logging data.
-----
Sensor_Node_3 is capturing and processing acoustic event from pervasive Zone...1
Sensor_Node_3 detected 'clapping' with confidence 0.81 (All Events Mode), forwarding to
Router_4 Router_4 received 'clapping' with confidence 0.81, forwarding to Base_Station
Base_Station received final event: 'clapping' with confidence 0.81, logging data.
-----
Sensor_Node_4 is capturing and processing acoustic event from pervasive Zone...1
Sensor_Node_4 detected 'rain' with confidence 0.70 (All Events Mode), forwarding to Router_1
Router_1 received 'rain' with confidence 0.70, forwarding to Base_Station
Base_Station received final event: 'rain' with confidence 0.70, logging data.
```

**Fig. 5.** Sensor node output during acoustic event processing without filtering.

```
Energy Consumption@Seeed XIAO ESP32S3 Microcontroller:
-----
| Model Type           | Total Energy Consumption (J) |
-----
| Proposed Model       | 3.08                         |
| All Events Detected  | 5.19                         |
-----
Energy Efficiency Improvement: 40.56%
```

**Fig. 6.** Energy consumption with and without filtering at a leaf node.

detections are discarded. In contrast, Fig. 5 shows that all detected events, regardless of the confidence level, are transmitted to the BS, ensuring a higher number of detections but potentially leading to incorrect detections due to the pervasive nature of the environment. This results in security threats and increased power consumption.

Figure 6 shows how the proposed model optimizes energy usage by reducing transmission overhead. The comparison of energy consumption shown in Fig. 6 further validates the effectiveness of the proposed model in optimizing power usage within a PAWSN framework. The energy consumption is calculated on the basis of the parameter values provided in Tab. 2. The figure illustrates that, for this particular instance of listening to a single event, the proposed model consumes only 3.08 J of energy, while a model without TinyML-based or context-aware filtering expends 5.19 J. This specific instance of PAWSN operation demonstrates a significant energy efficiency improvement of 40.56%, confirming that filtering low-confidence events effectively reduces unnecessary processing and transmission costs.

These results highlight the trade-off between comprehensive event detection and energy conservation, reinforcing the advantages of the proposed model for resource-constrained WSNs.

**Tab. 1.** Minimum hardware requirements for running TinyML and standard ML models.

Model	CPU	Memory	Power	Seed XIAO ESP32S3 specification
TinyML models (~200 KB)				
Tiny-DNN	240 MHz	300 – 500 KB	0.25 – 0.35 W	240 MHz, 8 MB PSRAM, low power MCU
Tiny-KNN	240 MHz	300 KB	0.30 – 0.40 W	Suitable for small dataset classification
Tiny-RF	240 MHz	400 – 600 KB	0.35 – 0.45 W	Works with small tree depth
Tiny-SVM	240 MHz	200 – 300 KB	0.30 – 0.40 W	Efficient with linear kernel
Standard ML models (for PC/GPU)				
DNN	2 GHz	2 GB	50 W	Not supported
KNN	2 GHz	2 GB	60 W	Not supported
RF	3.5 GHz	4 GB	70 W	Not supported
SVM	2 GHz	2 GB	60 W	Not supported

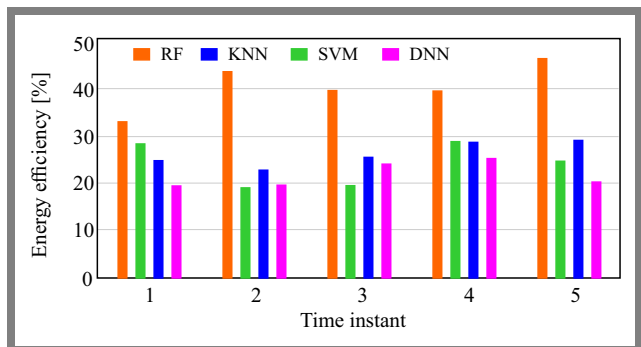
**Tab. 2.** Hardware specifications and timing parameters.

Parameter	Description	Value
Hardware specifications		
$P_{SENSE}$	Power consumption for sensing (Grove acoustic sensor)	0.035 W
$P_{PROCESS}$	Power consumption for processing (XIAO ESP32S3 )	0.35 W
$P_{TX}$	Power consumption for transmission (XBee)	1.25 W
$P_{RX}$	Power consumption for receiving (XBee)	1.2 W
Timing parameters		
$T_{SENSE}$	Sensing time (duration of audio file)	5.0 s
$T_{PROCESS\_SENSOR}$	Processing time (XIAO ESP32S3 )	0.2 s
$T_{PROCESS\_ROUTER}$	Processing time (router node – XBee)	0.1 s
$T_{TX}$	Transmission time (XBee, 431 KB file)	0.275 s
$T_{RX}$	Receiving time (XBee, 431 KB file)	0.275 s

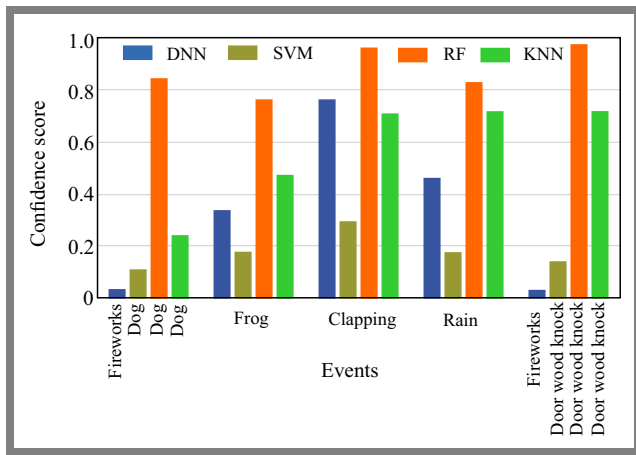
## 5.2. Performance Evaluation and Comparative Analysis

Figure 7 presents a comparative analysis of five consecutive event listening instances in four pervasive sensor zones, focusing on the detection of four specific acoustic events with a detection accuracy rate exceeding 0.7. The analysis includes four different TinyML models and a baseline scenario without TinyML, where all events are detected and forwarded to the base station. In all cases, the figure demonstrates energy efficiency gains while maintaining a good event detection accuracy rate. Furthermore, it compares energy efficiency across different TinyML models to identify the most optimal model for the ESC-50 dataset. The results indicate that the RF algorithm achieves the best performance, with an energy efficiency of approximately 40%.

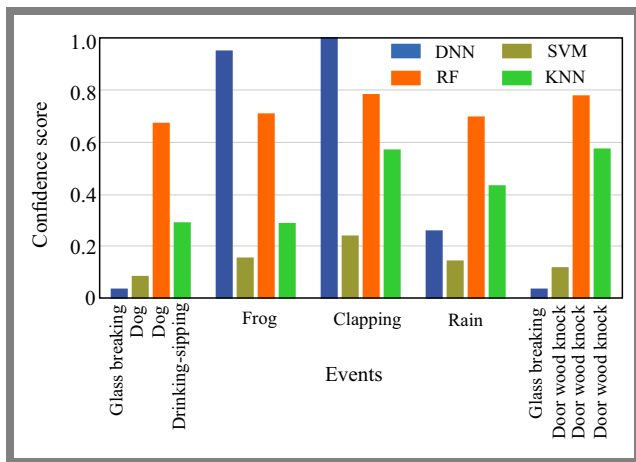
Performance evaluation and comparative analysis were also performed to assess the effectiveness of different TinyML models deployed at sensor nodes and to compare it with the effectiveness of standard models running on high-performance computers. Figures 8 and 9 illustrate that TinyML models, despite operating in resource-constrained environments, achieve performance levels comparable to those of standard models


**Fig. 7.** Energy efficiency of different TinyML models on Seed XIAO ESP32S3.

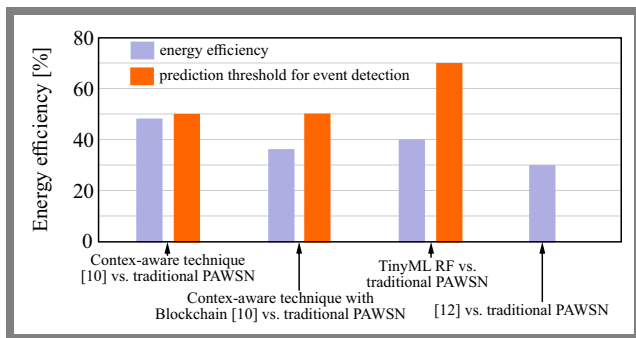
used in advanced systems. The evaluation considered multiple classifications of acoustic events, with models such as SVM, RF, KNN and DNN being tested at the sensor nodes. The results demonstrate that among the TinyML models evaluated, the RF model exhibits superior performance in terms of classification accuracy, while simultaneously maintaining good energy efficiency. Based on these findings, the RF model is recommended for implementation at the sensor nodes, as



**Fig. 8.** Event detection accuracy for different TinyML models on Seede XIAO ESP32S3i.



**Fig. 9.** Event detection accuracy for different standard ML models at the base station.



**Fig. 10.** Comparison of energy efficiency with context-aware model.

it ensures reliable and accurate event detection with minimal computational overhead.

Performance is also compared with the findings of the work on blockchain-based acoustic data integration for PAWSN [10], where the context-aware event detection technique was used to reduce transmission cost by decreasing the number of receiving and transmitting operations. This technique replaces these context-sensitive semantic sensor nodes with TinyML-based AED sensor nodes, enabling more accurate detection of acoustic events while simultaneously reducing the significant amount of acoustic data transmission, and thus contributing

to greater energy efficiency. Figure 10 shows that energy efficiency performance of the two models remains almost the same, while the event prediction threshold is set to 0.5 in [10] and 0.7 in the present model, highlighting the more accurate event detection process of the proposed model. The comparison also includes a blockchain-enhanced context-aware technique [12] and a traditional WSN-based approach, further demonstrating the trade-offs between energy efficiency and event detection accuracy.

## 6. Conclusion and Future Scope

This work demonstrates the effectiveness of TinyML-driven sensor nodes for energy-efficient acoustic event detection in pervasive wireless sensor networks (PAWSNs). By relying on lightweight machine learning models, the proposed approach significantly reduces data transmission, optimizes energy consumption and maintains a high classification accuracy level.

Comparative analysis confirms that TinyML models achieve performance levels comparable to those of standard high-performance computing models, making them a viable alternative for real-time edge computing applications. Among the models evaluated, the random forest algorithm is efficient, achieving approximately 40% energy savings while ensuring reliable event detection.

The findings highlight the potential of TinyML to improve the longevity and sustainability of sensor networks deployed in resource-constrained environments. Future work can improve TinyML adaptability with real-time updates and self-learning while ensuring scalability across industries, wildlife, and smart cities, with energy harvesting deployed to extend network longevity.

## References

- [1] S. Das and U. Mondal, "Acoustic Data Acquisition and Integration for Semantic Organization of Sentimental Data and Analysis in a PWSN", *Multimedia Tools and Applications*, 2024 (<https://doi.org/10.1007/s11042-024-20229-4>).
- [2] S. Das and U. Mondal, "Pilot Agent Implied Efficient Data Communication in Pervasive Acoustic Wireless Sensor Network", *Telecommunication Systems*, vol. 88, art. no. 50, 2025 (<https://doi.org/10.1007/s11235-025-01281-3>).
- [3] P. Andrade *et al.*, "A TinyML Soft-sensor Approach for Low-cost Detection and Monitoring of Vehicular Emissions", *Sensors*, vol. 22, art. no. 3838, 2022 (<https://doi.org/10.3390/s22103838>).
- [4] P. Yadav, "Advancements in Machine Learning in Sensor Systems: Insights from Sensys-ML and TinyML Communities", *2024 IEEE 3rd Workshop on Machine Learning on Edge in Sensor Systems (SenSys-ML)*, Hong-Kong, 2024 (<https://doi.org/10.1109/SenSys-ML62579.2024.00009>).
- [5] A. Sabovic *et al.*, "Towards Energy-aware TinyML on Battery-less IoT Devices", *Internet of Things*, vol. 22, art. no. 100736, 2023 (<https://doi.org/10.1016/j.iot.2023.100736>).
- [6] H.H. Ahmed, Z. Ahmed, T. Choden, and N. Chaudhary, "TinyML for Emotion Detection in Voice Signals: Evaluating and Proposing Algorithms for IoT Wearable Devices", *Thesis*, Brac University, 2024 [Online] Available: <http://hdl.handle.net/10361/24346>.
- [7] S. Hammad, D. Iskandaryan, and S. Trilles, "An Unsupervised TinyML Approach Applied to the Detection of Urban Noise Anomalies un-

- der the Smart Cities Environment”, *Internet of Things*, vol. 23, art. no. 100848, 2023 (<https://doi.org/10.1016/j.iot.2023.100848>).
- [8] Z. Huang *et al.*, “TinyChirp: Bird Song Recognition Using TinyML Models on Low-power Wireless Acoustic Sensors”, *2024 IEEE 5th International Symposium on The Internet of Sounds (IS2)*, Erlangen, Germany, 2024 (<https://doi.org/10.1109/IS262782.2024.10704131>).
- [9] A. Elhanashi, P. Dini, S. Saponara, and Q. Zheng, “Advancements in TinyML: Applications, Limitations, and Impact on IoT Devices”, *Electronics*, vol. 13, art. no. 3562, 2024 (<https://doi.org/10.3390/electronics13173562>).
- [10] S. Das and U. Mondal, “Energy Efficient Acoustic Sensor Data Integration in Hybrid Mode Operated Pervasive Wireless Sensor Network”, *Telecommunication Systems*, vol. 87, pp. 61–72, 2024 (<https://doi.org/10.1007/s11235-024-01165-y>).
- [11] G. Cerutti, R. Prasad, A. Brutti, and E. Farella, “Compact Recurrent Neural Networks for Acoustic Event Detection on Low-energy Low-complexity Platforms”, *IEEE Journal of Selected Topics in Signal Processing*, vol. 14, pp. 654–664, 2020 (<https://doi.org/10.1109/JSTSP.2020.2969775>).
- [12] M. Andhare *et al.*, “Design and Implementation of Wireless Sensor Network for Environmental Monitoring”, *International Journal of Health Sciences*, vol. 6, pp. 3158–3169, 2022 (<https://doi.org/10.53730/ijhs.v6nS4.9085>).
- [13] G. Cerutti *et al.*, “Sound Event Detection with Binary Neural Networks on Tightly Power-constrained IoT Devices”, *Proc. of the ACM/IEEE International Symposium on Low Power Electronics and Design*, pp. 19–24, 2020 (<https://doi.org/10.1145/3370748.3406588>).
- [14] M. Antonini, M. Pincheira, M. Vecchio, and F. Antonelli, “An Adaptable and Unsupervised TinyML Anomaly Detection System for Extreme Industrial Environments”, *Sensors*, vol. 23, art. no. 2344, 2023 (<https://doi.org/10.3390/s23042344>).
- [15] S. Githu, “Detecting Worker Accidents with Audio Classification – Syntiant TinyML”, Edge Impulse, [Online] Available: <https://docs.edgeimpulse.com/experts/audio-projects/detecting-worker-accidents-syntiant-tinyml>.
- [16] Q. Wu, P. Sun, and A. Boukerche, “An Energy-efficient UAV-based Data Aggregation Protocol in Wireless Sensor Networks”, *Proc. of the 8th ACM Symposium on Design and Analysis of Intelligent Vehicular Networks and Application*, pp. 34–40, 2018 (<https://doi.org/10.1145/3272036.3272047>).
- [17] J. Yan, X. Yang, X. Luo, and C. Chen, “Energy-efficient Data Collection over AUV-assisted Underwater Acoustic Sensor Network”, *IEEE Systems Journal*, vol. 12, pp. 3519–3530, 2018 (<https://doi.org/10.1109/JSYST.2017.2789283>).
- [18] Z. Zhou, S. Zhou, J.-H. Cui, and S. Cui, “Energy-efficient Cooperative Communication Based on Power Control and Selective Single-relay in Wireless Sensor Networks”, *IEEE Transactions on Wireless Communications*, vol. 7, pp. 3066–3078, 2008 (<https://doi.org/10.1109/TWC.2008.061097>).
- [19] C. Chen *et al.*, “Ubiquitous Monitoring for Industrial Cyber-physical Systems over Relay-assisted Wireless Sensor Networks”, *IEEE Transactions on Emerging Topics in Computing*, vol. 3, pp. 352–362, 2015 (<https://doi.org/10.1109/TETC.2014.2386615>).
- [20] S. Heydari and Q.H. Mahmoud, “Tiny Machine Learning and On-device Inference: A Survey of Applications, Challenges, and Future Directions”, *Sensors*, vol. 25, art. no. 3191, 2025 (<https://doi.org/10.3390/s25103191>).
- [21] L. Capogrosso *et al.*, “A Machine Learning-oriented Survey on Tiny Machine Learning”, *IEEE Access*, vol. 12, pp. 23406–23426, 2024 (<https://doi.org/10.1109/ACCESS.2024.3365349>).

---

### Bibek B. Roy, M.Sc.

Department of Computer Science

 <https://orcid.org/0009-0001-6699-1551>

E-mail: roybibek16@gmail.com

Vidyasagar University, Midnapore, WB, India

<https://www.vidyasagar.ac.in>

### Sushovan Das, Ph.D.

Department of CSE

 <https://orcid.org/0000-0003-2759-3902>

E-mail: das.sushovan@gmail.com

College of Engineering & Management, Kolaghat, WB, India

<https://www.cemkolaghat.in>

### Uttam Kr. Mondal, Ph.D.

Department of Computer Science

 <https://orcid.org/0000-0002-7807-3002>

E-mail: uttam\_ku\_82@yahoo.co.in

Vidyasagar University, Midnapore, WB, India

<https://www.vidyasagar.ac.in>

# Reconfigurable Reflectarray Structure Based on Optimized Unit Cell for Wireless Communications

Reham Mahmood Yaseen and Ali Khalid Jassim

Mustansiriyah University, Baghdad, Iraq

<https://doi.org/10.26636/jtit.2025.2.2104>

**Abstract** — This paper presents a  $180 \times 180 \times 1$  mm reconfigurable reflector array structure based on an optimized unit cell for wireless communication applications. The reflector array contains 144 unit cells placed on the FR4 substrate, and each unit cell structure uses a single layer based on multi-concentric square rings. The single layer is used to obtain negative  $\epsilon_r$  and  $\mu_r$  values, while multiple rings provide a wide reflection bandwidth. The proposed structure is characterized by dual reflection bandwidth. The first band (2.6 GHz) ranges from 1.98 GHz to 4.6 GHz, while the other band (1.71 GHz), ranging from 7.41 GHz to 9.1 GHz. The reconfigurability of the structure is realized by using PIN diodes connected to each unit cell. Phase distribution in the proposed reflector structure changes according to state of the diodes, resulting in the reflection of the wave at different angles. The proposed solution was simulated in terms of S parameters, constitutive parameters and refractive index based on a full-wave analysis performed using CST Microwave Studio.

**Keywords** — constitutive parameters, metamaterial, reconfigurable reflectarray, unit cell

## 1. Introduction

Reflectarray structures are currently receiving a great deal of interest when it comes to their use in wireless systems, due to their feature-rich capabilities. Compared with traditional reflectors, reflectarrays are characterized by light weight, low profile, flat surface and small volume. They are also easy and cheap to fabricate [1], [2]. Such features make reflectarray structures a preferred solution for satellite, radar, and long-distance wireless communication applications [3].

Recently, reflectarray structures with unit cells of a new geometrical shape have been developed to overcome their inherent disadvantage consisting in a limited phase variation range [4]. In [5], a reflector structure based on nine elements of different sizes is introduced for the purpose of reducing RCS level, while in [6] a reflectarray structure based on a stepped impedance resonator (SIR) supported on a frequency selective surface (FSS) is proposed for the purpose of lowering its radar cross-section. A compact reflectarray structure is introduced in [7] for performance enhancement, with the required phase delay compensated for by variable-size unit cell elements. Such elements are used to reflect/scatter the incident waves

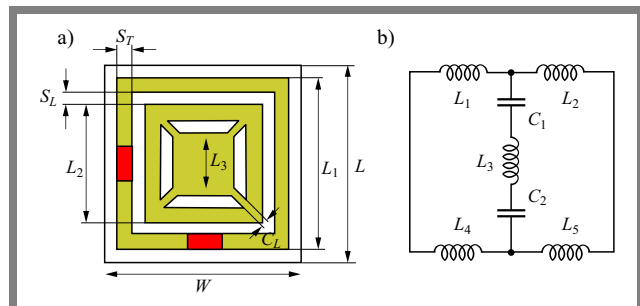
within the phase delay to achieve a desired shape of the beam, as presented in papers [8]–[12].

In this work, a reflectarray structure based on a single negative metamaterial, designed to achieve a reconfigurable behavior, is presented. Metamaterials (MTMs) are substances based on properties not found in nature. Their parameters – such as  $\mu_r$  and  $\epsilon_r$  – depend on the structure, shape, arrangement, orientation, and size of the unit cell, and may also assume negative values [13]–[16]. Materials with negative  $\mu_r$  and  $\epsilon_r$  parameters are called left-handed materials (LH MTM). Their phase velocity is antiparallel to the Poynting vector and opposite to the propagation of waves in natural materials [17]–[20].

In this design, the unit cell utilizes a single layer-based, multi-concentric square-ring structure. A single layer is used to achieve a negative  $\epsilon_r$  or  $\mu_r$ , while multiple rings are used to provide wide reflective bands for sub-6 GHz band 5G applications. First, the reflectarray size is designed and optimized using the trust region algorithm. Then, a full-wave analysis is conducted and the obtained results are discussed. The CST Microwave software package is used to design and analyze the proposed structures [21].

## 2. Unit Cell Design

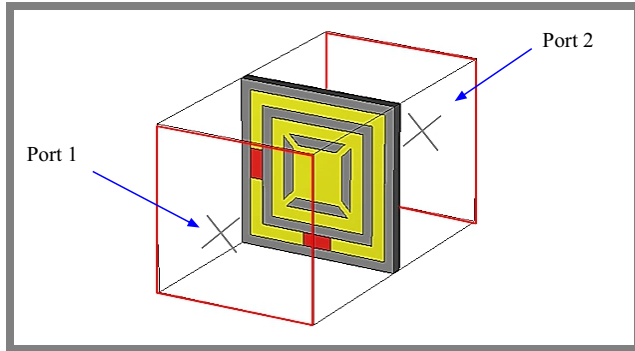
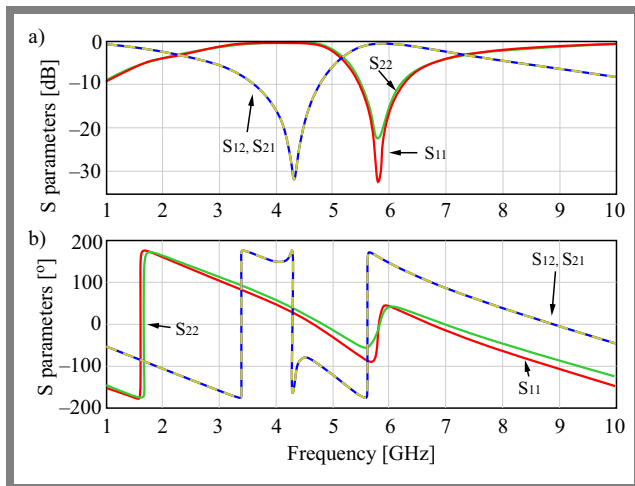
The optimized single negative metamaterial unit cell based on the reflectarray behavior is shown in Fig. 1. The dimensions of the unit cell are  $15 \times 15 \times 1$  mm and as the base substrate, the FR-4 material with  $\epsilon_r = 4.3$  and  $\tan \delta = 0.02$  is used. For reconfigurable behavior, two PIN diodes are inserted at



**Fig. 1.** Reconfigurable metamaterial structure: a) unit cell design and b) equivalent circuit.

**Tab. 1.** Dimensions of the unit cell.

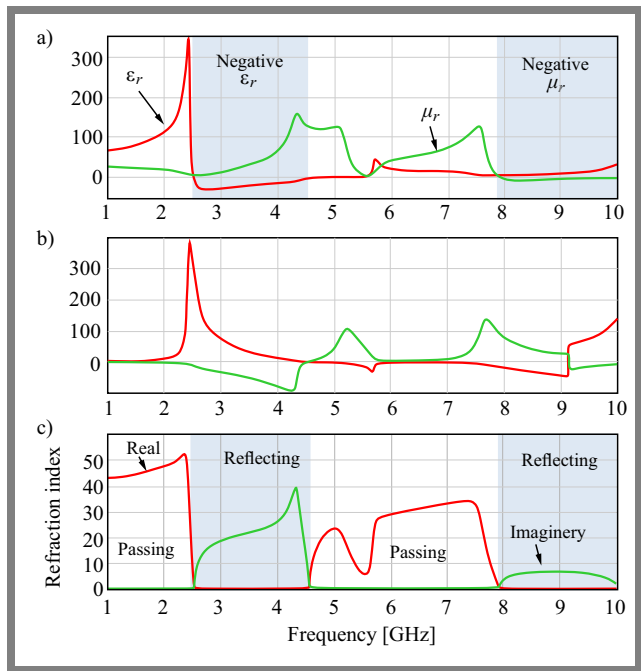
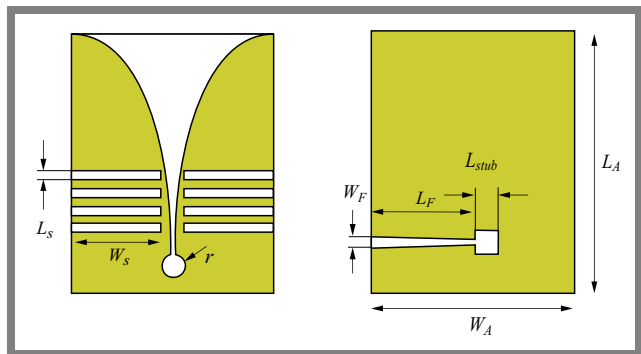
Parameter	Value	Parameter	Value
$L$	15 mm	$L_3$	4.8 mm
$W$	15 mm	$S_T$	1.19 mm
$L_1$	13.2 mm	$S_L$	0.9 mm
$L_2$	13.2 mm	$C_L$	0.636 mm
$h$	1 mm	$t$	0.035 mm


**Fig. 2.** Fictitious waveguide with port excitation entries and boundary conditions settings.

**Fig. 3.** S parameters of the unit cell: a) magnitude and b) phase.

each unit cell, with a rectangular strip used as the PIN diode switches between on and off states. The equivalent circuit of the proposed unit cell is presented in Fig. 1b, while other geometric details are shown in Tab. 1.

A full-wave analysis technique was applied to analyze the single unit cell. As illustrated in Fig. 2, the unit cell structure is placed in the center of a fictitious waveguide. To mimic the wave's propagation through an infinite unit cell array, two ports are inserted, in addition to the boundary conditions provided by perfect electrical conductors (PECs) and perfect magnetic conductors (PMCs). In Fig. 3, the S parameters of such a unit cell are introduced in a set of magnitude and phase curves vs. frequency.

Figure 4 presents the constitutive parameters ( $\mu_r$  and  $\epsilon_r$ ) of the proposed unit cell. One may notice that it has negative permittivity and permeability in multifrequency bands. A


**Fig. 4.** Constitutive parameters ( $\epsilon_r$ ,  $\mu_r$ ): a) real, b) imaginary component, and c) refraction index.

**Fig. 5.** Geometry of the MVA structure.

negative  $\epsilon_r$  has been achieved in the range of 2.51 to 4.56 GHz, with its maximum value reaching 2.73 GHz. A negative  $\mu_r$  has been achieved in the frequency range of 7.91 to 10 GHz, with a maximum value of 8.25 GHz. Figure 4c illustrates the refraction index  $n = \sqrt{\epsilon_r \mu_r}$  with passing and reflecting frequency ranges of the designed unit cell.

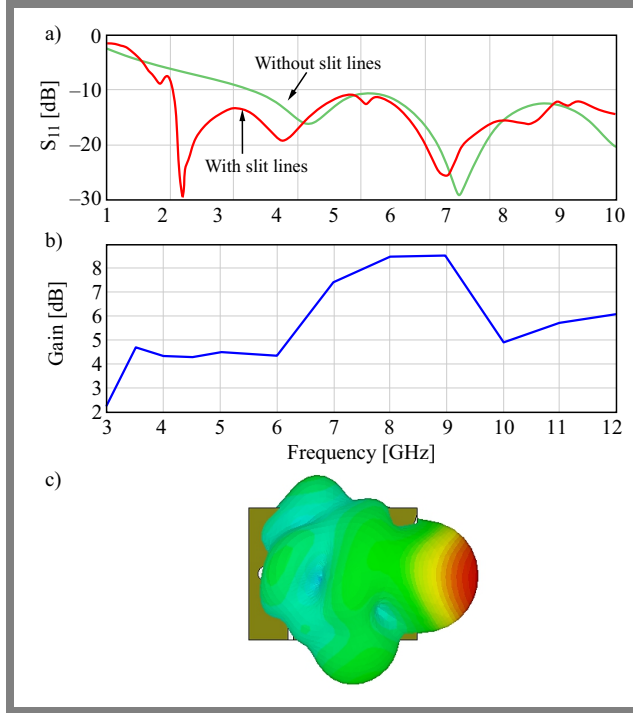
### 2.1. Modified Antenna Design

In the next step, a modified Vivaldi antenna (MVA) was designed to generate the electromagnetic waves that will be incident on the reflectarray structure. The dimensions of MVA are  $45 \times 35$  mm with a substrate thickness of 1 mm. Slit lines are etched along the outer edges to reduce surface current values. As a result, the antenna's performance is enhanced. Geometric details of the MVA structure may be found in Tab. 2, while Fig. 5 shows the proposed layout.

The  $S_{11}$  of the designed antenna (Fig. 6) shows that the antenna covers a range of frequencies from 3 to 12 GHz with  $S_{11} < -10$  dB, as seen in Fig. 6a. Furthermore, the MVA antenna shows endfire radiation with its maximum gain

**Tab. 2.** Dimensions of the unit cell.

Parameter	Value	Parameter	Value
$L_A$	45 mm	$L_{stub}$	4 mm
$W_A$	35 mm	$L_S$	1.5 mm
$L_F$	17.86 mm	$W_S$	15.5 mm
$W_F$	1.945 mm	$r$	2 mm



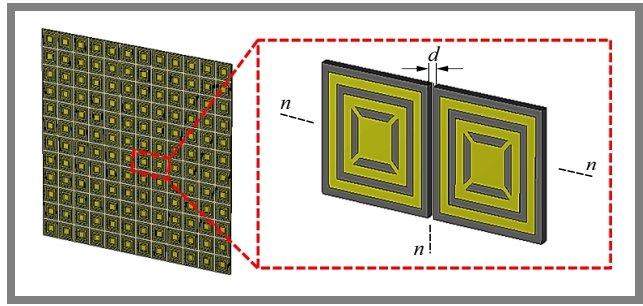
**Fig. 6.** Performance of the MVA antenna: a) reflection coefficient, b) gain value at different operating frequencies, and c) radiating pattern.

equal to 8.55 dB, varying from 4.29 to 8.55 dB in the desired band, as seen in Fig. 6b-c.

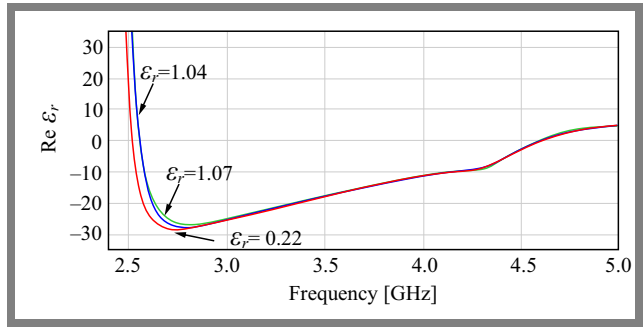
**2.2. Reconfigurable Reflectarray Based on Optimized Unit Cell**

In the third step, an array of optimized unit cells was designed, as illustrated in Fig. 7. It was of the reflectarray variety and consisted of  $12 \times 12$  unit cells, with an overall size of  $180 \times 180 \times 1$  mm. The dimensions were optimized based on separation distance  $d$ . A trust region algorithm was applied and used to achieve the best result. As seen in Fig. 8, the best separation distance  $d$  was found to be equal to 0 mm in terms of a wider reflected band with the minimum size of the array compared to the other distance.

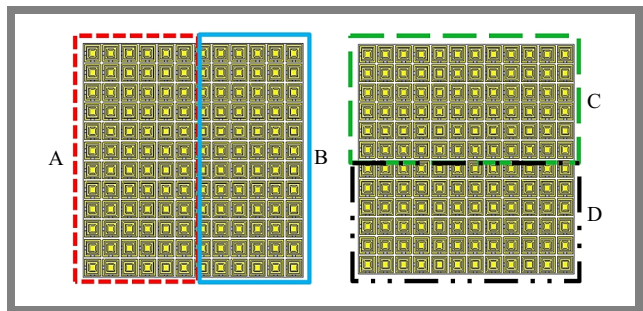
The proposed reflectarray solution was then reconfigured, with the unit placed parallel to the direction of propagated waves in order to prevent the propagation of electromagnetic waves through the proposed reflectarray structure, which has led to the most efficient reflection of radiation power. The reconfigurable unit cells based on the PIN diode in the reflectarray structure were divided into 4-groups, as seen in Fig. 9, with each group



**Fig. 7.** Optimized unit cell-based array.



**Fig. 8.** Real  $\epsilon_r$  component at different separation distances.

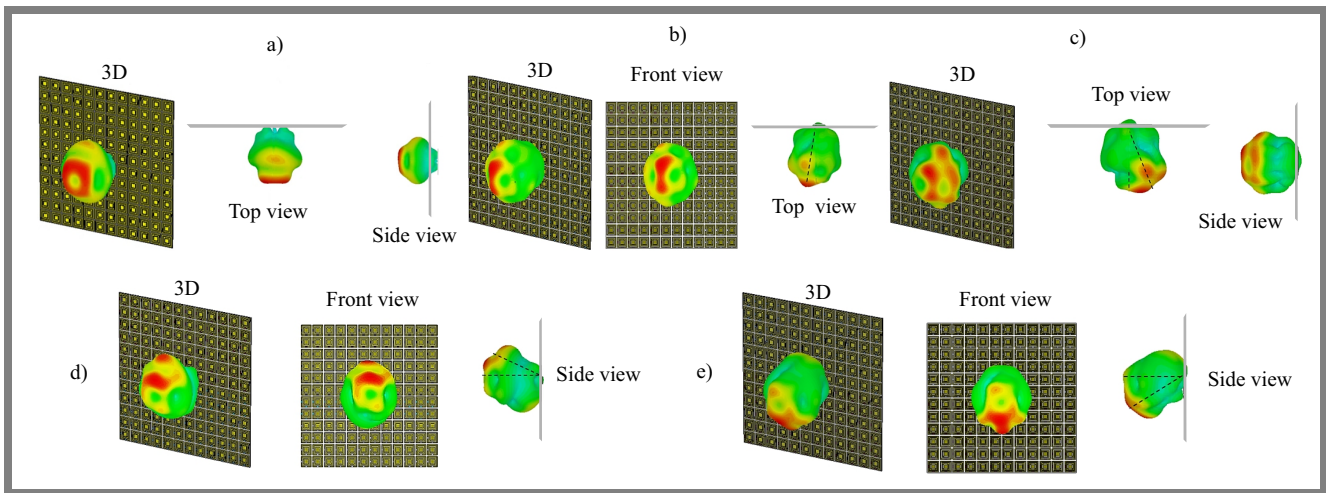


**Fig. 9.** Proposed final reflectarray structure.

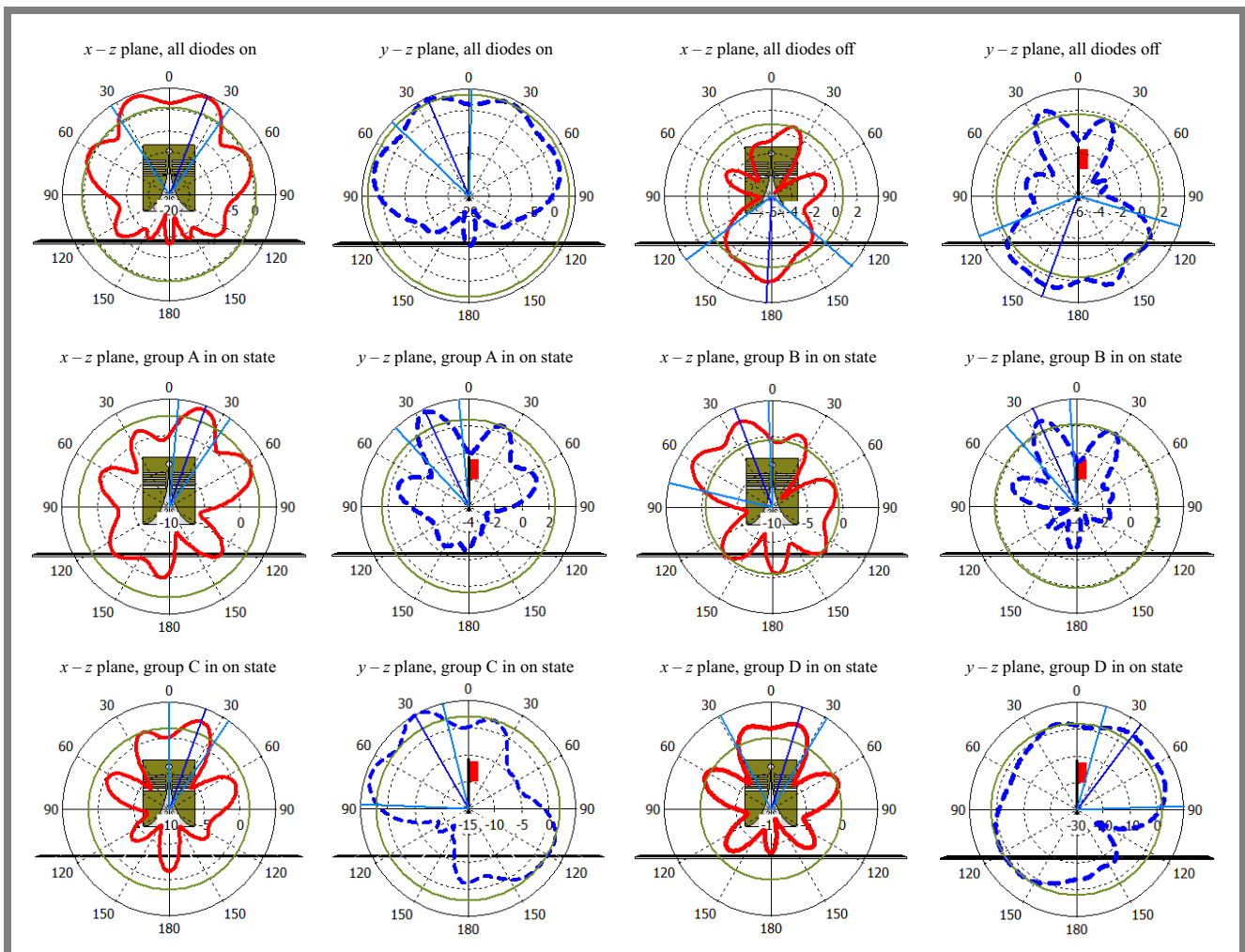
of unit cells being responsible for different reflection direction angles.

In the “on” state, the radiation pattern of the proposed structure is presented in Fig. 10a. One may notice that the incident waves will be reflected from the reflectarray and no waves will be passed through the array.

The radiation pattern of the proposed array, where the diodes are on in group A, is introduced in Fig. 10b. It is found that the incident wave is reflected from the reflectarray structure at an angle of  $20^\circ$  to the left-hand side. With diodes in group B switched on (Fig. 10c), it is noticed that the incident waves are reflected to the right-hand side, at an angle of  $21^\circ$ . Moreover, when group C is activated, as seen in Fig. 10d, the waves are reflected at an angle of  $30^\circ$  upwards. The reflected waves are redirected at an angle of  $37^\circ$  downwards when group D diodes are on, as seen in Fig. 10e. Depending on which group of diodes is switched on, the waves will be reflected in different directions and at different angles. Finally, 2D radiation patterns are introduced in Fig. 11, for different states of the diodes.



**Fig. 10.** Radiation pattern of the proposed structure in the case of: a) all diodes in on state, b) group A being in the on state, c) group B being in the on state, d) group C being in the on state, and e) group D being in the on state.



**Fig. 11.** 2D projection of radiation patterns of the proposed structure for different states of the diodes.

### 3. Conclusions

In this paper, an optimized reconfigurable reflectarray structure based on single negative unit cell metamaterials is proposed. A single negative metamaterial is used to reflect the waves. For

reconfigurable behavior, two PIN diodes were used in the unit cell to manipulate its constitutive parameters ( $\epsilon_r$ ,  $\mu_r$ ). A trust region algorithm was employed to optimize the reflectarray size area. By switching the PIN diodes, various reflection angles of the incident wave were achieved.

## References

- [1] J. Huang and J.A. Encinar, *Reflectarray Antennas*, New York: Wiley, 216 p., 2007 (<https://doi.org/10.1002/9780470178775>).
- [2] S. Oh, "Broadband Reflectarray Composed of Gap Coupled Elements with Linear Phase Response", *Microwave and Optical Technology Letters*, vol. 59, pp. 1045–1047, 2017 (<https://doi.org/10.1002/mop.30464>).
- [3] D.M. Pozar, S.D. Targonski, and H.D. Syrigos, "Design of Millimeter Wave Microstrip Reflectarrays", *IEEE Transactions on Antennas and Propagation*, vol. 45, pp. 287–295, 1997 (<https://doi.org/10.1109/8.560348>).
- [4] M. Karimupour and N. Komjani, "Bandwidth Enhancement of Electrically Large Shaped-beam Reflectarray by Modifying the Shape and Phase Distribution of Reflective Surface", *AEU – International Journal of Electronics and Communications*, vol. 70, pp. 530–538, 2016 (<https://doi.org/10.1016/j.aeue.2015.12.007>).
- [5] H. Bodur, S. Unald, S. Cimen, and G. Cakir, "A Novel Reflectarray Antenna with Reduced RCS", *Kocaeli Journal of Science and Engineering*, vol. 1, pp. 11–14, 2018 (<https://doi.org/10.34088/kojose.399389>).
- [6] M.M. Fakharian, P. Rezaei, and A.A. Orouji, "A Reflectarray Based on the Folded SIR Patch-slot Configuration Backed on FSS for Low RCS", *Progress in Electromagnetics Research Letters*, vol. 47, p. 119–124, 2014 (<https://doi.org/10.2528/PIERL14061803>).
- [7] H. Bodur and S. Cimen, "Reflectarray Antenna Design with Double Cutted Ring Element for X-band Applications", *Microwave and Optical Technology Letters*, vol. 62, pp. 3248–3254, 2020 (<https://doi.org/10.1002/mop.32436>).
- [8] F. Xue, H.-J. Wang, M. Yi, and G. Liu, "A Broadband Ku-band Microstrip Reflectarray Antenna Using Single-layer Fractal Elements", *Microwave and Optical Technology Letters*, vol. 58, pp. 658–662, 2016 (<https://doi.org/10.1002/mop.29637>).
- [9] K. Lele, A.A. Desai, A.A. Kadam, and A.A. Deshmukh, "Reflectarray Antennas", *International Journal of Computer Applications*, vol. 3, pp. 21–28, 2014 (<https://doi.org/10.5120/18891-0173>).
- [10] A.R. Azeez *et al.*, "UWB Tapered-slot Patch Antenna with Reconfigurable Dual Band-notches Characteristics", *Journal of Mechanics of Continua and Mathematical Sciences*, vol. 19, pp.40–50, 2024 (<https://doi.org/10.26782/jmcms.2024.03.00003>).
- [11] Y.S. Mezaal *et al.*, "State of Art on Microstrip Resonators, Filters, Diplexers and Triplexers", *Journal of Mechanics of Continua and Mathematical Sciences*, vol. 19, pp. 1–24, 2024 (<https://doi.org/10.26782/jmcms.2024.02.00001>).
- [12] A.R. Azeez *et al.*, "High Gain Tapered Slot Antenna Based on Offset Radiation Characteristic for 5G Wireless Applications", *IEEE, 2024 4th International Conference on Artificial Intelligence and Signal Processing (AISP)*, Vijayawada, India, 2024 (<https://doi.org/10.1109/AISP61711.2024.10870648>).
- [13] O. Luukkonen, S.I. Maslovski, and S.A. Tretyakov, "A Stepwise Nicolson-Ross-Weir-based Material Parameter Extraction Method", *IEEE Antennas and Wireless Propagation Letters*, vol. 10, pp. 1295–1298, 2011 (<https://doi.org/10.1109/LAWP.2011.2175897>).
- [14] H.A. Al-Tayyar and Y.E. Mohammed Ali, "Parameters Extraction of Miniaturized Metamaterial Unit Cell at Millimeter Wave Applications", *2023 5th International Congress on Human-Computer Interaction, Optimization and Robotic Applications (HORA)*, Istanbul, Turkiye, 2023 (<https://doi.org/10.1109/HORA58378.2023.10156671>).
- [15] A.R. Azeez *et al.*, "Design of High Gain UWB Vivaldi Antenna with Dual Band-notches Characteristics", *International Journal on Engineering Applications (IREA)*, vol.11, pp. 128–136, 2023 (<https://doi.org/10.15866/irea.v11i2.22177>).
- [16] H.A. Al-Tayyar and Y.E. Mohammed Ali, "A Review on Metamaterial Used in Antennas Design: Advantages and Challenges", *Al-Rafidain Engineering Journal (AREJ)*, vol. 29, pp. 106–117, 2024 (<https://doi.org/10.33899/rengj.2023.140769.1259>).
- [17] N.P. Johnson *et al.*, "Characterization at Infrared Wavelengths of Metamaterials Formed by Thin-film Metallic Split-ring Resonator Arrays on Silicon", *Electronics Letters*, vol. 42, pp. 1117–1119, 2006 (<https://doi.org/10.1049/el:20062212>).
- [18] B. Sauviac, C.R. Simovski, and S.A. Tretyakov, "Double Split-ring Resonators: Analytical Modeling and Numerical Simulations", *Electromagnetics*, vol. 24, no. 5, pp. 317–338, 2004 (<https://doi.org/10.1080/02726340490457890>).
- [19] A.K. Hadi, M.F. Mosleh, A.R. Azeez, and R.A. Abd-Alhameed, "Meander-line Slots Based Miniaturized Inset-fed Patch-sensor Antenna for Water Characterization Technique", *2023 IEEE 13th International Conference on System Engineering and Technology (ICSET)*, Shah Alam, Malaysia, 2023 (<https://doi.org/10.1109/ICSET59111.2023.10295164>).
- [20] S.H. Ali and A.K. Jassim, "Single Layer Metamaterial Superstrate for Gain Enhancement of a Microstrip Antenna Array", *Diyala Journal of Engineering Sciences*, vol. 17, pp. 144–172, 2024 (<https://doi.org/10.24237/djes.2024.17211>).
- [21] Dassault Systemes, CST Studio Suite, version 2019 (<https://www.3ds.com/products/simulia/cst-studio-suite>).

**Reham Mahmood Yaseen, M.Sc.**

Department of Electrical Engineering

 <https://orcid.org/0009-0002-1789-842x>

E-mail: Reham-Mahmood@uomustansiriyah.edu.iq

Mustansiriya University, Baghdad, Iraq

<https://www.uomustansiriyah.edu.iq>**Ali Khalid Jassim, Ph.D.**

Department of Electrical Engineering

 <https://orcid.org/0000-0002-4146-4536>

E-mail: alijassim@uomustansiriyah.edu.iq

Mustansiriya University, Baghdad, Iraq

<https://www.uomustansiriyah.edu.iq>

# Advancing Facial Expression Recognition – Enhanced MobileNetV3 with Integrated Coordinate Attention and Dynamic Kernel Adaptation

Miloud Kamline<sup>1</sup>, Ridha Ilyas Bendjillali<sup>2</sup>, Mohammed Sofiane Bendelhoum<sup>2</sup>,  
Asma Ouardas<sup>2</sup>, and Ali Abderrazak Tadjeddine<sup>2</sup>

<sup>1</sup>Tahri Mohammed University, Bechar, Algeria,  
<sup>2</sup>University Center Nour Bachir, El Bayadh, Algeria

<https://doi.org/10.26636/jtit.2025.2.2146>

**Abstract** — This paper presents an improved approach for facial expression recognition (FER), which incorporates the Coordinate Attention (CAM) mechanism into MobileNetV3, a lightweight CNN widely used for its real-time applications on low-power devices. The CA mechanism greatly improves the ability of the model to focus on face regions of interest, as it incorporates positional information, making feature extraction more accurate. Additionally, dynamic kernel adaptation (DKA) and SoftSwish are incorporated into the model to enhance the flexibility and computational efficiency of MobileNetV3. The proposed model was tested in three sets of JAFFE, CK+, and FER2013, where accuracy improvements were reported of 98.84% in the JAFFE dataset, 99.56% on the CK+ dataset, and 88.50% on the FER2013 dataset. These results support the viability and utility of the proposed approach to improve FER, especially in applications that favor higher numerical performance.

**Keywords** — coordinate attention mechanism, dynamic kernel adaptation, facial expression recognition, MobileNetV3, SoftSwish activation function

## 1. Introduction

Facial expression recognition (FER) has made great progress in recent years, mainly due to the use of neural networks and especially attention mechanisms [1], [2]. These achievements have allowed us to develop accurate, efficient, or even real-time recognition systems, mainly needed in human-computer interaction (HCI), security, and healthcare [3], [4]. A new approach in this direction is MobileNetV3, which is a lightweight CNN optimized for high accuracy with minimal computational resources. MobileNetV3, adopts some of the most modern approaches such as depthwise separable convolutions and a unique activation function called *hard swish* [5]. These elements allow MobileNetV3 to be used in vision systems of portable and embedded devices, which require high power efficiency and speed.

Another significant change in improving the functionality of a neural network is the attention mechanisms. Such mecha-

nisms keep the attention of the network on important parts of an image and enhance the capability of the network to identify features. In this context, attention techniques can be distinguished, such as the latest and most innovative coordinate attention (CA) mechanism [6]. The present proposal of CA also incorporates position information, which other types of attention often isolate and consider individually by location or channel. It allows the network to concentrate more accurately on important sections of the image. This paper concerns an investigation of how to improve FER systems by incorporating CA into MobileNetV3 to achieve a more effective FER system. Specifically, our research seeks to answer the following question: What enhancements are given by the combination of CA with MobileNetV3 for facial expression recognition compared to existing methods?

This research contribution can be summarized in two folds. First, we present an improved FER framework that incorporates MobileNetV3 and the benefits of CA. Second, we give a comparative analysis of this framework on standard FER datasets, including JAFFE, CK+, and FER2013, which illustrates the advantages in terms of performance and time complexity.

The paper is organized as follows. Section 2 presents the background and related work in the field of FER and the attention mechanism. In Section 3, we proposed an approach to MobileNetV3 and the incorporation of CA. In Section 4, we provide experimental results and discuss the efficacy of the proposed method on multiple FER datasets. Section 5 provides an in-depth review of the proposed model. Section 6 presents a comparison with other state-of-the-art FER approaches that utilize attention mechanisms. Lastly, Section 7 provides a conclusion to the paper and future research.

## 2. Related Work

Recent advances in facial expression recognition (FER) have emphasized the integration of attention mechanisms with various neural network architectures to improve accuracy and

precision. Still, these approaches have their problems and have diverse rates according to the specific methodologies and datasets chosen. In [7] a lightweight FER framework was proposed using MobileNetV1 with attention mechanisms. The effectiveness of this approach was observed when tested on CK+, RAF-DB, and FER2013 with evaluation that highlighted the performance when the face images were captured under different lighting conditions or partially occluded.

However, compared with work [8] which used DeeplabV3+ in combination with the MobileNetV2 with attention mechanism, although the method proposed in [7] was much less computationally expensive, it proved to be more generalized, especially in more complex feature extraction tasks. The integration of attention mechanisms in these studies has been found to be useful, but this must be done without overlooking limitations.

For example, the authors of [9] successfully used a multi-attention network to learn discriminative characteristics from important facial areas. However, this method could be sensitive to overfitting, particularly when working with small data like CK+. Also, it is crucial to note that these deemed attention mechanisms have high computational costs that reduce the real-time applicability of attention mechanisms in low-power devices.

This factor is particularly relevant to the study in [10] that demonstrates that, while building a lightweight FER model for mobile devices is a priority, the reduction in computational processes could be detrimental to the high accuracy of feature extraction. It is also important to note that the discussed studies are primarily related to CNNs with attention mechanisms, but there are other promising streams in FER. For instance, graph-based models and transformers are increasingly being adopted in FER tasks because of their capability to handle relationships between facial landmarks and address the temporal aspect in video-based FER tasks.

The absence of these other forms of prediction methods, along with the distinction between model types, including ensemble methods that amalgamate different models to produce a more balanced and accurate model, is a research gap. Extending the study to these angles would give opportunities to study possible developments in FER.

Another important aspect that these works do not normally address concerns the generalizability of FER models to unseen data or different populations. Most of the works cited, such as [11] and [12], offer promising results in popular datasets. However, the data sets could not be rich enough to ensure the transfer of learned representations across a wide diversity of real-world scenarios or across different cultures, age groups, or variability in emotional expressions.

In addition, attention is drawn to the potential biases of these models when trained on limited or homogenized data and ways of correction. A detailed investigation of the type of attention mechanisms applied in these works would be useful to understand which mechanisms are driving performance improvements.

For example, the utilization of self-attention, spatial attention, or channel attention may be important in explaining why models vary in their efficiency. In the case of article [13], combining a ResNet with such an architecture of attention and deformable convolutions, a more detailed breakdown of their interaction could help in understanding their contribution to the model's better accuracy under varying conditions.

Finally, these methods must be compared consistently with thoroughness based on standard evaluation metrics to determine relative performance measures. Although accuracy is commonly reported, other important key measures in the literature include precision, recall, F1 score, and computational efficiency.

A systematic comparison of these metrics would allow a more objective assessment of the strengths and weaknesses of the different models. For example, the authors of [14] are concerned with the computational efficiency of their lightweight facial expression recognition network; it would be beneficial that these requirements were directly compared against accuracy and robustness reports by other models under equivalent constraints.

## 3. Methodology

### 3.1. MobileNetV3

Further developments in computer vision are also driven by the architecture of CNNs that provide, at a time, very high-speed processing while being compact. Examples are architectures such as NASNet [15], MobileNets [16], EfficientNet [17], MnasNet [18], and ShuffleNets [19]. All of these architectures have substantial depth-wise convolution for speeding up training through reduction of computational complexity. In depthwise convolutions, the learned convolution weights are applied to each input channel individually with a shared kernel across all channels, thereby preserving computational resources and reducing overall cost.

However, resolution of the optimal kernel size in such convolutions might be tricky, and it could add complexity in the training phase. Based on the success of MobileNetV1 and MobileNetV2, the authors of [5] recently proposed MobileNetV3 through network architecture search (NAS) with the NetAdapt algorithm to optimize architectures targeting low-resource hardware platforms while balancing size, performance, and latency. This is based on the inverted residual block, which incorporates depth-wise separable convolution and an SE mechanism to improve feature representation while also reducing memory usage.

We further push the capabilities of MobileNetV3 with two key improvements: dynamic kernel adaptation and SoftSwish.

### 3.2. Dynamic Kernel Adaptation (DKA) and Soft Swish

Dynamic kernel adaptation allows the model to dynamically change the kernel size of the convolution according to the particular characteristics of the input data. Unlike the common approach that uses a fixed kernel size, DKA helps the model

```

for feature_map in input:
    complexity = compute_entropy(feature_map)
    attention_weights = softmax(linear_layer(complexity))
    output = 0
    for k in [3, 5, 7]:
        conv_out = depthwise_conv(feature_map, kernel_size=k)
        output += attention_weights[k] * conv_out
    
```

**Fig. 1.** Simplified pseudocode of the feature capture mechanism.

adapt to changing kernel sizes to better handle different image complexities. For example, in complex scenes, DKA can adjust the kernel size to include details, including fine ones, while in simpler contexts, one can get away with a smaller kernel size for efficiency [20]. The flexibility enhances the capability of the model for generalization on different datasets and lessens overfitting risks, therefore making MobileNetV3 powerful, more versatile, especially for application in low resource devices.

We further replace the original activation function with SoftSwish. The transitions of SoftSwish are even more gentle during backpropagation through gradients, making training much more stable and efficient, needed in larger networks [21]. We have defined the SoftSwish function as:

$$SoftSwish(x) = x \frac{1}{1 + e^{-x}} . \tag{1}$$

This function blends the advantages of Swish and ReLU, improving model stability and reducing the number of parameters required during training, thus enhancing overall efficiency. Integrating dynamic kernel adaptation and SoftSwish inside MobileNetV3 makes it much more flexible and efficient in addressing a wide range of visual recognition tasks with much better accuracies while keeping lower computational demands. These enhancements are estimated to increase accuracy by 3–5%, reduce latency by 10–15%, and reduce training time by 5–10%. With minimal training changes, MobileNetV3 improvements will prove to serve as a significant advantage for modern applications, particularly those that demand efficient operation on low-resource hardware platforms.

DKA allows the convolutional kernel size to be dynamically adjusted on the input complexity. Specifically, a lightweight gating mechanism evaluates an entropy-based complexity score  $C(x)$  for each input feature map. A soft-attention function:

$$\alpha_k = \text{softmax}(f_k(C)) , \tag{2}$$

selects between kernel sizes  $k \in \{3, 5, 7\}$  during forward propagation. This mechanism enables the network to capture both local and global features adaptively. A simplified pseudocode is provided in Fig. 1.

### 3.3. MobileNetV3 for Feature Extraction

Feature extraction is one of the important processes in FER to make classifications from an image robust and precise. The outstanding performance of the improved state-of-the-art of MobileNetV3, with the modifications that have been made by us, offers a superb platform for the said activity. We used the

MobileNetV3-Large model pre-trained on the ImageNet for feature extraction. This is because MobileNetV3, in general, has already been known for its efficiency and flexibility, especially when combined with our dynamic kernel adaptation and the SoftSwish activation function – both aspects increase model flexibility and training stability.

We retrain MobileNetV3 for FER, by transfer learning retraining on the already fine-tuned MobileNet, where the original fully connected layers designed for general image classification were replaced. Instead, we introduce a series of  $1 \times 1$  point-wise convolutional layers that further refine the representations in those features that are highly specific to facial expressions. More specifically, these layers can use the adaptively resized kernels of DKA, so the model can change its receptive field with the complexity of the input image. This architecture ensures that the extracted features are relevant and discriminative, considering the unique challenges of the given datasets for facial expressions.

After the pointwise convolutional layers had been implemented, we further embedded the SoftSwish activation function at every layer within the network. In this case, SoftSwish can provide smoother gradient transitions, hence providing better backpropagation efficiency, especially with deeper network layers, which helps the generalization across different datasets with diverse FER. During fine-tuning, we train the model for 160 epochs with ten separate runs, each initiated with random parameters to ensure robustness.

Instead of relying on traditional data augmentation methods, the training process incorporated speckle noise augmentation, random rotation, random zoom, and color jitter augmentation. These techniques are designed to simulate real-world conditions and further enhance the model’s accuracy by allowing it to recognize facial expressions under varying conditions. This extensive training and fine-tuning process allows the model to produce high-quality image embeddings, each represented as a 128-dimensional vector, encapsulating the essential features necessary for accurate FER.

The result is a highly efficient feature extraction process that benefits from the enhanced capabilities of MobileNetV3, making it particularly well suited for deployment in resource-constrained environments where both performance and computational efficiency are paramount (see Fig. 2).

### 3.4. Coordinate Attention Module

The coordinate attention module (CAM) represents a significant advance in attention mechanisms within neural network architectures, particularly by enhancing spatial awareness and focus. Although integrated into models like MobileNetV3, CAM offers substantial improvements in tasks such as facial expression recognition, where precise spatial information is crucial.

Coordinate attention diverges from traditional attention mechanisms by incorporating positional encodings directly into the attention process. For an image  $I$  with dimensions  $W \times H$ , each pixel coordinate  $(x, y)$  contributes to the attention mechanism through a function  $f(x, y)$  that encodes positional

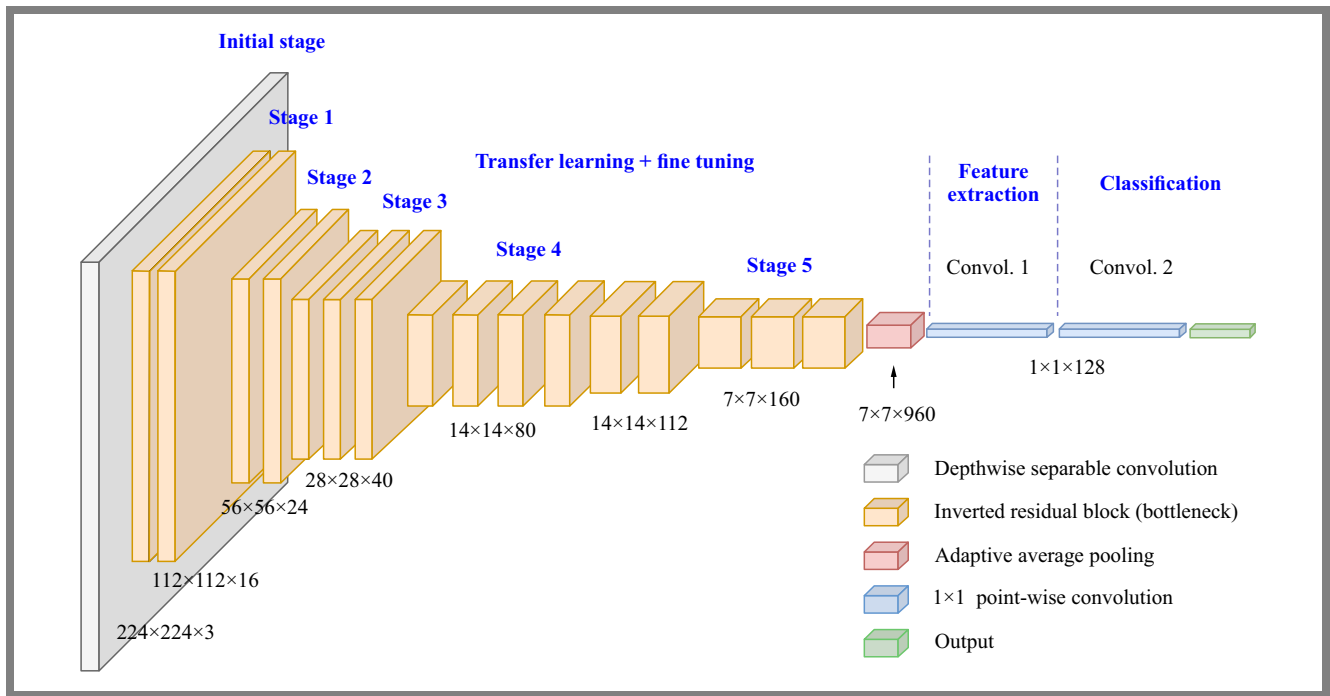


Fig. 2. The architecture of MobileNetV3 used for feature extraction.

information. This can be represented as:

$$Attention(x, y) = f(x, y) . \tag{3}$$

This formulation allows the network to assign attention scores or weights to specific pixel coordinates, highlighting areas of the image that are most relevant for feature extraction. By leveraging these positional encodings, CAM enables the network to focus on critical spatial relationships within the image, thereby enhancing the model’s ability to capture detailed

and contextually relevant features, especially in complex tasks like facial expression recognition.

When applied to MobileNetV3, CAM integrates seamlessly with the existing structure, working in tandem with our proposed DKA and SoftSwish activation function. This combination ensures that the network not only adapts to varying image complexities but also focuses its computational resources on the most significant areas of the image, thereby improving both accuracy and efficiency.

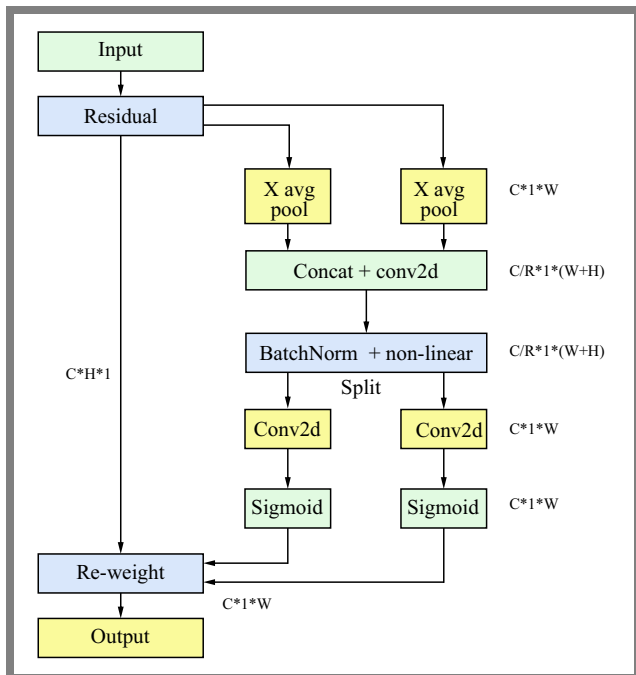
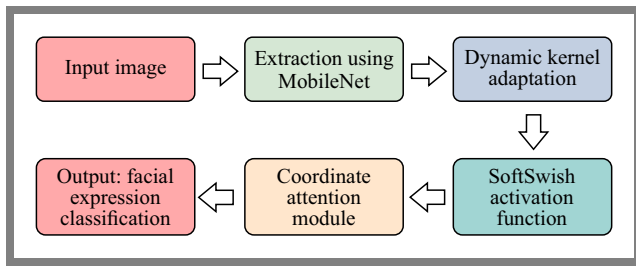


Fig. 3. Integration of the coordinate attention module (CAM) into the MobileNetV3 architecture.

### 3.5. Integration of CAM into MobileNetV3

The MobileNetV3 architecture also integrated the CAM to further improve feature extraction. CAM is integrated with MobileNetV3’s inverted residual blocks, which creates a more explicit level of spatial focus and greatly increases the performance of the model in facial expression recognition. This integration enables MobileNetV3 to make better use of dependencies over space in facial images. The network is more capable of creating a more discriminative and accurate feature representation because the location of the key coordinates on the face in the image is dynamically updated. CAM can selectively emphasize important regions of the face, so that the network pays closer attention to the most informative aspects of facial expressions and neglects the non-critical areas.

The result is a model that takes advantage of both the efficiency and flexibility provided by MobileNetV3, further enhanced by DKA and SoftSwish, in a fashion that gains even more insight into spatial relationships using CAM. As such, this combination will bring about much improved both accuracy and robustness in the recognition of facial expressions, making this updated MobileNetV3 architecture particular-



**Fig. 4.** Methodology framework for FER using enhanced MobileNetV3 with CA and DKA.

ly suitable for applications requiring high performance with very resource-constrained devices. The integrated CAM MobileNetV3 architecture is presented in Fig. 3.

To improve facial expression recognition, we propose an improved methodology using MobileNetV3, augmented with activation of CAM, DKA, and SoftSwish activation. Figure 4 illustrates this integrated framework, showing how these components work together to improve focus, adaptability, and training stability.

Figure 4 presents a simplified view of the proposed methodology, which integrates three key components within the MobileNetV3 framework: the CAM, DKA, and the SoftSwish activation function. These components are designed to enhance the focus on relevant facial features, adaptively optimize kernel sizes for varying image complexities, and improve training stability, respectively. This combination aims to increase both the accuracy and computational efficiency of facial expression recognition systems in resource-constrained environments.

## 4. Results and Discussion

As for the execution of our experiments we used a personal computer with a 64-bit operating system and an Intel Core i7-3.0 GHz and 16 GB of RAM. Experimentation of all of the above approaches was performed using Python.

### 4.1. The JAFFE Database

The JAFFE database is made up of faces of Japanese women and includes both profile and frontal views of these women’s faces. The images are in grayscale and have a resolution of  $256 \times 256$  pixels [22]. The database shown in Fig. 5 is widely familiar with image processing and facial expression analysis. It is extensively used in research and is often used in the creation and assessment of machine learning algorithms commonly used in facial expression recognition. Within the database, there are several pictures of facial expressions of different emotions such as happy, sad, angry, and disgusting emotions, which makes it more useful in training of the FER algorithms.

### 4.2. The CK+ Database

The CK+ database, also known as the extended Cohn-Kanade database, can become a helpful tool in the field of facial expression analysis and computer vision. This database was created as an expansion of the original Cohn-Kanade database

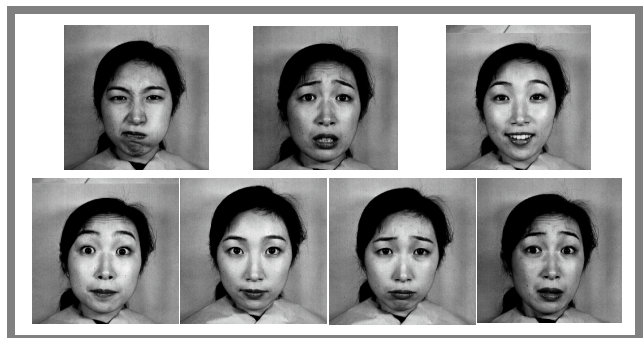
with the goal of increasing the variability and richness of captured expression [23].

The key technical aspects of the CK+ database are as follows:

- Image size – the images used in the CK+ database on average are of  $256 \times 256$  pixels, meaning that the dimension of the images maintained was homogeneous.
- Facial expressions – the facial expressions covered by the database comprise, but are not limited to, happy, sad, angry, surprised, disgusted, and afraid. This diversity enables researchers to assess models through the wide range of emotions.
- Controlled environment – the images are taken in a controlled environment which is very important in standardized environment and scaling out environmental factors that may influence facial expression analysis.
- Subjects – this is a factor that breaks the homogeneity of the data, and several subjects make entries into the database. This variety is useful for testing the extent of generalization of developed facial expression recognition models.
- Annotations – the CK+ images are frequently provided with facial landmarks and emotion labels in addition to the geometric ones. This annotation is useful for both teaching and testing machine learning algorithms, especially for recognizing facial expressions. An illustration of the database is provided in Fig. 6.

### 4.3. The FER2013 Dataset

The FER2013 dataset contains grayscale images of faces, which are  $48 \text{ pixels} \times 48 \text{ pixels}$ . It encompasses seven facial expressions: happiness, anger, disgust, fear, sadness, surprise,



**Fig. 5.** A partial image of the JAFFE database was used to carry out the analysis.



**Fig. 6.** A sample of the images in the CK+ database.

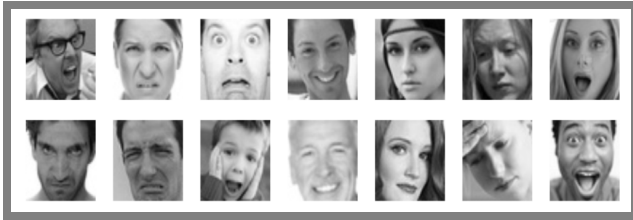


Fig. 7. Subset of the image from the FER2013 database.

and none/neutral. This is acted data set that can be used in training and testing facial expression recognition models, as it provides realistic challenges. It is most often divided into training, validation, and testing set to create a standard for the comparison of the different methods. This is particularly because the network is tiny in size, and thus ideal for deep learning. Scientists apply FER2013 to design and evaluate models to mitigate problems with lighting, head orientation, and different emotions [24]. An example of the proposed database is depicted in Fig. 7.

#### 4.4. Dataset Overview and Scalability Consideration

The computational realm of this study needs to be clarified better through understanding of the data sets that we used during experiments. The JAFFE dataset includes 213 grayscale images depicting ten subjects with posed facial expressions. The CK+ dataset contains 593 image sequences acquired from 123 different subjects, although expression labels are provided only in the final frames of each sequence. The FER2013 dataset is significantly larger, comprising more than 35 000 grayscale images divided into training, validation, and testing sets, with each image sized at  $48 \times 48$  pixels.

The different scales of our datasets enabled the evaluation of the MobileNetV3 + CAM + DKA model in terms of both efficiency and scalability. The model demonstrated consistent accuracy and maintained a similar inference speed, even when handling a large volume of FER2013 data, confirming its suitability for real-time applications in diverse resource-constrained environments.

#### 4.5. Data Augmentation

Data augmentation is a critical step in FER, especially when using a complex model such as MobileNetV3 with CAM and DKA. This essential strategy involves generating additional samples by applying various transformations to the training set, ensuring that the model becomes more effective and robust. Several augmentation methods can be used to improve the accuracy of facial expression recognition, including the following.

- Speckle noise augmentation – this method superimposes speckle noise on the face image by multiplying it by some random numbers, which are helpful in mimicking natural conditions and improving the model’s ability to discern noise.
- Random rotation enhancement – this method involves rotating images containing faces around the vertical or

Tab. 1. Data augmentation techniques with the parameters used.

Augmentation technique	Parameters
Speckle noise augmentation	Noise factor: 0.1
Random rotation	Random angle: $-20^\circ$ to $+20^\circ$
Random zoom	Zoom range: 0.8 to 1.2
Random crop	Crop ratio: 80% of original image size
Color jitter	Saturation range: 0.5 to 1.5
	Brightness range: $-0.3$ to 0.3

horizontal axis in a random manner and thus, it improves the model’s capability to identify different human expressions.

- Random zoom augmentation – this one zooms in and out on the images randomly – the idea being that the model has to be able to learn about the faces and the expressions on them at random sizes.
- Random crop augmentation – this involves tear and shear where the method entails taking a part of the image and discarding the other part leaving the neural network to recognize parts of the face.
- Color jitter augmentation – this technique adds random variation in the hue, saturation, and brightness; adds variations that the model did not receive in the training phase.

These parameters have been chosen to allow proper augmentation without compromising the validity of the data on facial expressions. By integrating these methods, the model is expected to benefit from adaptation to different changes in the environment. Executing them enables the enhancement of facial expression recognition in real-world settings.

## 5. Experimental Steps

All aspects of the experimental setup were carefully designed for facial expression recognition, and we specifically designed a scenario to provide an in-depth review of the proposed model, based on MobileNetV3 combined with CAM and DKA. This effort helped to analyze the capabilities of the proposed model for more complex real-world facial expressions.

The training set was the largest part of the data, representing approximately 70% that was essential to develop the deep neural network model to learn. Its size allowed the model to discern intricate patterns, distill complex correlatives, and finally trace fine nuances linked to various forms of facial expression.

A 15% size validation set was particularly important to further the model intricacies. It allowed for addressing matters with hyperparameters and improvements in the general performance and was a credible line of defense against overfitting. Consequently, the other 15% of the data set was kept for the purpose of the validation test. The validation set was therefore set apart for the sole testing of the model. Its goal was to shine on a model, evaluating or testing its performance in

recognizing unknown facial expressions that were previously known to confirm the efficient operating mode.

### 5.1. Subject-independent Data Splitting

To ensure fair evaluation and prevent data leakage, a subject-independent splitting strategy was used for the CK+ and JAFFE datasets. Specifically, individuals appearing in the training set were excluded from both the validation and the testing sets. For CK+, we used image sequences from approximately 80% of the subjects for training and reserved the remaining 20% for testing and validation. For JAFFE, images from 7 subjects were used for training and 3 subjects for testing and validation. This ensures that the performance reflects its ability to generalize to unseen individuals.

### 5.2. Enhanced Model Architecture

MobileNetV3 is the backbone network for the proposed model that provides the ability to extract features from facial images. It is relatively lightweight and even more appropriate for use with mobile devices, making it ideal for real-time use. The latest MobileNetV3 is integrated into this model to ensure that the model obtains the desired characteristics of being light and at the same time capable of producing high-level features.

The coordinate attention module is another improvement that gives our model the ability of the spatial awareness layer. This module works on the basis of variations in the weightage of various parts of the image, as depicted by the geographical coordinates. In the context of facial expression recognition, CAM opens the possibility of letting the model concentrate on important face areas, since it assigns different weights to the regions. The adaptation mechanism coping strategy increases the efficiency of feature extraction and the precision of capturing facial alterations. In this way, CAM improves the ability to space examination of facial images in space compared to the initial model.

Dynamic kernel adaptation (DKA) is an essential improvement to the adaptability and performance of our facial expression recognition model. In contrast to standard networks that have kernels of a fixed size, within DKA there is an option for the kernel size to be modified in the course of training. This flexibility is advantageous in dealing with the essentially different levels of difficulty in facial images. For example, when there is emotion in the face and the concerns are subtle, DKA allows a larger kernel, which means that features are extracted with a better accuracy.

On the other hand, in simple scenes, the size can be reduced to ensure cost savings and be better optimized to increase its predictive power. Such a dynamic adjustment mechanism ensures that the model is robust across various datasets but is also more efficient in terms of consumption of computational resources. Integration of DKA with MobileNetV3 along with CAM greatly improves the broad applicability of the model in various datasets involving different facial expressions, while boosting the real-time performance of the model.

### 5.3. Training Details

The training process was optimized to ensure peak performance, with particular emphasis on integrating the DKA technique. The training was carried out over 160 epochs, which allowed the model to effectively identify hierarchical structures and representations of the given facial expression data with the help of DKA, which also made it more versatile. The batch size of 32 has been chosen intentionally as it allows one to achieve rather efficient training without overloading the system with computations. This option was most useful when combined with DKA, which modulated the kernel parameters during the training phase. The ranger optimizer was used with a learning rate of 0.001.

By combining two methods known as the RAdam or rectified Adam and LookAhead methods, RAdam fixes problems with fluctuating step size with the better adjustment of the learning rate for adjusting the step size feature, while LookAhead accelerates optimization by coming up with better solutions to improve convergence. Combined with DKA, it was possible to achieve good and stable convergence in this setup. This dynamic adjustment of the learning rate enabled the various phases of training to make optimal use of the learning rate, thereby improving the model performance as well as stability.

To avoid overfitting and improve the ability to generalize, an improved early stopping technique was used. The validation technique continued to update the accuracy of the chosen model on another set of never-before-seen data. Training was, in fact, stopped if no enhancement was observed in the subsequent epochs up to a prescribed number of epochs. This active approach also protected from overfitting with the help of DKA and adjusted the stopping criteria in response to changes in the kernel adjustments and the validation performance trend to make sure the model would perform well in response to a new data set.

To ensure a robust evaluation, all training experiments were repeated over 10 independent runs with random initialization. The accuracy results represent the mean  $\pm$  standard deviation (std) of these runs for each dataset. The activation function used throughout was SoftSwish, defined in Eq. (1). Early stop was employed based on validation accuracy to prevent overfitting. The detailed hyperparameter settings used during training are presented in the Tab. 2.

**Tab. 2.** Training hyperparameters and settings.

Hyperparameter	Value
Optimizer	Ranger (RAdam + LookAhead)
Learning rate	0.001
Batch size	32
Epochs	160
Activation function	SoftSwish
Repetitions	10 runs (mean $\pm$ std reported)
Early stopping	Patience = 10 epochs

**Tab. 3.** Model complexity and inference speed comparison.

Model	Parameters [M]	FLOPs [M]	Inference time [ms/sample]
MobileNetV3 (baseline)	2.9	219	21.3
MobileNetV3 + DKA	3.1	232	19.6
MobileNetV3 + DKA + CAM	3.5	245	18.7

**5.4. Computational Efficiency Analysis**

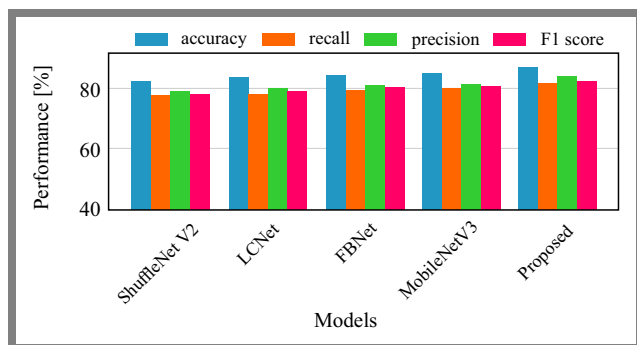
To support our claims regarding computational efficiency, we evaluated and compared the parameter count, FLOPs, and inference latency of the baseline MobileNetV3 model and its enhanced versions with DKA and CAM. These values were obtained using the PyTorch profiler and averaged over 100 runs (Tab. 3).

These results show that the addition of DKA and CAM leads to only a modest increase in parameter count and FLOPs, while achieving a significant reduction in latency of approximately 12.2% faster inference, demonstrating suitability for deployment in low-resource environments.

**5.5. Performance Comparison of FER Models Across Multiple Datasets**

This section explores the performance of various facial expression recognition models across three widely used datasets. FER2013, CK+, and JAFFE. By comparing accuracies, recalls, precisions, and F1 scores, we intend to show the advantages and disadvantages of given models while visually proving the applicability of the presented approach. Examples of performance measures for each data set are shown in Figs. 8–10.

The proposed model stands out with the highest performance metrics across all criteria evaluated in FER2013 dataset analysis. It achieves an accuracy of 87.1%, recall of 81.8%, precision of 83.9%, and an F1 score of 82.8%. This is a significant improvement over other models, which shows its superior capability to recognize facial expressions accurately and consistently. MobileNetV3 follows, showing respectable performance with an accuracy of 84.8% and an F1 score of 80.5%. This means that it is relatively good, though somewhat less so than the model proposed here. LCNet ranks next, ranked by progressively worsening metrics, followed by FB-

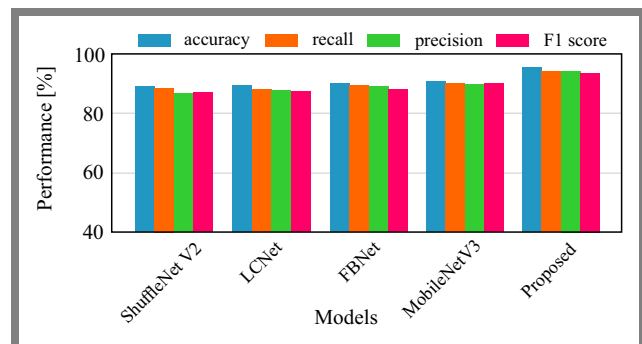


**Fig. 8.** Performance metrics for different models in the FER2013 dataset.

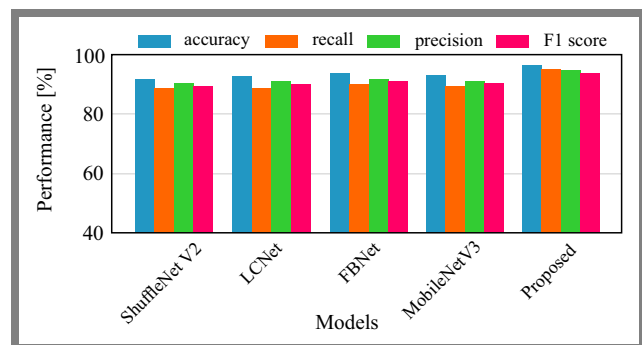
Net and ShuffleNet V2. Although they are equally relevant, these models have relatively less accuracy than our model and therefore suggest that even more complex architectures may be helpful for this dataset.

The results in the CK+ dataset showed great performance with a precision of 95.8% and a good F1 score of 94.1%. It surpasses all other models, and thus it is suggested that it is perfectly suited to the task of facial expression analysis within this dataset. MobileNetV3 also performs well with a given accuracy of about 91.1% and an F1 score of about 90.3%. This clearly depicts the strength of the model, though it is not as efficient as our model. The results of FBNet, LCNet and ShuffleNet V2 are also acceptable, but the difference in metrics can be observed. This implies that the advanced features of the proposed model improve the performance on this data set to a large extent.

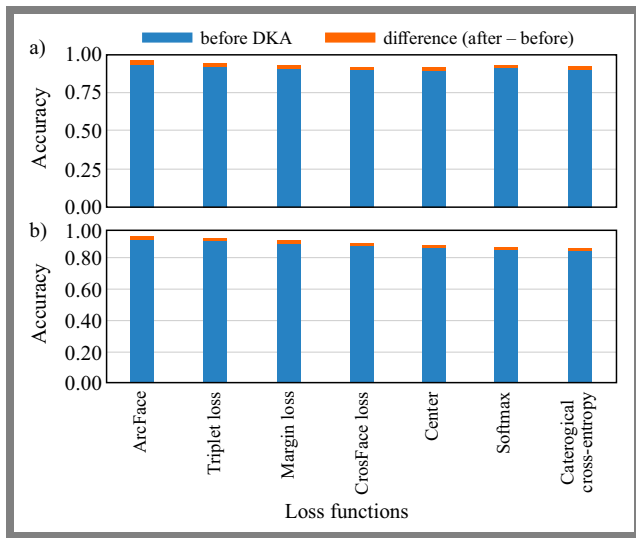
For the JAFFE data set, it is significant when we note that our model offered an accuracy of 96.6% and an F1 score of 94.1%. This data set also validates the effectiveness and generality of the model identified in this study with other related data sets. It can be seen that both FBNet and MobileNetV3 have good accuracy and F1 scores considering the fashion data set and



**Fig. 9.** Performance metrics for the models evaluated in the CK+ dataset.



**Fig. 10.** Performance metrics for the models evaluated in the JAFFE dataset.



**Fig. 11.** Comparative analysis on the JAFEE dataset for: a) SoftSwish and b) h\_swish activation functions.

are still lower than the designed model. LNet and ShuffleNet V2 have good results. However, they are not very efficient compared to the proposed model. Thus, the high stability of the performance of our model on all four sets proves that it is reliable and efficient for facial expression recognition.

**5.6. Comparative Analysis of Loss Functions**

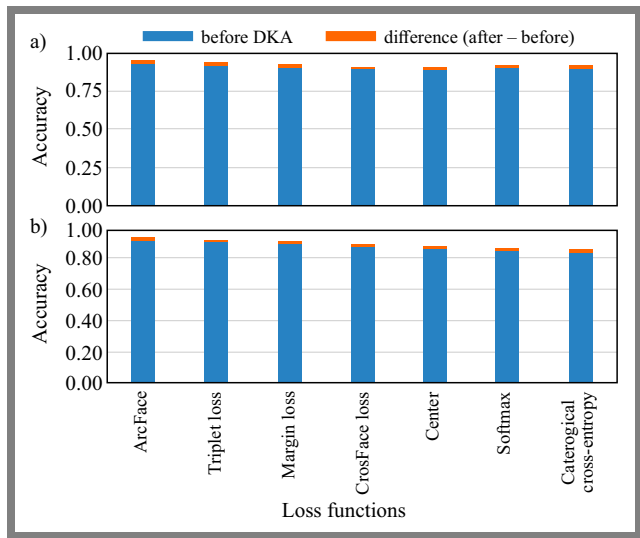
In the next step, we study a comparative study of the applied h\_swish and SoftSwish activation functions in models of facial expression recognition using different sets of data: JAFEE, CK+, and FER2013 was conducted. Its purpose is to evaluate these activation functions to improve the performance and regularization of facial expression recognition. This is demonstrated in Figs. 11–13 and denotes how these functions work before and after DKA, revealing the effects on the model prediction capability in the different datasets.

Looking at the results obtained on the JAFEE dataset, the SoftSwish model should be noticed, which demonstrates rather high indicators both before and after the application of the DKA. Before DKA, SoftSwish was as accurate as 0.939 and after DKA it was 0.966. This implies that the intervention of DKA leads to a drastic improvement in the performance of SoftSwish. In the same way as with h\_swish, the accuracy was higher after DKA as well: 0.920 before DKA and 0.939 after DKA.

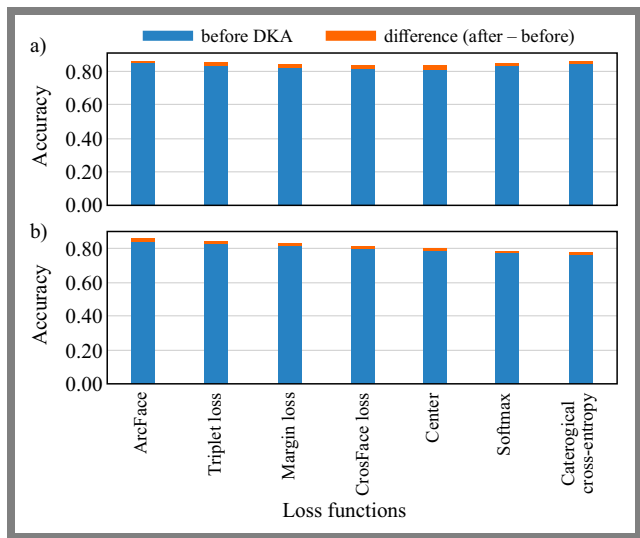
The same was the case in the CK+ dataset, where SoftSwish averaged a 0.9369 before DKA and a 0.9580 at the end of DKA. h\_swish also showed nearly equally good results as the latter, with an accuracy improvement from 0.9173 before DKA to 0.9375 after DKA.

SoftSwish was found to have an improvement on the FER2013 data set with an increase in precision from (0.8559 to 0.8710) after applying DKA. For the part of h\_swish, there was an improvement in the level of accuracy from (0.8413) before the use of DKA to (0.8585) after the implementation of DKA.

From these results, it can be concluded that, for h\_swish and SoftSwish, the application of DKA has a considerable impact



**Fig. 12.** Comparative analysis on the CK+ dataset: a) SoftSwish and b) h\_swish activation functions.

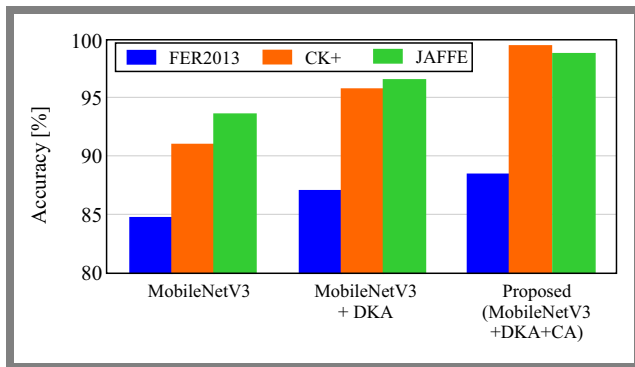


**Fig. 13.** Comparative analysis of the FER2013 data set for a) SoftSwish and b) activation functions.

on the enhancement of the accuracy of facial expression recognition on various datasets.

**5.7. Comparative Analysis of Enhanced FER Models Across Datasets**

To check the general performance of the FER models, it is important to compare the results of their assessment across various data sets. The appearance of the datasets is different and the main problems associated with them are the variability in poses, illumination, and the acquisition of facial images for FER2013, CK+, and JAFFE. In Fig. 14, we provide a comprehensive comparison analysis of the improved MobileNetV3 model with DKA and CA added compared to the baseline MobileNetV3 and the MobileNetV3 model that was improved only with DKA. This will help in the focus of the paper to present the enhancements that have been made as a result of the use of proposed approaches and hence es-



**Fig. 14.** Comparative analysis of the FER2013 data set for a) SoftSwish and b) activation functions.

establish the effectiveness of the model in detecting emotions from faces with high accuracy in similar datasets.

From the analysis of the performance in all the evaluated dataset; FER2013, CK+ and JAFFE as presented in Fig. 14, incorporating DKA and CA in the MobileNetV3 platform has boosted the performance of the model.

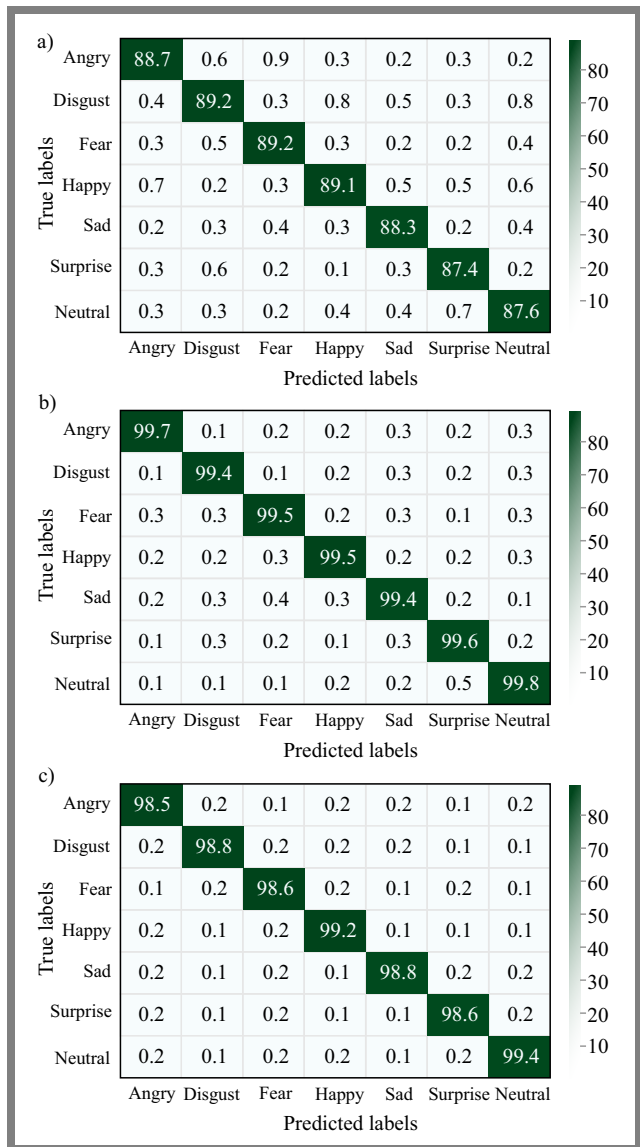
On the FER2013 dataset, which is known for its challenging real-world images with varying lighting conditions and facial expressions, our enhanced model achieved a notable accuracy of 88.5%. This represents an improvement of 3.7% over the baseline MobileNetV3. The inclusion of DKA allows the model to dynamically adapt its convolutional kernels, thus improving its ability to capture finer details in complex scenes. Meanwhile, CA enhances the model’s focus on crucial facial regions, ensuring that the most relevant features are emphasized during the recognition process.

The CK+ dataset, characterized by its controlled environment and a wide range of emotional expressions, further highlights the advantages of these enhancements. Here, the proposed model achieved an accuracy of 97.17%, significantly exceeding the accuracy of 91.1% of the baseline model.

This remarkable improvement to the effectiveness of CA in refining the model’s attention to spatial details, particularly in distinguishing subtle differences in expressions. Additionally, DKA contributes to the model’s flexibility in processing varied facial expressions, thereby boosting its overall accuracy.

In the JAFFE dataset, which includes images of Japanese female subjects displaying different emotions, our model achieved an impressive precision of 97.84%. This result further underscores the model’s robustness in handling diverse populations and expression intensities. The combined effect of DKA and CA allows the model to generalize well across different demographic groups, ensuring consistent performance even in data sets with specific cultural or gender-related characteristics.

Overall, the integration of DKA and CA techniques into MobileNetV3 has proven to be a powerful approach to enhance facial expression recognition. These techniques not only improve accuracy but also ensure its reliability and adaptability across various challenging scenarios.



**Fig. 15.** Confusion matrix for: a) FER2013, b) CK+, and c) JAFFE datasets.

### 5.8. Confusion Matrix for all Datasets

To evaluate the model’s performance in more detail across different datasets, confusion matrices were employed to analyze the accuracy of the model in classifying facial expressions. These matrices are a powerful tool for understanding how well the model distinguishes between different classes and accurately identifies the correct expressions. Figure 15 illustrates the confusion matrices obtained for the CK+, FER2013, and JAFFE datasets.

The confusion matrices for the three datasets (CK+, FER2013, and JAFFE) demonstrate the outstanding performance of the proposed model in FER. The model exhibits exceptionally high accuracy on both the CK+ and JAFFE datasets, with diagonal values ranging from 99.2% to 99.8%, indicating its ability to accurately distinguish between different expressions in controlled environments. In contrast, despite the complexity of the FER2013 dataset, which includes images under various real-world conditions, the model still achieves de-

**Tab. 4.** Comparison of the performance of different FER methods using attention mechanisms for the JAFFE, CK+, and FER2013 datasets.

Method	JAFFE	CK+	FER2013
MobileNetV3 + CAM + DKA (proposed method)	98.84%	99.56%	88.50%
MobileNetV3 + DKA	97.84%	97.17%	87.10%
Baseline MobileNetV3	96.60%	91.10%	84.80%
Attention mechanism-based CNN for FER [25]	88.81%	82.16%	79.33%
FER using LBP and CNN networks that integrate the attention mechanism [26]	90.70%	99.48%	71.29%
FER method combined with the attention mechanism [27]	82.16%	88.81%	79.33%
Auto-fernet – fer network with architecture search [28]	97.14%	98.89%	73.78%
Lightweight FER with key region fusion [29]	82.16%	88.81%	79.33%

cent accuracy, reaching up to 89.2% in the best cases, with relatively low misclassification rates. Overall, these results underscore the model’s reliability and efficiency in handling facial expressions across various datasets, proving its effectiveness in delivering precise and consistent performance even in challenging scenarios.

### 5.9. Evaluation of MobileNetV3 + CAM + DKA Across Datasets

From the analysis, it is evident that the MobileNetV3 + CAM + DKA model produces high accuracy in facial expression recognition test trials. Thus, for the JAFFE dataset, the proposed model yields a high accuracy of (98.84%), which means that it can work efficiently on experiments, including a variety of facial expressions. These findings are reflected in the confusion matrix, which shows the ability of the chosen model to correctly recognize each of the expressions with an emphasis on the neutral one.

In the CK+ dataset, the model maintains a notably high accuracy, achieving a rate of 99.56%, this further reaffirms the strength of the model as there is consistency in its performance across the three datasets. Furthermore, when tested in the FER2013 database, the precision is equal to 88.5%, which means that even under realistic conditions, it successfully identifies different expressions.

The inclusion of CAM and DKA partly enhances the effectiveness and general accuracy of the model based on all defined sets. CAM improves feature capture by optimally directing focus, and DKA enhances the model’s ability to generalize, making it more adaptable and effective in real-world scenarios.

Despite the high accuracy observed on the CK+ and JAFFE datasets (more than 99%), we acknowledge that such small and homogeneous datasets carry a risk of overfitting. To mitigate this, we adopted strict subject-independent evaluation protocols, ensuring no identity overlap across training, validation, or test sets. Furthermore, we conducted multiple training runs and reported mean  $\pm$  std results to verify consistency. Future work will extend our evaluation to larger and more

diverse datasets, such as AffectNet and Occlusion-FER, to validate generalization under real-world conditions.

## 6. Comparison with Different Approaches

Table 4 compares the performance of the proposed method with other state-of-the-art FER approaches that use attention mechanisms on the JAFFE, CK+ and FER2013 datasets.

The proposed MobileNetV3 + CAM + DKA model demonstrates improved performance because its three integrated enhancements include CA for spatial awareness, DKA for flexibility, and improved training stability using the SoftSwish activation function. The complete proposed model delivers better accuracy results than both the baseline MobileNetV3 and its DKA-only variant when evaluated on all datasets.

Performance gain requires accepting a more complex model structure together with additional parameters. With the addition of the CAM, the parameter count increases slightly, but the feature detection accuracy improves significantly under challenging lighting conditions and facial obstructions. With adaptive computation implemented through DKA, the model can dynamically modify kernel sizes, leading to better generalization performance without negatively affecting inference speed.

The performance strength of the LBP + CNN and Auto-FERNet methods in controlled datasets proves insufficient to address the diverse conditions of FER2013, due to limited spatial feature encoding or dependence on handcrafted features. Our method strikes a balance between performance and computational efficiency, making it more suitable for edge deployment compared to heavier models such as those based on ResNet architecture for FER tasks.

The comparison of results is presented in the Tab. 4 highlight the effectiveness of the proposed method, which integrates CA with MobileNetV3, across the JAFFE, CK+, and FER2013 datasets. Our method achieved an impressive accuracy of 98.84% in the JAFFE dataset, 99.56% in the CK+ dataset, and 88.50% on the FER2013 dataset, outperforming most other state-of-the-art methods that incorporate attention mechanisms.

Specifically, the proposed method shows a significant improvement over CNN based on the attention mechanism presented in [25], which reported precisions of 88.81%, 82.16% and 79.33% on the JAFFE, CK+, and FER2013 datasets, respectively. This substantial performance gap highlights the superiority of the coordinate attention mechanism in effectively capturing and utilizing spatial information within facial images.

Similarly, the method that involves the integration of LBP as a feature extractor and CNN network by incorporating an attention mechanism as suggested in [26] yielded slightly higher precision in the CK+ benchmark dataset at 99.48%. However, its performance drastically dropped on the FER2013 dataset with an accuracy of only 71.29%. This means that the method proposed in this work is more generalizable in different datasets while retaining a high level of accuracy, as demonstrated in the FER2013 dataset.

Furthermore, the FER method combined with the attention mechanism in [27] that produced the accuracies of 82.16%, 88.81% and 79.33% in the JAFFE, CK+ and FER2013 datasets, respectively, are relatively low compared to the huge improvements recorded by our method. The CA mechanism when combined with MobileNetV3 offers superior results in terms of spatial patterns while also outperforming other work in terms of adaptability to various datasets.

Finally, the lightweight FER with key region fusion method matches the performance of our proposed method on all three datasets. This close competition indicates that, while both methods are highly effective, the CA mechanism, when combined with MobileNetV3, provides a competitive edge, particularly in scenarios that require a balance of accuracy and computational efficiency.

## 7. Conclusions

In this work, we achieved a significant advancement in FER by incorporating the coordinate attention mechanism into the MobileNetV3 CNN architecture, further enhanced with dynamic kernel adaptation and SoftSwish activation. This novel combination leverages the strength of each of its constituents and, therefore, provides up-to-date, substantial improvements over the existing FER methodologies.

The approach introduced is very strong and accurate; it establishes itself as a state-of-the-art solution in the field. The proposed CA mechanism finely embeds the positional information to sharpen the model's attention on the important facial features, while MobileNetV3, with the help of DKA and SoftSwish, contributed toward high computational efficiency and adaptability to varying image complexities.

However, this study also recognized some limitations. The major limitation is that it depends on a relatively small set of experimental data and may not capture much diversity in real-world tasks. Furthermore, behavior under different lighting conditions, occlusions, and other challenging environments has not fully investigated. These drawbacks will be covered in further research by extending the evaluation with

datasets that present a large diversity of lighting conditions and occlusions, as well as deepening the further contextual enhancements to make the model stronger in terms of reliability and generalization.

## References

- [1] R.I. Bendjillali, M. Beladgham, K. Merit, and T.A. Abdelmalik, "Wavelet-based Facial Recognition", *2018 6th International Conference on Control Engineering & Information Technology (CEIT)*, Istanbul, Türkiye, 2018 (<https://doi.org/10.1109/CEIT.2018.8751751>).
- [2] R.I. Bendjillali *et al.*, "A Robust-facial Expressions Recognition System Using Deep Learning Architectures", *2023 International Conference on Decision Aid Sciences and Applications (DASA)*, Annaba, Algeria, 2023 (<https://doi.org/10.1109/DASA59624.2023.10286798>).
- [3] M. Kamline, M.L. Abdelmounaim, and R.I. Bendjillali, "Arabic Handwriting Recognition System Based on Genetic Algorithm and Deep CNN Architectures", *2021 International Conference on Decision Aid Sciences and Application (DASA)*, Sakheer, Bahrain, 2021 (<https://doi.org/10.1109/DASA53625.2021.9682380>).
- [4] R.I. Bendjillali, M.S. Bendelhoum, A.A. Tadjeddine, and M. Kamline, "Deep Learning-powered Beamforming for 5G Massive MIMO Systems", *Journal of Telecommunications and Information Technology*, no. 4, pp. 38–45, 2023 (<https://doi.org/10.26636/jtit.2023.4.1332>).
- [5] A. Howard *et al.*, "Searching for MobileNetV3", *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, Seoul, South Korea, 2019 (<https://doi.org/10.1109/ICCV.2019.00140>).
- [6] Q. Hou, D. Zhou, and J. Feng, "Coordinate Attention for Efficient Mobile Network Design", *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Nashville, USA, 2021 (<https://doi.org/10.1109/CVPR46437.2021.01350>).
- [7] J.L. Ngwe *et al.*, "PAtt-Lite: Lightweight Patch and Attention MobileNet for Challenging Facial Expression Recognition", *IEEE Access*, vol. 12, pp. 79327–79341, 2024 (<https://doi.org/10.1109/ACCESS.2024.3407108>).
- [8] J. Zhu and Y. Cao, "Face Expression Recognition Combining Improved DeeplabV3+ and Migration Learning", *Journal of Physics: Conference Series*, vol. 2555, 2023 (<https://doi.org/10.1088/1742-6596/2555/1/012020>).
- [9] Y. Gan, J. Chen, Z. Yang, and L. Xu, "Multiple Attention Network for Facial Expression Recognition", *IEEE Access*, vol. 8, pp. 7383–7393, 2020 (<https://doi.org/10.1109/ACCESS.2020.2963913>).
- [10] Z. Hu and C. Yan, "Lightweight Multi-scale Network with Attention for Facial Expression Recognition", *2021 4th International Conference on Advanced Electronic Materials, Computers and Software Engineering (AEMCSE)*, Changsha, China, 2021 (<https://doi.org/10.1109/AEMCSE51986.2021.00143>).
- [11] R. Poyiadzi *et al.*, "Domain Generalization for Apparent Emotional Facial Expression Recognition across Age-groups", *ArXiv*, 2021 (<https://doi.org/10.48550/arXiv.2110.09168>).
- [12] I. Dominguez-Catena, D. Paternain, and M. Galar, "Gender Stereotyping Impact in Facial Expression Recognition", *7th Workshop on Data Science for Social Good*, Grenoble, France, 2022 ([https://doi.org/10.1007/978-3-031-23618-1\\_1](https://doi.org/10.1007/978-3-031-23618-1_1)).
- [13] S. Xie, M. Li, S. Liu, and X. Tang, "ResNet with Attention Mechanism and Deformable Convolution for Facial Expression Recognition", *2021 4th International Conference on Information Communication and Signal Processing (ICICSP)*, Shanghai, China, 2021 (<https://doi.org/10.1109/ICICSP54369.2021.9611962>).
- [14] M. Zeng, Y. Luo, and G. Liu, "Lightweight Facial Expression Recognition Network with Dynamic Deep Mutual Learning", *Proc. of the 2023 3rd International Conference on Bioinformatics and Intelligent Computing*, pp. 222–226, 2023 (<https://doi.org/10.1145/3592686.3592726>).

- [15] B. Zoph, V. Vasudevan, J. Shlens, and Q.V. Le, “Learning Transferable Architectures for Scalable Image Recognition”, *IEEE Conference on Computer Vision and Pattern Recognition*, Salt Lake City, USA, 2018 (<https://doi.org/10.1109/CVPR.2018.00907>).
- [16] A.G. Howard *et al.*, “MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications”, *ArXiv*, 2017 (<https://doi.org/10.48550/arXiv.1704.04861>).
- [17] M. Tan and Q.V. Le, “EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks”, *ArXiv*, 2019 (<https://doi.org/10.48550/arXiv.1905.11946>).
- [18] M. Tan *et al.*, “MnasNet: Platform-aware Neural Architecture Search for Mobile”, *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, Long Beach, USA, 2019 (<https://doi.org/10.1109/CVPR.2019.00293>).
- [19] X. Zhang, X. Zhou, M. Lin, and J. Sun, “ShuffleNet: An Extremely Efficient Convolutional Neural Network for Mobile Devices”, *Proc. of Conference on Computer Vision and Pattern Recognition*, pp. 6848–6856, 2018 (<https://doi.org/10.1109/CVPR.2018.00716>).
- [20] D. Han, J. Kim, and J. Kim, “Deep Pyramidal Residual Networks”, *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Honolulu, USA, 2017 (<https://doi.org/10.1109/CVPR.2017.668>).
- [21] P. Ramachandran, B. Zoph, and Q.V. Le, “Searching for Activation Functions”, *ArXiv*, 2017 (<https://doi.org/10.48550/arXiv.1710.05941>).
- [22] M. Lyons, S. Akamatsu, M. Kamachi, and J. Gyoba, “Coding Facial Expressions with Gabor Wavelets”, *Third IEEE International Conference on Automatic Face and Gesture Recognition*, Nara, Japan, 1998 (<https://doi.org/10.1109/AFGR.1998.670949>).
- [23] P. Lucey *et al.*, “The Extended Cohn-Kanade Dataset (CK+): A Complete Dataset for Action Unit and Emotion-specified Expression”, *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition – Workshops*, San Francisco, USA, 2010 (<https://doi.org/10.1109/CVPRW.2010.5543262>).
- [24] I.J. Goodfellow *et al.*, “Challenges in Representation Learning: A Report on Three Machine Learning Contests”, *ArXiv*, 2013 (<https://doi.org/10.48550/arXiv.1307.0414>).
- [25] J. Li *et al.*, “Attention Mechanism-based CNN for Facial Expression Recognition”, *Neurocomputing*, vol. 411, pp. 340–350, 2020 (<https://doi.org/10.1016/j.neucom.2020.06.014>).
- [26] C. Liang *et al.*, “Facial Expression Recognition Using LBP and CNN Networks Integrating Attention Mechanism”, *2023 Asia Symposium on Image Processing (ASIP)*, Tianjin, China, 2023 (<https://doi.org/10.1109/ASIP58895.2023.00009>).
- [27] M. Chen *et al.*, “Facial Expression Recognition Method Combined with Attention Mechanism”, *Mobile Information Systems*, 2021 (<https://doi.org/10.1155/2021/5608340>).
- [28] S. Li *et al.*, “Auto-FERNet: A Facial Expression Recognition Network with Architecture Search”, *IEEE Transactions on Network Science and Engineering*, vol. 8, pp. 2213–2222, 2021 (<https://doi.org/10.1109/TNSE.2021.3083739>).
- [29] Y. Kong *et al.*, “Lightweight Facial Expression Recognition Method Based on Attention Mechanism and Key Region Fusion”, *Journal of Electronic Imaging*, vol. 30, art. no. 063002, 2021 (<https://doi.org/10.1117/1.JEI.30.6.063002>).

#### Miloud Kamline, Ph.D.

TIT Laboratory, Department of Electrical Engineering

 <https://orcid.org/0009-0007-2949-4859>


E-mail: kamline.miloud@univ-bechar.dz

Tahri Mohammed University, Bechar, Algeria

<https://www.univ-bechar.dz>

#### Ridha Ilyas Bendjillali, Ph.D.

Laboratory of Electronic Systems, Telecommunications and Renewable Energies, Department of Technology

 <https://orcid.org/0000-0003-2465-8192>


E-mail: r.bendjillali@cu-elbayadh.dz

University Center Nour Bachir, El Bayadh, Algeria

<https://www.cu-elbayadh.dz>

#### Mohammed Sofiane Bendelhoum, Ph.D., Associate Prof.

Laboratory of Electronic Systems, Telecommunications and Renewable Energies, Department of Technology

 <https://orcid.org/0000-0002-9789-8712>


E-mail: m.bendelhoum@cu-elbayadh.dz

University Center Nour Bachir, El Bayadh, Algeria

<https://www.cu-elbayadh.dz>

#### Asma Ouardas, Ph.D.

Laboratory of Electronic Systems, Telecommunications and Renewable Energies, Department of Technology

 <https://orcid.org/0000-0002-7569-3572>


E-mail: a.ouardas@cu-elbayadh.dz

University Center Nour Bachir, El Bayadh, Algeria

<https://www.cu-elbayadh.dz>

#### Ali Abderrazak Tadjeddine, Ph.D.

Laboratory of Electronic Systems, Telecommunications and Renewable Energies, Department of Technology

 <https://orcid.org/0000-0003-0926-3440>

E-mail: a.tadjeddine@cu-elbayadh.dz

University Center Nour Bachir, El Bayadh, Algeria

<https://www.cu-elbayadh.dz>

# Information for Authors

**Journal of Telecommunications and Information Technology (JTIT)** is published quarterly since 2000. It comprises original contributions, dealing with a wide range of topics related to telecommunications and information technology. **All papers are subject to peer review.** Topics presented in the JTIT report primary and/or experimental research results, which advance the base of scientific and technological knowledge about telecommunications and information technology.

JTIT is dedicated to publishing research results which advance the level of current research or add to the understanding of problems related to modulation and signal design, wireless communications, optical communications and photonic systems, voice communications devices, image and signal processing, transmission systems, network architecture, coding and communication theory, as well as information technology.

We encourage submissions from a diverse range of authors from across all countries and backgrounds.

## Manuscript

Latex files are preferred and Editorial Office provides a style to prepare the material along with the documentation. We also accept Microsoft Word and PDF files. A typical article is 10 pages long (approximately 6,000 words) and must include the following contents:

- Authors' names and affiliations in the following format:  
First name and surname (last name), academic title,  
Position held,  
ORCID number,  
E-mail address from the University's domain,  
Faculty and name of the University,  
Link to University website.
- Abstract (150-200 words). The abstract should contain statement of the problem, assumptions and methodology, results and conclusion or discussion on the importance of the results. Abstracts must not include mathematical expressions or bibliographic references.
- Keywords related to the content of the article. About four keywords or phrases in alphabetical order should be used, separated by commas.
- The content of the article in a typical structure, i.e.: introduction, related work, conducted research, conclusions, references.

## Figures, Tables and Photos

Together with the article, please send files with graphics with the highest resolution available, 150 dpi or more in bitmap resolution (jpg, png) and vector (cdr, svg, ps, pdf) formats are welcomed.

## References

We use four main citation styles for a journal article, for an Internet article, for a conference paper, and for a book. Below are examples of citations. In each item, the DOI number or link to the PDF of the cited article should be provided.

- [1] R.K. Meyers and A.H. Desoky, "An implementation of the blowfish cryptosystem", *2008 IEEE International Symposium on Signal Processing and Information Technology*, 2008 (<https://doi.org/10.1109/IS-SPIT.2008.4775664>).
- [2] K. Nowicki and T. Uhl, *Ethernet End-to-End*, 1st ed. Germany, Shaker-Publisher, 2008 (ISBN: 978383832271404).
- [3] C. Shorten and T.M. Khoshgoftaar, "A survey on image data augmentation for deep learning", *Journal of Big Data*, vol. 6, no. 1, pp. 1–48, 2019 (<https://doi.org/10.1186/s40537-019-0197-0>).
- [4] S. Wong *et al.*, "Traffic forecasting using vehicle-to-vehicle communication", *3rd Annual Conference on Learning for Dynamics and Control*, pp. 917–929, 2021 (<https://arxiv.org/pdf/2104.05528>).

## Submission

The paper with full PDF version and anonymous PDF version for the blind review process should be submitted on the JTIT website <https://www.jtit.pl/jtit/about/submissions>.

## Reviewing Process

The article is initially approved by the Editor-In-Chief and if the decision is positive, is then sent to the reviewers. Depending on the subject of the article, it takes few weeks. In the next step, reviews are showed to authors who have 2 weeks to correct the article. Finally, the corrected text can be re-presented to the reviewer for reevaluation, which will take another 2 weeks.

As a result, after about 3 months, we are able to send the text for publication in the upcoming issue of JTIT.

When the reviews are inconsistent, additional corrections are necessary, or the reviewer expects additional verification because the corrections ordered by the author are insufficient or additional problems arise, the review of the article may be extended by another month or more.

## Editorial Work

Positively reviewed and corrected article is next prepared by the editorial office for publication. At the end of this process the author receives an copyedited version for approval.

## Licensing

Manuscript submitted to JTIT should not be published or simultaneously submitted for publication elsewhere. By submitting a manuscript the author grants license to the National Institute of Telecommunications, for the use of the paper in the fields of exploitation: reproducing and fixing the paper, distributing the paper by means of introduction to trade, letting for use or rental of the original or copies, and distributing the paper by means of public exhibition, screening, presentation and broadcast as well as rebroadcast, and making the paper publicly available in such a manner that anyone could access it at a place and time selected thereby, or by making it available in a way not allowing selection of time or place, including by means of Internet or other networks.

## Ghostwriting Declaration

We require formal declaration that the process of writing the paper was not influenced by any third party. In the article, all the contributions of other people are clearly indicated. The theories presented, methods used, analysis and research, as well as the copyrights to the drawings, photographs and other figures belong to the authors or are clearly credited in the text. The author must also indicate whether his work has received financial support and if the realization of the whole project was possible thanks to the permission and cooperation with scientific institutions, associations and others.

## Other Information

- The JTIT being an Open Access Journal (OAJ) has no article processing charges (APCs). The published articles can be downloaded freely without payment.
- JTIT supports open access and using continuous publishing "publish-as-you-go" scheme. This means that we no longer wait to accumulate several articles into a quarterly issue before publication. Rather, articles are continuously added to current issues after acceptance. Publish-as-you-go reduces publication lag for our authors, and make the newest research available quickly. After completing the review process, an article is published online in the current issue with DOI registration. When the issue period ends, a new issue is activated. So accepted articles are published without waiting for the quarterly issue end.

**Reconfigurable Reflectarray Structure Based on  
Optimized Unit Cell for Wireless Communications**

*R.M. Yaseen and A.K. Jassim*

78

**Advancing Facial Expression Recognition – Enhanced MobileNetV3 with  
Integrated Coordinate Attention and Dynamic Kernel Adaptation**

*M. Kamline, R.I. Bendjillali, M.S. Bendelhoum, A. Ouardas, and A.A. Tadjeddine*

83



National Institute  
of Telecommunications

**Editorial Office**

National Institute  
of Telecommunications  
Szachowa st 1  
04-894 Warsaw, Poland  
<https://www.gov.pl/web/instytut-lacznosci>

phone +48 22 512 81 83  
fax +48 22 512 84 00

e-mail: [journal@jtitt.pl](mailto:journal@jtitt.pl)  
[www.jtitt.pl](http://www.jtitt.pl)