



ISSN 2224-087X

ЕЛЕКТРОНІКА ТА ІНФОРМАЦІЙНІ ТЕХНОЛОГІЇ

ELECTRONICS
AND INFORMATION TECHNOLOGIES

Збірник наукових праць

Випуск 33



2026

**ELECTRONICS
AND
INFORMATION
TECHNOLOGIES**

Issue 33

Scientific journal

Published 4 issue per year

Published since 1966

**ЕЛЕКТРОНІКА
ТА
ІНФОРМАЦІЙНІ
ТЕХНОЛОГІЇ**

Випуск 33

Збірник наукових праць

Виходить 4 рази на рік

Видається з 1966 р.

**Ivan Franko National
University of Lviv**

**Львівський національний
університет імені Івана Франка**

2026

ЗАСНОВНИК: ЛЬВІВСЬКИЙ НАЦІОНАЛЬНИЙ УНІВЕРСИТЕТ ІМЕНІ ІВАНА ФРАНКА

Друкується за ухвалою Вченої Ради
Львівського національного університету
імені Івана Франка
протокол №5/3 від 25.03.2026 р.

У 1966–2010 рр. збірник виходив під назвою «Теоретична електротехніка»

Збірник «Електроніка та інформаційні технології» містить оригінальні результати досліджень з електронного матеріалознавства, моделювання фізичних явищ, процесів і систем електроніки, обробки сигналів і зображень, інформаційних технологій.

Редактори

Проф. *І. Карбовник* – головний співредактор
Проф. *І. Бордун* – головний співредактор
Проф. *О. Крупиш* – відповідальний редактор
С.н.с. *Я. Шмигельський* – відповідальний секретар

Редакційна колегія

д-р техн. наук, проф. *О. Андрейків*
д-р габіл., проф. *Б. Андрієвський*
д-р фіз.-мат. наук, проф. *І. Болеста*
канд. фіз.-мат. наук, доц. *С. Вельгош*
д-р фіз.-мат. наук, проф. *П. Венгерський*
д-р техн. наук, проф. *Р. Воробель*
д-р фіз.-мат. наук, проф. *Р. Головчак*
д-р техн. наук, д-р габіл., проф. *Ф. Івацішин*
д-р фіз.-мат. наук, проф. *О. Кушнір*
д-р фіз.-мат. наук, проф. *А. Лучечко*
д-р техн. наук, проф. *Л. Муравський*
д-р техн. наук, проф. *М. Назаркевич*
д-р фіз.-мат. наук, проф. *І. Оленич*
д-р фіз.-мат. наук, проф. *Б. Павлик*
д-р техн. наук, доц. *Б. Павлишенко*
д-р техн. наук, проф. *С. Рендзіняк*
д-р фіз.-мат. наук, проф. *Б. Русин*
д-р фіз.-мат. наук, проф. *М. Прутула*
д-р габіл., проф. *Ц. Славінські*
канд. фіз.-мат. наук, доц. *Ю. Фургала*
д-р габіл., проф. *Б. Ціж*,
д-р габіл., проф. *Бушита Шахраї*
д-р фіз.-мат. наук, проф. *Г. Шинкаренко*
канд. фіз.-мат. наук, доц. *Р. Шувар*
д-р фіз.-мат. наук, проф. *І. Яворський*

Адреса редакційної колегії:

Львівський національний університет імені Івана
Франка, факультет електроніки та інформаційних
технологій, вул. Ген. М. Тарнавського, 107,
79017, Львів, Україна
тел. (+38) (093) 864-01-19
e-mail: elit@lnu.edu.ua

web-сайт: <http://publications.lnu.edu.ua/collections/index.php/electronics/index>

Реєстрація суб'єкта у сфері друкованих медіа:
Рішення Національної ради України з питань
телебачення і радіомовлення № 1877 від
30.05.2024 р. Ідентифікатор медіа R30-04912

“Electronics and information technologies” journal contains original research results on electronics material science, modelling of physical phenomena, processes and electronic systems, signal and image processing and information technologies.

Editors

Prof. *I. Karbovnyk* – Chief Co-Editor
Prof. *I. Bordun* – Chief Co-Editor
Prof. *O. Krupych* – Managing Editor
Sen. Res. *Ya. Shmygelsky* – Technical Editor

Editorial Board

O. Andreykiv, Dr. Sc., Prof.
B. Andrievsky, Dr. Habil., Prof.
I. Bolesta, Dr. Sc., Prof.
S. Velgosh, PhD, Assoc. Prof.
P. Vengersky, Dr. Sc., Prof.
R. Vorobel, Dr. Sc., Prof.
R. Holovchak, Dr. Sc., Prof.
F. Ivachyshyn, Dr. Sc., Dr. Habil., Prof.
O. Kushnir, Dr. Sc., Prof.
A. Luchechko, Dr. Sc., Prof.
L. Muravsky, Dr. Sc., Prof.
M. Nazarkevych, Dr. Sc., Prof.
I. Olenych, Dr. Sc., Prof.
B. Pavlyk, Dr. Sc., Prof.
B. Pavlyshenko, Dr. Sc., Assoc. Prof.
S. Rendzinyak, Dr. Sc., Prof.
B. Rusyn, Dr. Sc., Prof.
M. Prytula, Dr. Sc., Prof.
C. Slawinski, Dr. Habil., Prof.
Yu. Furgala, PhD, Assoc. Prof.
B. Tszh, Dr. Habil., Prof.
Bouchta SAHRAOUI, Dr. Habil., Prof.
G. Shynkarenko, Dr. Sc., Prof.
R. Shuvar, PhD, Assoc. Prof.
I. Yavorsky, Dr. Sc., Dr. Habil., Prof.

Editorial office address:

Ivan Franko National University of L'viv,
Faculty of Electronics and Computer Technologies
107 Tarnavsky St., UA–79017,
Lviv, Ukraine
tel. (+38) (093) 864-01-19

АДРЕСА РЕДАКЦІЇ, ВИДАВЦЯ І ВИГОТОВЛЮВАЧА:

Львівський національний університет імені Івана Франка
вул. Університетська, 1, 79000 Львів, Україна

Свідоцтво про внесення суб'єкта видавничої справи до Державного реєстру видавців,
виготівників і розповсюджувачів видавничої продукції. Серія ДК № 3059 від 13.12.2007 р.

© Львівський національний університет імені Івана Франка, 2026

CONTENTS

INFORMATION SYSTEMS AND TECHNOLOGIES

Using large language models for text analysis in the evaluation of university educational programs. <i>Mykola Stasiuk, Vitaliy Kukharskyy, Bohdan Pavlyshenko (12)</i>	5
Phased integration of neural networks of different architectures in mathematical computing systems. <i>Mykhailo Bavdys, Oleksii Kushnir (14)</i>	17
Indoor positioning with bluetooth low energy and extended kalman filter. <i>Tadei-Nazarii Kalynchuk (12)</i>	31
Communication architectures for cloud-integrated spectrometric laboratories: ESP-NOW vs MQTT. <i>Andriy Krupych, Pavlo Levush (14)</i>	43
Evaluation of the performance of a multi-level model for anomalous DNS query detection. <i>Andrii Senyk (14)</i>	57
Application of <i>Penetration Testing</i> for assessing the information security level of web-oriented information systems. <i>Sergiy Sveleba, Ivan Katerynychuk, Ivan Kunyo, Oleh Krupych, Yaroslav Shmyhelskyy, Marta Dufanets, Natalia Sveleba, Lucjan Pelc, Volodymyr Brygilevych (16)</i>	71
Long Short-Term Memory recurrent neural network for state prediction and resource allocation optimization in distributed systems. <i>Zinovii Liubun, Oleh Tereshchuk (10)</i>	87
Intelligent analysis of performance results based on object-relational mapping strategies and foreign key constraints in SQL databases. <i>Oleksandra Ryback, Oleh Husak, Roman Mysiuk (16)</i>	97

MODELING OF PROCESSES AND EFFECTS

Spatial-geometric evaluation of local features in monocular visual odometry. <i>Andriy Fesiuk, Yuriy Furgala (18)</i>	113
Entropy-guided tracker switching method for unmanned aerial vehicle real-time tracking. <i>Volodymyr Oleksiuk, Serhiy Velhosh (14)</i>	131
Predictive thermal management in embedded electronics using deep reinforcement learning. <i>Oleh Yatskiv, Bohdan Koman (20)</i>	145
Vision transformer-based fall detection: a spatial temporal attention mechanism for robust video analysis. <i>Ivan Ursul, Andriy Pereymybidia (16)</i>	165

MATERIALS FOR ELECTRONIC ENGINEERING

The effect of preparation conditions on the electrical conductivity of thin films of $(Y_{0.06}Ga_{0.94})_2O_3$. <i>Ihor Kukharskyy, Iryna Kofliuk, Ivanna Medvid, Ihor Kuz, Oleh Bordun (8)</i>	181
Photodetectors based on field effect in porous silicon – reduced graphene oxide structures. <i>Igor Olenych, Andrii Kozak (10)</i>	191

INTERNET OF THINGS SYSTEMS

Integration of decentralized performance verification in hybrid architectures edge-fog-cloud to increase IoT systems reliability. <i>Roman Diachok, Halyna Klym (8)</i>	201
---	-----

ЗМІСТ

ІНФОРМАЦІЙНІ СИСТЕМИ ТА ТЕХНОЛОГІЇ

Використання великих мовних моделей для текстового аналізу в оцінюванні університетських освітніх програм. <i>Микола Стасюк, Віталій Кухарський, Богдан Павлишенко (12)</i>	5
Поетапна інтеграція нейронних мереж різної архітектури в системах математичних обчислень. <i>Михайло Бавдис, Олексій Кушнір (14)</i>	17
Визначення положення об'єктів в приміщенні з використанням радіотехнології bluetooth low energy та розширеного фільтра Калмана. <i>Тадей-Назарій Калинчук (12)</i>	31
Комунікаційні архітектури для хмарно-інтегрованих спектрометричних лабораторій: порівняння ESP-NOW та MQTT. <i>Андрій Крупич, Павло Левуш (14)</i>	43
Оцінювання продуктивності багаторівневої моделі для виявлення аномальних DNS-запитів. <i>Андрій Сенік (14)</i>	57
Застосування <i>Penetration Testing</i> для оцінювання рівня інформаційної безпеки веб-орієнтованих інформаційних систем. <i>Сергій Свелеба, Іван Катеринчук, Іван Куньо, Олег Крупич, Ярослав Шмигельський, Марта Дуфанець, Наталя Свелеба, Люціан Пельці, Володимир Бригілевич (16)</i>	71
Рекурентна нейронна мережа на основі довгої короткочасної пам'яті для прогнозування станів та оптимізації розподілу ресурсів у розподілених системах. <i>Зіновій Любунь, Олег Терещук (10)</i>	87
Інтелектуальний аналіз результатів продуктивності на основі стратегій об'єктно-реляційного відображення та обмежень зовнішнього ключа в базах даних SQL. <i>Олександр Рибак, Олег Гусак. Роман Мисюк (16)</i>	97

МОДЕЛЮВАННЯ ПРОЦЕСІВ ТА ЯВИЩ

Просторово-геометричне оцінювання локальних ознак у монокулярній візуальній одометрії. <i>Андрій Фесюк, Юрій Фургала (18)</i>	113
Керування перемиканням засобів супроводу на основі ентропійного аналізу для відстеження цілей безпілотних літальних апаратів у реальному часі. <i>Володимир Олексюк, Сергій Вельгош (14)</i>	131
Прогнозне керування тепловими процесами у вбудованих електронних пристроях з використанням глибокого навчання з підкріпленням. <i>Олег Яцків, Богдан Кома (20)</i>	145
Виявлення падінь на основі зорового трансформера: просторово-часовий механізм уваги для надійного аналізу відео. <i>Іван Урсул, Андрій Переймибіда (16)</i>	165

МАТЕРІАЛИ ЕЛЕКТРОННОЇ ТЕХНІКИ

Вплив умов одержання на електропровідні властивості тонких плівок $(Y_{0.06}Ga_{0.94})_2O_3$. <i>Ігор Кухарський, Ірина Кофлюк, Іванна Медвідь, Ігор Кузь, Олег Бордун (10)</i>	181
Фотодетектори на основі ефекту поля у структурах поруватий кремній – відновлений оксид графену. <i>Ігор Оленіч, Андрій Козак (10)</i>	191

СИСТЕМИ ІНТЕРНЕТУ РЕЧЕЙ

Інтеграція децентралізованої перевірки продуктивності в гібридних архітектурах «периферія-туман-хмара» для підвищення надійності систем інтернету речей. <i>Роман Дячок, Галина Клим (8)</i>	201
--	-----

UDC: 004.89

USING LARGE LANGUAGE MODELS FOR TEXT ANALYSIS IN THE EVALUATION OF UNIVERSITY EDUCATIONAL PROGRAMS

Mykola Stasiuk¹  , Vitaliy Kukharskyy²  , Bohdan Pavlyshenko¹  ,

¹System Design Department, Ivan Franko National University of Lviv,
50, Drahomanova St., Lviv, 79005, Ukraine

²Applied Mathematics Department, Ivan Franko National University of Lviv,
1, Universytetska St., Lviv, 79000, Ukraine

Stasiuk, M. I., Kukharskyy, V. M., Pavlyshenko, B. M. (2026). Using Large Language Models for Text Analysis in the Evaluation of University Educational Programs, *Electronics and Information Technologies*, 33, 5–16. <https://doi.org/10.30970/eli.33.1>

ABSTRACT

Background. Large language models (LLMs) are increasingly used in educational analytics, particularly for processing large volumes of accreditation-related documents. However, it remains unclear how reliably such models can assess the quality of self-evaluation reports for educational programs, and which textual characteristics influence how models form their assessments.

Materials and Methods. In the study, ten self-evaluation reports of educational programs were analyzed: five identified by the expert assessment as the strongest within the higher education institution over the last three years, and five as the weakest over the same period. GPT-5 and Gemini-2.5 models independently evaluated each document using the official ten Ukrainian National Agency for Higher Education Quality Assurance (NAQA) criteria and eight textual metrics reflecting structural, semantic, argumentative, and factual properties of the text. All evaluation grades were generated directly by the models on a unified scale from 1 to 10. To analyze the relationships between NAQA and textual criteria, Pearson's and Spearman's correlation coefficients were used.

Results and Discussion. LLMs demonstrated limited alignment with the NAQA criteria, yielding weak correlations. In contrast, textual criteria, primarily factual density, argumentativeness, semantic coherence, and lexical diversity, consistently differentiated between stronger and weaker reports. GPT-5 exhibited lower variability and reduced sensitivity to stylistic noise, while Gemini-2.5 reacted more strongly to structural and stylistic deficiencies. Correlation matrices showed that textual criteria better capture the latent quality characteristics of documents than the direct application of NAQA criteria.

Conclusion. The results show that LLMs currently do not accurately reproduce expert evaluations based on the formal NAQA criteria but effectively analyze the structural and content-related characteristics of reports using textual metrics. These metrics complement the NAQA criteria by accelerating expert workflows and enhancing document monitoring. Future research will focus on expanding the dataset, standardizing prompts, and comparing a broader range of models.

Keywords: large language models, educational programs, quality assessment.

INTRODUCTION

Large Language Models (LLMs) are gradually being integrated into educational analytics, from automated reviews to the creation of self-assessments and accreditation reports [1, 2]. At the same time, the question arises not only of the accuracy of such



© 2025 Mykola Stasiuk et al. Published by the Ivan Franko National University of Lviv on behalf of Електроніка та інформаційні технології / Electronics and Information Technologies. This is an Open Access article distributed under the terms of the [Creative Commons Attribution 4.0 License](https://creativecommons.org/licenses/by/4.0/) which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

assessments but also of the quality of the assessment process itself. This includes how effectively the model explains its own decisions, refers to sources, and demonstrates reasoned motivation. Studies indicate that models can differ significantly in their ability to build logical chains of reasoning and engage in reasoned thinking, which affects the stability and interpretability of their evaluations [3].

Educational program self-assessment texts are the central element of the accreditation process and reflect the aspects of the educational program's functioning required by the regulatory framework of the National Agency for Higher Education Quality Assurance (NAQA). In such documents, the higher education institution must demonstrate the completeness and systematic nature of student preparation, justify the rationale for the curriculum structure, and describe the faculty, the learning environment, resources, and quality assurance mechanisms.

The analysis of large volumes of texts is a significant burden for expert groups, especially amid the growing number of programs. Meanwhile, modern LLMs demonstrate the ability to work with long texts, summarize information, search for patterns, and measure structural and stylistic characteristics. Previous research, particularly approaches such as G-Eval, shows that LLMs can serve as tools for evaluating and analyzing text quality [4]. In this study, textual criteria are considered an additional tool for the analysis of the document's structural and semantic properties, which are not always explicitly represented in the formal criteria. This opens the possibility of using LLMs to analyze self-assessment materials.

Despite growing interest in the use of LLMs in educational analytics processes, little is currently known about their ability to work with educational program self-assessment texts. In the context of Ukrainian-language research, the first attempts to create specialized benchmarks for evaluating LLM capabilities are already emerging [5]. However, it has scarcely been explored whether models can consistently evaluate such texts, differentiate their quality characteristics, and how their assessments correlate with the official NAQA criteria. In this research, the capabilities of GPT-5 and Gemini-2.5 were evaluated by comparing their results on real self-assessment texts.

The purpose of this study is to investigate how modern LLMs interpret and apply the official NAQA criteria when analyzing educational program self-assessment texts, as well as the extent to which LLM assessments are consistent and justified. The models GPT-5 and Gemini-2.5 were used to assess their ability to work with long normative texts, distinguish between strong and weak reports, and identify the key document properties that influence the evaluation. Additionally, it was investigated whether textual criteria can complement the official evaluation system and serve as a tool for preliminary analytical auditing of such documents. A separate focus was placed on analyzing which elements of the documents' content influence the LLM assessments and correlate with program quality.

MATERIALS AND METHODS

In this work, educational programs are understood as higher education degree programs evaluated within formal accreditation frameworks. The study focuses on the analysis of self-evaluation texts and does not address the evaluation of software or programming code.

Ten self-assessment texts of educational programs that are publicly available on the official website of Ivan Franko National University of Lviv were used in the experiments. The sample included five programs that received higher expert evaluations based on accreditation results, and five programs that received the lowest evaluations. All documents were structured in accordance with NAQA requirements, ensuring their comparability. Each text was processed separately by two LLMs: GPT-5 [6] and Gemini-2.5 [7].

Two sets of criteria were applied for evaluating the self-assessment texts. The first contained 10 official NAQA criteria [8], each with a score ranging from 1 to 10. The criteria

cover educational program design, its structure, staffing, learning environment, internal quality assurance system, and other components, as shown in **Table 1**. Since Criterion 10 applies only to third-level educational programs (PhD level), its scores were not included in the subsequent analysis.

Table 1. NAQA criteria.

#	Criterion	Short Description
1	Educational Program Design	Corresponds to the standard of higher education and professional standards. Clear goal, aligned with the HEI's mission. Consideration of labor market needs and trends, and the regional context.
2	Structure and Content of the Educational Program	ECTS credit volume complies with legislation. Logical structure, relevance to the subject area. Possibility of an individual educational trajectory. Practical training, social skills. Consideration of the UN Sustainable Development Goals.
3	Access to the Educational Program and Recognition of Learning Outcomes	Transparent and non-discriminatory admission rules. Procedures for recognizing learning outcomes acquired in other institutions, as well as results of non-formal and informal education.
4	Learning and Teaching within the Educational Program	Student-centered approach, academic freedom. Availability of syllabi. A combination of learning with research. Updating content based on scientific achievements. Internationalization.
5	Monitoring Measures, Evaluation of Higher Education Students, and Academic Integrity	Clear and published evaluation criteria. Clear rules for conducting exams, preventing conflicts of interest, and appeal procedures. Effective policy and culture of academic integrity.
6	Human Resources	Faculty qualifications meet requirements. Transparent competitive selection. Involvement of employers and practitioners in teaching. Promotion of faculty professional development.
7	Educational Environment and Material Resources	Sufficiency of material and technical base, libraries, and software. Safe and inclusive environment. Support for physical and mental health. Anti-corruption policies and anti-discrimination measures.
8	Internal Quality Assurance of the Educational Program	Procedures for program monitoring and review. Involvement of students and employers as partners. Tracking graduate careers. Response to surveys and feedback from previous accreditations.
9	Transparency and Publicity	Clear rules for all participants. Public discussion of program drafts before approval. Publication of complete information on the website.
10	Learning through Research	Preparation for solving complex problems. Alignment of research with the supervisor's field. Presence of specialized academic councils. Resource provision for research. Integration into the international community. Integrity of supervisors.

The formation of the set of textual metrics was based on modern approaches to the qualitative properties of text analysis used in the field of Natural Language Generation (NLG) [4] and in the evaluation of language models. These properties include structural organization, logical coherence, argumentative completeness, semantic connectivity, factual correctness, and lexical variability.

In this study, all metrics were computed directly by the LLMs, which were instructed to assign each metric a value on a scale from 1 to 10 based on the document's content, without any subsequent manual adjustments. The textual criteria used are described in **Table 2**.

Table 2. Textual Criteria for the Assessment of Self-Evaluation Reports.

Metric	Abbreviation	Essence of the Metric
Structural Consistency	S	Integrity and logical coherence of the document's structure; presence of all key sections.
Criteria Coverage	C	The proportion of NAQA criteria for which substantive explanations are provided.
Argumentative Saturation	A	The quantity and quality of arguments, examples, and evidence-based statements in the text.
Lexical Diversity	L	The degree of vocabulary variety, including the proportion of unique lexemes, the presence of terminology, synonyms, and different language registers.
Evaluation Variability	V	Diversity of types of evaluative judgments in the text (positive, negative, neutral), as well as the balance between them and their sources of origin.
Semantic Connectivity	SC	Logical and linguistic coherence between adjacent sentences in the document.
Factual Saturation	F	The number of references to regulatory documents or standards.
Model Motivation Clarity	M	The presence of an explicit explanation for the logic the model followed when forming the evaluation.

The metrics for structural consistency and argumentative saturation rely on the principles of logical organization of text analysis described in modern NLG evaluation systems [4]. Completeness of coverage reflects the model's ability to encompass all structural components of the document correctly and is based on coverage-based assessment approaches in NLG. The factual saturation metric is based on practices for verifying the reliability of statements in language model responses and related research on factuality [9]. The semantic connectivity metric applies approaches to sentence similarity in embedding space analysis, a typical approach for evaluating the coherence of generated texts [10].

The variability metric reflects the diversity of evaluative judgments present in the document's text. It takes into account the balance between positive, negative, and neutral statements, as well as different sources of evaluation, such as claims about strengths, acknowledgment of shortcomings, and neutral descriptions of processes. The lexical diversity metric assesses the degree of vocabulary variation, including the proportion of unique lexemes, the use of terminology, synonyms, and different linguistic registers. It reflects the richness of the language and the stylistic saturation of the text and operates based on the Type-Token Ratio [11]. Clarity of motivation assesses the presence of explicit explanations regarding the logic of evaluation formation and is related to approaches for evaluating the transparency of reasoning processes in LLMs.

Thus, the set of metrics combines generally accepted approaches for studying academic texts with modern methods for evaluating the quality of LLM responses, and is used in this study as an additional analytical tool.

The research focused on three aspects:

- Relationship between NAQA criteria scores and textual metrics. It was checked whether self-assessment text metrics are associated with high or low scores on the NAQA criteria.
- Inter-model consistency. The correlation between the two models' scores was analyzed separately for the NAQA criteria and the textual metrics. This allowed for determining how stable and reproducible the decisions of different LLMs are when working with identical documents.
- Internal structure of each model's evaluations. It was investigated which textual properties most strongly influence the NAQA scores within each model, and whether the models form similar patterns.

The Pearson correlation coefficient [12] was used for quantitative analysis, enabling assessment of the linear relationships between the criterion values on a unified scale. Additionally, Spearman's rank correlation coefficient [13] was employed to analyze relationships based on score rankings. Together, these correlation measures were used to examine whether the textual metrics and NAQA scores generated by the models are sufficient to distinguish between program self-assessment documents previously identified by experts as stronger or weaker.

Since the study sample comprises only 10 documents, the obtained Pearson correlation coefficients should be considered indicative. With so few observations, a single atypical document can significantly influence the values of the correlation indicators. Therefore, the interpretation of correlations is presented as identified tendencies rather than statistically confirmed patterns.

To simplify the analysis and comparison, the results were structured using conditional notation. The five programs that experts previously identified as higher quality are denoted as B1-B5 (Best). The five programs that received the lowest expert evaluations are denoted as W1-W5 (Weakest). All graphical materials, tables, and correlation analysis results use these notations for data unification.

RESULTS AND DISCUSSION

The results of the models' work, presented in **Fig. 1**, demonstrate varying behavior between the NAQA and textual criteria. Both LLMs are generally prone to assigning high scores based on the NAQA criteria. However, significantly more variation is observed in the textual criteria, especially for metrics such as lexical diversity, semantic connectivity, and evaluation variability. The NAQA scores are less sensitive to the text's actual properties than the textual criteria, which more accurately capture the quality of the document's content.

Notably, the Gemini-2.5 model shows a significantly wider spread of textual criteria scores, whereas the GPT-5 model assigns scores less frequently at 5-6 and below. This points to different internal text analysis strategies: GPT-5 focuses on structural integrity and general style, while Gemini-2.5 concentrates more on argumentation, facts, and logical connectivity.

The correlation analysis showed that the relationship between the NAQA criteria and the textual characteristics is weak or unstable in most programs. Even in cases where the text contains a significant number of facts, examples, and clear arguments, the model may assign very high scores for structural criteria but react weakly to the quality of the content. In program W5, GPT-5 assigns 10 points for Criterion 1, "Educational Program Design" but the textual criteria indicate poor semantic connectivity (SC = 6). This highlights a disconnect between the formal fulfillment of the structure and the low quality of the content.

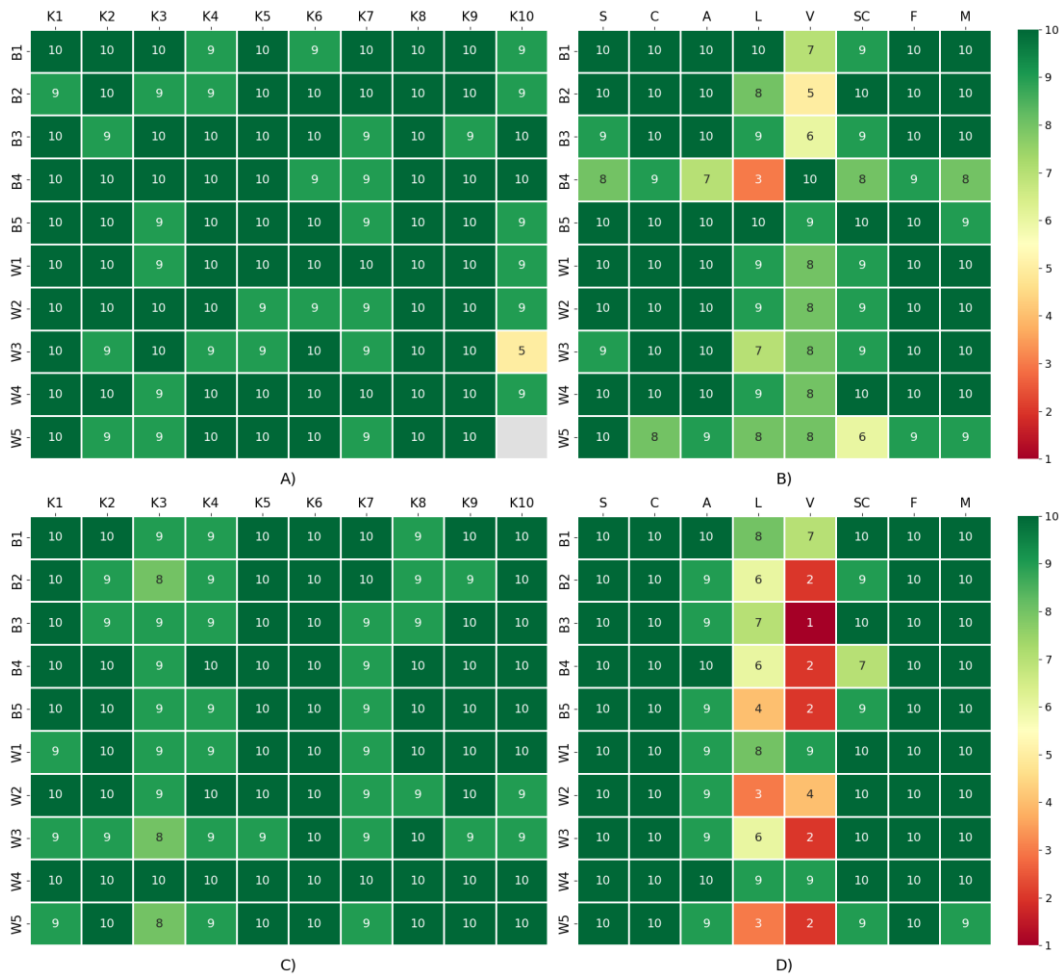


Fig. 1. Self-evaluation report evaluation scores obtained from GPT-5 and Gemini-2.5 models: (A) NAQA criteria evaluation by GPT-5; (B) Textual criteria evaluation by GPT-5; (C) NAQA criteria evaluation by Gemini-2.5; (D) Textual criteria evaluation by Gemini-2.5.

Furthermore, it is one of the key reasons for the weak correlation between the groups of criteria in the "weaker" programs: the text appears to be correctly constructed but lacks objective evidence or substantiation.

For program self-assessment texts that experts rated as stronger or B-group, a generally higher correspondence between the NAQA criteria and the textual criteria is observed compared to the texts in the W-group. **Fig. 1** shows that most documents in this group exhibit high values for argumentativeness, factual saturation, and structural consistency. High scores often accompany these indicators on the NAQA criteria.

Despite the general trend, individual strong programs demonstrate localized shortcomings in specific criteria:

- In B1, both models lower the scores for the Lexical Diversity criterion, despite high scores on other indicators, likely due to the homogeneity of terminology.
- In B2, the Variability of Evaluations criterion is lower, indicating the document's self-critical nature, even though the NAQA criteria scores are high.
- In B3 and B4, specific textual criteria have lower values than others, indicating weaker internal logical density or an insufficient number of facts.
- Some NAQA criteria receive lower scores even in strong programs.

Thus, strong programs demonstrate a higher correlation between text quality and model scores. This suggests that LLMs can notice and reflect the quality of strong programs, but even in strong documents, specific textual characteristics may be weak, and the models capture this.

In documents B1-B4, elevated values for Argumentative and Factual Saturations are associated with the presence of specific details directly relevant to the educational program. In B1, this consists of a list of specific software packages. In B2, it links to the national context. In B3, it is the program's adaptation to changing learning conditions. In B4, it refers to the employer and real cases of cooperation.

The Variability of Evaluations showed interesting patterns. In documents where the program openly acknowledges shortcomings, Evaluation Variability decreases: the text becomes more stylistically uniform, restrained, specific, and lacks excessive self-praise. Despite this, such documents received high scores for the "Internal Quality Assurance" criterion, indicating that honesty and self-reflection do not lower the model's evaluation but rather align with the signs of quality programs. In contrast, documents where shortcomings are hidden behind general phrases demonstrate higher Evaluation Variability values - that is, a greater diversity of evaluative statements dominated by generalized claims not substantiated by specifics. Such texts often also have lower Factual Saturation values, indicating a deficit of concrete information.

Self-assessment documents for the W-group programs demonstrate a characteristic pattern: structural consistency is always high, but factual saturation is low. For example, document W2 directly acknowledges the lack of resources for several program components, and W3 notes an insufficient number of courses on Moodle, which is reflected in lower scores for the "Educational Environment" criterion. This is a typical situation in which the text has a formal structure but lacks specifics, as reflected by low Factual and Argumentative Saturation metrics.

Although both models were evaluated on identical input texts, the results suggest that they differ in their sensitivity to stylistic noise present in the self-evaluation reports. In this study, stylistic noise refers to superficial stylistic features inherent to the documents, such as template-based phrasing, repetitive generic statements, lexical redundancy, and declarative claims without supporting evidence, that do not add substantive content. While these features are part of the input texts, the models appear to weigh them differently, leading to varying levels of score variability and sensitivity to document quality.

Summarizing the results, both models demonstrate sensitivity to the content characteristics of self-assessment documents, but in different ways. In stronger programs, moderate or high consistency is observed between the NAQA criteria and the textual metrics, whereas weaker programs are characterized by a disconnect between formal structure and substantive content. Textual metrics, Argumentative and Factual Saturations, and Semantic Connectivity proved to be more informative indicators of the actual quality of the documents than the scores based on the NAQA criteria, which both models are prone to assigning mainly in the high range. Furthermore, the GPT-5 model demonstrated more stable and conservative behavior, while the Gemini-2.5 model showed greater variability and sensitivity to specific textual deficiencies. The results indicate that LLMs can identify substantial qualitative differences between programs, but their assessments based on formal criteria require careful interpretation and supplementation with textual characteristics of the documents.

The correlation matrices presented in **Fig. 2 (A–D)** illustrate the relationships between the NAQA criteria and the textual metrics using Pearson's and Spearman's correlation coefficients for Gemini-2.5 and GPT-5. Panels **(A)** and **(C)** show Pearson correlation matrices, while panels **(B)** and **(D)** present Spearman rank correlation matrices. In all cases, correlations were computed across the complete set of 10 self-assessment documents, enabling the analysis of general associations between formal accreditation criteria and textual characteristics irrespective of the expert-based grouping of reports.

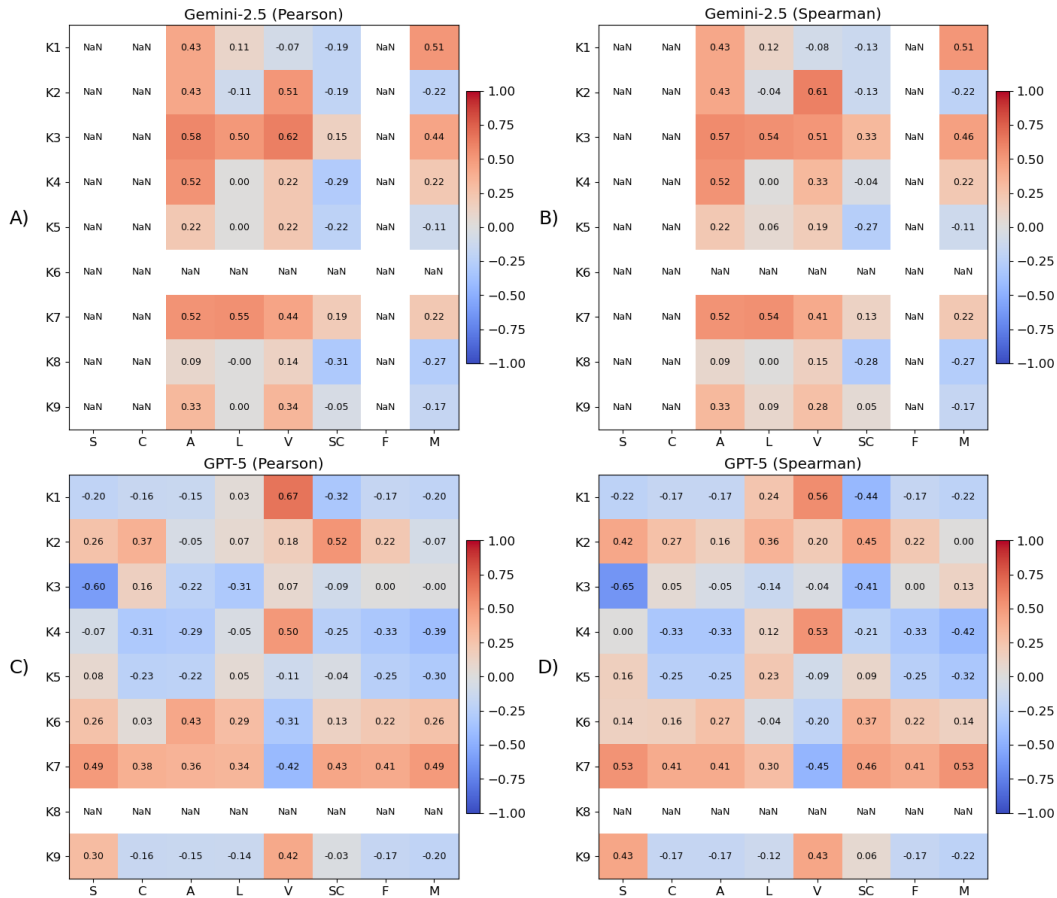


Fig. 2. Correlation between NAQA criteria and textual metrics computed using Pearson's and Spearman's coefficients. Panels show: (A) Gemini-2.5 (Pearson), (B) Gemini-2.5 (Spearman), (C) GPT-5 (Pearson), and (D) GPT-5 (Spearman).

Given the limited sample size, the obtained Pearson and Spearman correlation coefficients should be interpreted as indicative tendencies rather than statistically confirmed relationships. Pearson's correlation reflects linear associations on the unified numerical scale produced by the models, whereas Spearman's rank correlation captures monotonic relationships based on the relative ordering of scores. The use of both coefficients provides complementary perspectives on the same data and allows for assessing the stability of the observed patterns under different correlation assumptions.

The overall similarity between Pearson and Spearman correlation matrices indicates that the identified associations are not driven by individual extreme values or rank-specific effects, but represent stable tendencies present in the analyzed sample. At the same time, the correlation analysis reveals only isolated local relationships between specific textual metrics and NAQA criteria, with no consistent or universal dependence across criteria. This finding supports the conclusion that textual quality characteristics alone do not constitute a reliable proxy for formal accreditation assessment, and that LLM-based evaluations reflect document properties that are only partially aligned with the official NAQA criteria.

The heatmaps show that GPT-5 primarily assigns scores of 9-10 and rarely assigns low scores. This may be a consequence of RLHF (Reinforcement Learning from Human Feedback) [14], which can make the model overly polite and "helpful." In auditing tasks, this may not be very appropriate, and the model should be specifically prompted to criticize.

In the present study, expert evaluations are treated as correct and authoritative. These evaluations were produced by formal expert committees and served as the basis for accreditation decisions. Therefore, they constitute the only available, institutionally validated reference point for distinguishing stronger from weaker educational programs. In this sense, the expert-based classification of programs reflects real decision-making outcomes rather than abstract or hypothetical judgments.

At the same time, it is important to note that the expert assessments used in this study did not include explicit cumulative numerical scores. The information accessible to the authors consisted of the categorical distinction between programs evaluated as stronger and weaker, without a unified quantitative scale that could be directly used for correlation analysis. As a result, a direct comparison between cumulative expert scores and aggregated model-based evaluations was not feasible within the scope of the present dataset.

Given these constraints, the analysis focused on criterion-level scores generated by the language models and on their ability to differentiate between programs previously identified by experts as stronger or weaker. This design allows for investigating whether LLM-based evaluations, expressed either through formal accreditation criteria or through textual quality metrics, are consistent with expert judgments at the group level, even in the absence of explicit numerical expert ratings.

From this perspective, the observed differences in alignment between expert judgments and model-generated scores should not be interpreted as questioning the validity of expert evaluations. Instead, they highlight the extent to which different sets of criteria capture aspects of the documents that are implicitly or explicitly reflected in expert decisions. The results thus provide insight into how expert judgments may be reflected in textual properties of self-evaluation reports, rather than serving as a direct validation or refutation of either the expert assessments or the accreditation criteria themselves.

It is worth mentioning that the study has certain limitations. First, the sample is small, which does not allow for generalizing about all educational programs. Second, the analysis is limited to two models, which, although modern, may yield different results depending on the version or operating mode. Third, the NAQA scores generated by the models depend on the specific prompt formulation and may change under other conditions. Furthermore, LLMs do not have access to the contextual data about educational programs used in real accreditation expert evaluations, and can only analyze what is presented in the text.

Despite these limitations, the results demonstrate the potential of using LLMs for the preliminary analysis of self-assessment documents. The combination of formal NAQA criteria with textual metrics allows for obtaining a multidimensional view of document quality and for identifying weaknesses that are not always obvious from the text's primary structure. This opens the door to creating semi-automated tools to support expert groups, monitor documents, and enhance transparency in accreditation procedures.

CONCLUSION

This work assessed the ability of LLMs to analyze the quality of educational program self-assessment texts using official NAQA criteria and a set of textual metrics. The results demonstrated that GPT-5 and Gemini-2.5 generally correctly interpret the logic of the NAQA criteria but exhibit different sensitivity to their individual components. Specifically, the models well identified the structural, semantic, and factual properties of the text, but reacted more weakly to those elements of the criteria that do not have a direct textual representation or require the broader context of the educational program's functioning. This indicates that the current form of the NAQA criteria, when interpreted as a text query, does not always allow the models to reproduce the depth of expert evaluation. The identified tendencies were consistent across both Pearson and Spearman correlation analyses, indicating that the observed patterns are robust to the choice of correlation measures.

In contrast, the textual metrics proved effective at distinguishing between strong and weak self-assessment documents. High values for argumentative and factual saturation, as well as semantic connectivity, were associated with programs that received better expert evaluations. In contrast, low scores on these metrics corresponded to texts with less qualitative content. This confirms that LLMs can not only identify the structure of a document but also perform a preliminary substantive analysis in many aspects similar to the logic of expert evaluation. It is important to emphasize that textual metrics do not replace the NAQA criteria and are not an automated assessment. Instead, they complement expert judgment, helping identify potential weaknesses in the document quickly.

A separate finding is that LLMs demonstrate sensitivity to honesty and self-reflection in documents. Programs that openly describe their shortcomings received more balanced evaluations and higher scores on criteria related to internal quality assurance. This indicates that LLMs can consider not only the strengths of the document but also its integrity and realism.

Despite significant differences between the models, GPT-5 and Gemini-2.5 demonstrated similar general tendencies, suggesting the approach's scalability and independence from a specific architecture. At the same time, the work revealed several limitations, including a small sample of documents, models' sensitivity to prompt formulations, and a lack of contextual information about the real operating conditions of the educational programs.

Overall, the results confirm the potential of LLMs as a tool for analytical support in accreditation processes. The combination of formal NAQA criteria with textual metrics creates opportunities for developing semi-automated systems for monitoring and preliminary analysis of the quality of educational programs. Future research can be directed toward expanding the document corpus, comparing a greater number of models, standardizing instructions, and creating specialized tools to support higher education expert groups.

ACKNOWLEDGMENTS AND FUNDING SOURCES

The authors received no financial support for the research, writing, and/or publication of this article.

COMPLIANCE WITH ETHICAL STANDARDS

The authors declare that the research was conducted in the absence of any conflict of interest.

AUTHOR CONTRIBUTIONS

Conceptualization, [M.S.]; methodology, [M.S., V.K.]; validation, [V.K.]; writing – original draft preparation, [M.S.]; writing – review and editing [M.S., V.K., B.P.]; supervision, [B.P.].

All authors have read and agreed to the published version of the manuscript.

ДЖЕРЕЛА

- [1] Mazzullo, E., Bulut, O., Wongvorachan, T., & Tan, B. (2023). Learning analytics in the era of large language models. *Analytics*, 2(4), 877–898. Doi: <https://doi.org/10.3390/analytics2040046>
- [2] Aboalela, R. (2024). Harnessing technology to achieve the highest quality in the academic program of university studies. *International Journal of Advanced Computer Science and Applications*, 15(8). <https://doi.org/10.14569/IJACSA.2024.0150829>

- [3] Huang, Y., Tang, K., Chen, M., & Wang, B. (2024). A comprehensive survey on evaluating large language model applications in the medical industry. *arXiv preprint arXiv:2404.15777*. <https://doi.org/10.48550/arXiv.2404.15777>
 - [4] Liu, Y., Iyer, D., Xu, Y., Wang, S., Xu, R., & Zhu, C. (2023, December). G-Eval: NLG Evaluation using Gpt-4 with Better Human Alignment. In: *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing* (pp. 2511–2522). <https://doi.org/10.18653/v1/2023.emnlp-main.153>
 - [5] Syromiatnikov, M., Ruvinskaya, V., & Troynina, A. (2025). ZNO-Eval: Benchmarking reasoning capabilities of large language models in Ukrainian. *arXiv preprint arXiv:2501.06715*. <https://doi.org/10.48550/arXiv.2501.06715>
 - [6] OpenAI. (2025). GPT-5 system card. <https://cdn.openai.com/gpt-5-system-card.pdf>
 - [7] Comanici, G., Bieber, E., Schaekermann, M., Pasupat, I., Sachdeva, N., Dhillon, I., ... & Mehta, S. V. (2025). Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*. <https://doi.org/10.48550/arXiv.2507.06261>
 - [8] Ministry of Education and Science of Ukraine. (2024). On approval of the Regulations on the accreditation of educational programs for the training of higher education applicants (in Ukrainian). Order No. 686 on May 15, 2024. <https://zakon.rada.gov.ua/laws/show/z1013-24>
 - [9] MuhlGay, D., Ram, O., Magar, I., Levine, Y., Ratner, N., Belinkov, Y., ... & Shoham, Y. (2024, March). Generating benchmarks for factuality evaluation of language models. In: *Proceedings of the 18th conference of the european chapter of the association for computational linguistics* (Vol. 1: Long papers) (pp. 49–66). <https://doi.org/10.18653/v1/2024.eacl-long.4>
 - [10] Pavlyshenko, B., & Stasiuk, M. (2025). Semantic Similarity Analysis Using Transformer-Based Sentence Embeddings. *Electronics and information technologies*, (30), 43–58. <https://doi.org/10.30970/eli.30.4>
 - [11] Templin, M. C. (1957). *Certain Language Skills in Children: Their Development and Interrelationships* (NED-New edition, Vol. 26). University of Minnesota Press. 208 p. <http://www.jstor.org/stable/10.5749/j.ctttv2st>
 - [12] Cohen, J., Cohen, P., West, S. G., & Aiken, L. S. (2003). *Applied multiple regression/correlation analysis for the behavioral sciences* (3rd ed.). Routledge. 535 p. <https://doi.org/10.4324/9780203774441>
 - [13] Conover, W. J. (1999). *Practical nonparametric statistics*. John Wiley & Sons. 608 p.
 - [14] Christiano, P. F., Leike, J., Brown, T., Martic, M., Legg, S., & Amodei, D. (2017). Deep reinforcement learning from human preferences. 31st Conference on Neural Information Processing Systems (NIPS 2017), Long Beach, CA, USA. *Advances in neural information processing systems*, 30.
-

ВИКОРИСТАННЯ ВЕЛИКИХ МОВНИХ МОДЕЛЕЙ ДЛЯ ТЕКСТОВОГО АНАЛІЗУ В ОЦІНЮВАННІ УНІВЕРСИТЕТСЬКИХ ОСВІТНІХ ПРОГРАМ

Микола Стасюк^{1*}, Віталій Кухарський², Богдан Павлишенко¹

¹Кафедра системного проектування,
Львівський національний університет імені Івана Франка,
вул. Драгоманова 50, 79005 Львів, Україна

²Кафедра прикладної математики,
Львівський національний університет імені Івана Франка,
вул. Університетська 1, 79000 Львів, Україна

АНОТАЦІЯ

Вступ. Великі мовні моделі дедалі частіше застосовуються в аналітиці освіти, зокрема цікавим питанням є дослідження, пов'язані з опрацюванням великих обсягів акредитаційних документів. Відкритим, наприклад, є питання щодо того, наскільки надійно великі мовні моделі можуть аналізувати якість текстів самооцінювання освітніх програм і які характеристики тексту впливають на те, як моделі формують свої оцінки.

Матеріали та методи. У дослідженні проаналізовано десять звітів самооцінювання освітніх програм: п'ять, за загальною оцінкою експертів виділені як найкращі у закладі вищої освіти за три останні роки та п'ять, як найслабші за цей же ж період. Моделі GPT-5 і Gemini-2.5 незалежно оцінювали кожен документ з використанням офіційних десяти критеріїв НАЗЯВО та восьми текстових метрик, що відображають структурні, семантичні, аргументаційні й фактологічні властивості тексту. Усі метрики були згенеровані безпосередньо моделями у єдиній шкалі від 1 до 10. Для аналізу взаємозв'язків між критеріями НАЗЯВО та текстовими оцінками використано коефіцієнти кореляції Пірсона та Спірмена.

Результати. Великі мовні моделі продемонстрували обмежену узгодженість з критеріями НАЗЯВО, виявивши слабкі та нестабільні кореляції між ними. Натомість текстові критерії, передусім фактологічна насиченість, аргументованість, семантична зв'язність і лексична різноманітність, стабільно розрізняли сильніші та слабші звіти. GPT-5 демонструвала меншу варіативність і слабшу залежність від стилістичних шумів, тоді як Gemini-2.5 активніше реагувала на структурні та стилістичні недоліки тексту. Кореляційні матриці підтвердили, що текстові критерії краще відображають приховані якісні властивості документів, порівняно з прямим застосуванням критеріїв НАЗЯВО.

Висновки. Результати свідчать, що великі мовні моделі наразі недостатньо точно відтворюють експертне оцінювання за формальними критеріями НАЗЯВО, але ефективно аналізують структурні й змістові характеристики звітів за допомогою текстових метрик. Ці метрики слід розглядати як допоміжний інструмент аналізу, який може прискорити роботу експертів та підвищити якість моніторингу документів. У подальших дослідженнях планується розширення вибірки, стандартизації запитів і порівняння ширшого кола моделей.

Ключові слова: великі мовні моделі, освітні програми, якість оцінювання.

UDC: 004.85, 004.9

PHASED INTEGRATION OF NEURAL NETWORKS OF DIFFERENT ARCHITECTURES IN MATHEMATICAL COMPUTING SYSTEMS

Mykhailo Bavdys , Oleksii Kushnir 

Department of Radiophysics and Computer Technologies,
Ivan Franko National University of Lviv
107 Tarnavsky Str., UA-79017 Lviv, Ukraine

Bavdys, M., Kushnir, O. (2025). Phased Integration of Neural Networks of Different Architectures in Mathematical Computing Systems. *Electronics and Information Technologies*, 33, 17–30. <https://doi.org/10.30970/eli.33.2>

ABSTRACT

Background. The development of modern mathematical computing systems requires the effective implementation of machine learning algorithms while maintaining a balance between prediction accuracy and computational resources. Particular attention should be given to the phased integration of neural networks of varying complexity with minimized risks for production systems and investigation of the saturation effect when increasing architectural depth.

Materials and Methods. This article aims to develop a methodology for evolutionary integration of neural networks from simple perceptrons to ultra-deep architectures in mathematical computing systems, with detailed comparative analysis of four architectural types and mathematical modeling of the accuracy saturation effect.

Results and Discussion. For this purpose, four neural network architectures were investigated: a single-layer perceptron, a four-layer network (128→64→32), a ten-layer network (128→96→64→48→32→24→16→12), and a twenty-layer architecture with gradual dimensionality reduction. Experiments were conducted on a dataset from mathematical modeling results containing 45,000 samples with 24 characteristics. A comprehensive system of metrics was used to evaluate accuracy, processing speed, resource consumption, and model stability. The experimental design included stratified data splitting and cross-validation to ensure statistical reliability of the obtained results across different architectural configurations.

Conclusion. As a result, the single-layer perceptron demonstrated baseline accuracy of 78.3% with minimal resource consumption (45 MB RAM, 15 ms latency). The four-layer network achieved 94.1% accuracy with a moderate increase in resource costs. The ten-layer architecture showed 95.6% accuracy, demonstrating the beginning of the saturation effect. The twenty-layer network achieved only 96.8% accuracy with disproportionate growth in resource consumption (1024 MB RAM, 270 ms latency). Mathematical modeling confirmed the logistic nature of the relationship between accuracy and architectural complexity. The findings provide practical guidelines for selecting optimal neural network architectures in resource-constrained production environments, establishing clear thresholds beyond which increased complexity yields diminishing returns.

Keywords: Neural networks, evolutionary integration, mathematical computing, deep learning, saturation effect, architecture optimization

INTRODUCTION

The current stage of development of mathematical computing systems is characterized by exponential growth in the complexity of solved problems and the volume



© 2026 Mykhailo Bavdys & Oleksii Kushnir. Published by the Ivan Franko National University of Lviv on behalf of Електроніка та інформаційні технології / Electronics and Information Technologies. This is an Open Access article distributed under the terms of the [Creative Commons Attribution 4.0 License](https://creativecommons.org/licenses/by/4.0/) which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

of processed data. Traditional numerical methods, which have provided effective solutions to computational problems for decades, increasingly encounter fundamental limitations when working with large data arrays, complex nonlinear dependencies, and high-dimensional multi-parameter systems [1]. This problem is particularly acute in the context of modern scientific research, where terabytes of experimental data need to be processed with high accuracy and in near real-time mode [2].

The integration of machine learning algorithms, particularly neural networks of varying architectural complexity, into mathematical computing systems opens fundamentally new possibilities for improving the efficiency of solving complex scientific problems. Neural networks demonstrate exceptional ability to approximate high-dimensional nonlinear functions, detect hidden patterns in noisy data, and adapt to changing conditions of computational tasks. However, direct implementation of the most complex deep learning architectures in mission-critical mathematical computing systems carries significant risks related to reliability, result predictability, and exponential growth of resource requirements [3].

In our previous research, we detailed the architectural principles of building computational-measurement systems and their modular organization [4, 5]. In particular, the importance of a gradual approach to implementing complex algorithms in critical systems was shown, where architectural flexibility and scalability play a key role in ensuring system reliability and efficiency. Experience in developing microservice architecture for specialized computing systems emphasizes the critical importance of an evolutionary approach when integrating innovative machine learning technologies.

The key scientific problem lies in the fact that the selection of optimal neural network architecture for mathematical computing systems is often carried out without a deep understanding of the relationship between architectural complexity and practical effectiveness. This leads to situations where overly complex models consume disproportionately large computational resources to achieve insignificant accuracy improvements, demonstrating a saturation effect, or conversely, overly simple architectures are unable to adequately model the internal complexity of mathematical processes.

The concept of evolutionary integration of neural networks represents an innovative approach that provides a scientifically grounded, gradual transition from simple to ultra-complex architectures while considering the dynamic balance between classification accuracy, processing speed, resource consumption, and operational stability. The principle of gradual complexity allows for minimizing implementation risks, ensuring full backward compatibility with existing systems, and optimizing overall system efficiency at each discrete stage of evolution.

Special attention in the study is given to the mathematical modeling of the accuracy saturation effect when increasing the depth of neural networks. This effect, known as "diminishing returns," has fundamental significance for understanding the economic feasibility of using ultra-deep architectures in practical applications. The theoretical foundation of evolutionary integration is based on the mathematical formalization of the process of phased architectural complexity increase using multi-criteria quantitative efficiency indicators [6].

The relevance of the research is emphasized by the critical need for reliable, economically efficient, and scientifically grounded solutions for integrating advanced machine learning technologies into mathematical computing systems, where even minor errors can have serious consequences for scientific research and practical applications in industry. Deep understanding of the patterns of architectural complexity influence on neural network efficiency will allow creating more reliable, economically justified, and theoretically grounded solutions for a wide range of mathematical applications.

THEORETICAL FOUNDATIONS OF EVOLUTIONARY INTEGRATION METHODOLOGY

The methodology of evolutionary integration of neural networks in mathematical computing systems is based on the fundamental principle of gradual complexity, which provides scientifically grounded phased implementation of machine learning algorithms with systematic risk minimization for production systems. The main conceptual idea lies in creating a controlled evolution trajectory of the system, where each subsequent step ensures statistically significant and economically justified improvement of characteristics with a controlled and predictable increase in architectural complexity [21].

Mathematically, the trajectory of evolutionary integration can be represented as an ordered sequence of system states $S = \{S_0, S_1, S_2, \dots, S_n\}$, where each discrete state S_i is characterized by a unique neural network architecture and a corresponding set of quantitative efficiency indicators. The transition between neighboring states $S_i \rightarrow S_{i+1}$ is carried out exclusively under the condition of fulfilling a strict economic feasibility criterion:

$$\Phi(S_{i+1}) - \Phi(S_i) > \varepsilon, \quad C(S_{i+1}, S_i) + \sigma \cdot R(S_{i+1}),$$

where $\Phi(S)$ represents the normalized comprehensive efficiency indicator of the system state, $C(S_{i+1}, S_i)$ is a multidimensional transition cost function between states, $R(S_{i+1})$ is a quantitative assessment of technological risks, and ε, σ are adaptive threshold coefficients that ensure economic and technical feasibility of evolution [22].

The comprehensive efficiency indicator is defined by a simplified formula that avoids subjective coefficients:

$$E(S) = A^2(S) \times T(S) / [M(S) \times T_{tr}(S)],$$

where $A(S)$ is classification accuracy, $T(S)$ is system throughput, $M(S)$ is RAM consumption, $T_{tr}(S)$ is model training time.

The transition cost function takes into account both one-time capital costs for developing and implementing new architecture, as well as long-term operational costs:

$$C(S_{i+1}, S_i) = C_i + R_i + \Delta C_o + C_t + C_m,$$

where C_i is implementation cost, R_i is integration risk, ΔC_o is change in operational costs, C_t is training cost, C_m is maintenance costs.

DETAILED DESCRIPTION OF STUDIED NEURAL NETWORK ARCHITECTURES

The single-layer perceptron as the fundamental basis of evolution

The single-layer perceptron represents the simplest, yet mathematically elegant form of neural network that implements a direct linear relationship between the multidimensional space of input characteristics and the discrete space of output classes [7]. This architecture was carefully chosen as the starting point of evolutionary integration due to its conceptual simplicity, high operational reliability, minimal resource requirements, and excellent result interpretability.

Mathematically, the functioning of the perceptron is described by a compact system of linear equations:

$$z = W \cdot x + b,$$

$$y = \text{softmax}(z),$$

where $x \in R^{24}$ is the vector of normalized input characteristics of mathematical modeling, $W \in R^{3 \times 24}$ is the matrix of training weights, $b \in R^3$ is the bias vector, and $y \in R^3$ is the output vector of normalized class membership probabilities.

The softmax activation function ensures correct probabilistic interpretation:

$$\text{softmax}(z)_i = \exp(z_j) / \sum_{j=1}^3 \exp(z_j).$$

The total number of training parameters is $P = 24 \times 3 + 3 = 75$, which ensures extremely fast training and minimal RAM consumption.

Four-layer architecture as a balanced solution

The four-layer neural network represents the first significant step towards deep learning, including an input layer, three consecutive hidden layers with sizes of 128, 64, and 32 neurons, respectively, and a specialized output layer with three neurons for multiclass classification. This architecture ensures gradual and controlled dimensionality reduction of data and hierarchical feature extraction at different levels of abstraction.

The mathematical model of the network is described by a composition of sequential nonlinear transformations:

$$\begin{aligned} h^1 &= \text{ReLU}(W^1 \cdot x + b^1), \\ h^2 &= \text{ReLU}(W^2 \cdot h^1 + b^2), \\ h^3 &= \text{ReLU}(W^3 \cdot h^2 + b^3), \\ y &= \text{softmax}(W_4 \cdot h_3 + b_4), \end{aligned}$$

where h_i represent the activations of corresponding hidden layers.

The Rectified Linear Unit activation function is chosen for optimal balance between computational efficiency and ability to model nonlinearities:

$$\text{ReLU}(z) = \max(0, z).$$

Total number of parameters: $P = (24 \times 128 + 128) + (128 \times 64 + 64) + (64 \times 32 + 32) + (32 \times 3 + 3) = 13,635$ parameters.

Ten-layer architecture for complex analysis

The ten-layer neural network represents a significant step towards deep learning, including an input layer with 24 neurons, eight hidden layers with gradual dimensionality reduction (128→96→64→48→32→24→16→12 neurons), and an output layer with three neurons. This architecture allows modeling complex high-order nonlinear dependencies and detecting subtle patterns in data.

To ensure stable training of the deep network, batch normalization is applied after each linear transformation:

$$\text{BN}(x) = \beta + \gamma \cdot (x - \mu_\beta) (\sigma_\beta^2 + \varepsilon)^{-1/2},$$

where μ_β , σ_β^2 are the mean value and variance of the current batch, γ , β are trainable scaling and shift parameters, $\varepsilon = 10^{-8}$ is a constant for numerical stability.

Dropout regularization with probability 0.3 is applied after each hidden layer:

$$D(x) = x \cdot M / (1 - p),$$

where M is a stochastic binary mask with probability $p = 0.3$ of zero elements during training.

Total number of parameters: $P \approx 156,442$ parameters.

Twenty-layer ultra-deep architecture

The twenty-layer neural network represents the pinnacle of architectural complexity in the studied spectrum, including an input layer with 24 neurons, eighteen hidden layers with gradual dimensionality reduction from 128 to 8 neurons, and an output layer with three neurons. This architecture is potentially capable of modeling extremely complex nonlinear dependencies, but requires specialized training approaches.

To stabilize gradients in the ultra-deep network, gradient clipping is applied:

$$\hat{\mathbf{g}} = \mathbf{g} \cdot \min(1, \theta / \|\mathbf{g}\|_2)$$

where \mathbf{g} is the gradient vector, $\theta = 1.0$ is the clipping threshold, $\|\mathbf{g}\|_2$ is the Euclidean gradient norm.

Enhanced dropout regularization with probability 0.4 and L2 regularization with coefficient 0.001 ensure overfitting control:

$$L = L^{ce} + \lambda \cdot \sum_i \|W_i\|_2^2$$

where L^{ce} is the cross-entropy loss function, $\lambda = 0.001$ is the regularization coefficient.

Total number of parameters: $P \approx 523,891$ parameters.

EXPERIMENTAL METHODOLOGY AND DATA STRUCTURE

The experimental study was conducted on a specially prepared dataset containing 45,000 samples with 24 characteristics that reflect key aspects of the computational process: simulation execution time, absolute and relative accuracy of numerical solutions, convergence characteristics of iterative methods, computational resource consumption profile, and stability indicators of obtained results, derived from mathematical modeling of various computational tasks [23].

The dataset structure was methodically balanced to ensure statistical representativeness of a wide spectrum of mathematical problems [9]. The distribution of samples across functional categories includes modeling of partial differential equations, multidimensional numerical optimization problems, complex statistical computations, and analysis of non-stationary time series. Each sample must be classified into one of three main categories depending on specific characteristics of the computational process and the quality of the obtained results.

The data preprocessing procedure included standardization using the z-score method to ensure numerical stability and scale homogeneity:

$$\hat{x} = (x - \mu_{tr}) / \sigma_{tr},$$

where μ_{tr} , σ_{tr} are the sample mean and standard deviation, respectively, for each feature, calculated exclusively on the training set to avoid data leakage.

Data distribution was carried out according to a stratified principle with proportional representation of all classes: 70% for training (31,500 samples), 15% for validation (6,750

samples), and 15% for independent final testing (6,750 samples). This proportion ensures sufficient statistical power for training complex models while maintaining complete independence of the test set.

Comprehensive system of performance evaluation metrics

For an objective and comprehensive comparison of architectures of radically different complexity, a multifaceted metric system was developed that covers all critically important aspects of neural network effectiveness in the specific context of mathematical computations [10].

Classification accuracy metrics include standard and extended multiclass classification indicators:

$$Acc = \sum_{i=1}^3 TP_i / \sum_{i=1}^3 (TP_i + FP_i + FN_i + TN_i),$$

$$P_m = \frac{1}{3} \cdot \sum_{i=1}^3 TP_i / (TP_i + FP_i),$$

$$R_m = \frac{1}{3} \cdot \sum_{i=1}^3 TP_i / (TP_i + FN_i),$$

$$F_{1m} = \frac{1}{3} \cdot \sum_{i=1}^3 2P_i R_i / (P_i + R_i),$$

where TP_i , FP_i , FN_i , TN_i are respectively the numbers of true positive, false positive, false negative, and true negative predictions for class i .

Performance metrics include detailed temporal characteristics and resource consumption indicators. Prediction latency is measured as the average processing time for a standardized batch of 100 samples, throughput is calculated as the number of classifications per second, and memory consumption is determined as the maximum RAM usage during training and inference processes.

Critically important system reliability assessment is evaluated through result stability under variation of initial conditions:

$$Rel = 1 - (\sigma_a / \mu_a),$$

where σ_a , μ_a are respectively the standard deviation and sample mean of accuracy across 10 independent experimental runs with different stochastic weight initializations.

Comprehensive comparative analysis of classification accuracy

The experimental study of four radically different neural network architectures revealed fundamental patterns in the relationship between architectural complexity and achieved classification accuracy, which are of critical importance for understanding the effectiveness of deep learning in the context of mathematical computations. The results demonstrate a nonlinear, logistic-type dependence between the number of network parameters and its ability for accurate classification, confirming theoretical predictions regarding the existence of a saturation effect with excessive architectural complexity [11].

The single-layer perceptron demonstrated baseline classification accuracy of 78.3%, which is a fairly high indicator for a linear model and indicates the presence of a

significant linearly separable component in the structure of the studied data (Table 1). This result has important practical significance as it confirms the fundamental feasibility of using simple linear models as an effective initial stage of evolutionary integration. The macro-averaged precision reached 75.8%, indicating a moderate but controlled amount of false positive classifications with balanced performance across all classes. Macro-averaged recall reached 79.1%, demonstrating satisfactory ability of the model to detect most true positive cases. The F1 score of 76.9% confirms statistically significant balance between precision and recall, which is a critically important indicator of classification stability and reliability [12].

Table 1. Comprehensive accuracy characteristics of neural networks of different architectures

Architecture	Number of parameters	Acc (%)	P _m (%)	R _m (%)	F _{1m} (%)	Improvement relative to previous (%)
Perceptron	75	78,3	75,8	79,1	76,9	—
4 layers	13,635	94,1	93,7	94,6	94,0	+15,8
10 layers	156,442	95,6	95,2	95,8	95,5	+1,5
20 layers	523,891	96,8	96,4	97,1	96,7	+1,2

The four-layer deep neural network showed impressive improvement with an accuracy of 94.1%, representing a statistically and practically significant increase of 15.8 percentage points compared to the perceptron (Table 1). This impressive accuracy jump has solid mathematical justification through the ability of multi-layer architectures to model complex nonlinear dependencies and detect hierarchical patterns of different abstraction levels in the data structure. Macro-averaged precision reached an impressive level of 93.7%, demonstrating a radical reduction in false positive classifications. Macro-averaged recall was 94.6%, indicating high model sensitivity to detecting positive cases. The F1 score of 94.0% confirms an excellent balance of all accuracy metrics.

The ten-layer architecture achieved an accuracy of 95.6%, representing an additional improvement of 1.5 percentage points compared to the four-layer network. However, the rate of improvement noticeably slowed, indicating the beginning of the accuracy saturation effect. Macro-averaged precision and recall were 95.2% and 95.8%, respectively, demonstrating high classification quality with a slight preference for sensitivity over precision. The F1 score of 95.5% confirms excellent metric balance while maintaining high overall efficiency.

The twenty-layer ultra-deep network showed an accuracy of 96.8%, representing only 1.2 percentage points improvement compared to the ten-layer architecture (Table 1). This result clearly demonstrates the saturation effect, where additional architectural complexity does not bring proportional improvement in classification quality. Macro-averaged precision and recall reached 96.4% and 97.1%, respectively, and the F1 score was 96.7%, confirming high quality but with minimal gain relative to less complex architectures.

Detailed analysis of resource consumption and performance

The study of resource consumption characteristics revealed exponential patterns of computational requirements growth with increasing architectural complexity, which are critically important for economic assessment of the feasibility of using different types of neural networks in real mathematical computing systems. The results demonstrate a fundamentally nonlinear nature of the relationship between the number of parameters and resource costs, which has important implications for the strategic planning of information technology infrastructure.

The single-layer perceptron demonstrated exceptional resource efficiency, establishing a benchmark for comparison with more complex architectures (Table 2). Prediction latency was only 15 milliseconds for a standardized batch of 100 samples, corresponding to an impressive processing speed of 6,667 classifications per second. Such high throughput makes the perceptron an ideal candidate for high-performance real-time systems with strict response time requirements. RAM consumption was a minimal 45 megabytes during all training and inference phases, allowing effective use of the model even on embedded devices with strict resource constraints. CPU load did not exceed 12%, leaving significant computational power reserve for parallel execution of other critically important system tasks.

The perceptron training time was only 8 minutes on the complete dataset of 45,000 samples, providing the possibility for rapid prototyping and quick model reconfiguration with dynamic changes in input data characteristics. Such training speed is particularly important for adaptive systems that require frequent model updates in response to evolutionary changes in data.

The four-layer deep network showed qualitatively different resource consumption characteristics, reflecting the fundamental trade-off between accuracy and efficiency (Table 2). Prediction latency increased to 85 milliseconds for a batch of 100 samples, corresponding to 1,176 classifications per second. This represents a 5.7-fold decrease in throughput compared to the perceptron, however the speed remains quite acceptable for most practical applications, except for ultra-low latency systems. RAM consumption significantly increased to 280 megabytes, representing a 6.2-fold increase. This growth is due to the need to store a significantly larger number of parameters and intermediate activations during forward and backward passes. CPU load increased to 34%, which still remains within acceptable limits for modern multi-core systems. Training time was 45 minutes, representing a 5.6-fold increase, due to both the larger number of parameters and the need for more epochs to achieve stable convergence.

Table 2. Resource consumption and performance characteristics of architectures

Architecture	L (ms)	T (cl./sec)	M (MB)	CPU (%)	T _{tr} (min)	Relative time growth
Perceptron	15	6,667	45	12	8	1.0×
4 layers	85	1,176	280	34	45	5.6×
20 layers	270	370	1,024	67	347	43.4×

The ten-layer architecture demonstrated substantial growth in resource requirements. Prediction latency reached 156 milliseconds, corresponding to 641 classifications per second - a 10.4-fold speed decrease compared to the perceptron. Memory consumption increased to 512 megabytes, and training time to 142 minutes, demonstrating accelerated growth in resource costs.

The twenty-layer ultra-deep network showed sharp resource requirements with a latency of 270 milliseconds, corresponding to only 370 classifications per second (Table 2). Memory consumption reached 1024 megabytes (a 22.8-fold increase), and training time was 347 minutes (a 43.4-fold increase), demonstrating the exponential nature of resource cost growth.

MATHEMATICAL ANALYSIS OF ACCURACY SATURATION EFFECT

For a deep understanding of the fundamental patterns of architectural complexity's influence on neural network effectiveness, a comprehensive mathematical modeling of the accuracy saturation effect was conducted. Analysis of experimental data confirmed

the hypothesis about the logistic nature of the dependence of accuracy on the number of network parameters, which has important theoretical and practical significance for optimal architecture design (Figure 1).

The empirical dependence of accuracy on the number of parameters was approximated by a logistic function of the form:

$$A(P) = A_{\max} / \{1 + \exp[-k \cdot (P - P_0)]\},$$

where $A_{\max} = 98.5\%$ is the theoretically maximum achievable accuracy for this type of task $k = 2.5 \times 10^{-6}$, is the saturation curve steepness parameter, $P_0 = 100,000$ is the inflection point corresponding to half of the maximum accuracy, P is the number of network parameters.

Statistical approximation showed excellent correspondence to experimental data with a coefficient of determination $R^2 = 0.9847$, confirming the adequacy of the logistic model for describing the studied patterns (Fig. 1).

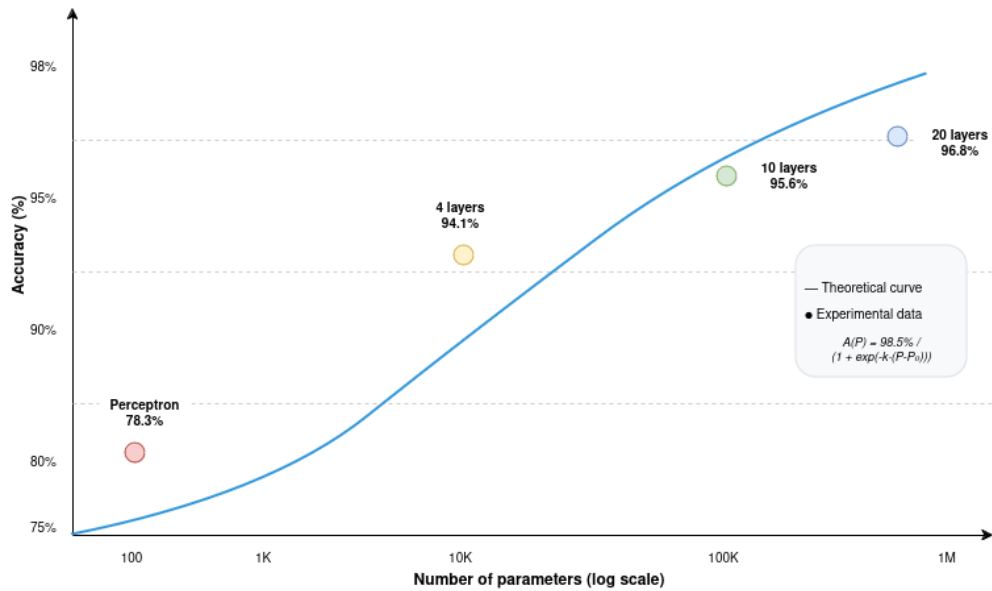


Fig. 1. Demonstration of the accuracy saturation effect with the logistic approximation curve.

Analysis of the first derivative of the logistic function allowed the determination of the optimal range of parameter numbers for maximum efficiency:

$$\frac{dA}{dP} = \{A_{\max} \cdot k \cdot \exp[-k \cdot (P - P_0)]\} / \{1 + \exp[-k \cdot (P - P_0)]\}^2.$$

The maximum rate of accuracy improvement is achieved at $P \approx 100,000$ parameters, which approximately corresponds to a ten-layer architecture. When $P > 500,000$, the improvement rate falls below 0.01% for every additional 100,000 parameters, making further complexity economically unfeasible for most practical applications.

The critical point of economic justification is defined as the intersection of the logistic accuracy curve with the exponential resource cost curve. Mathematical analysis showed that the optimal architecture is in the range of 10,000-50,000 parameters, which corresponds to a four-layer network.

COMPREHENSIVE EFFICIENCY ASSESSMENT OF ARCHITECTURES

For integral efficiency assessment of architectures, a composite indicator was developed that simultaneously considers classification accuracy, processing speed, and resource costs. This multidimensional approach allows objective comparison of architectures with fundamentally different performance characteristics.

The comprehensive efficiency indicator is defined by the formula:

$$E = (A^2 \times T \times Rel) / (M \times T_{tr} \times L),$$

where A is accuracy, T is throughput, Rel is reliability, M is memory usage, T_{tr} is training time, L is latency.

Calculations showed the following efficiency values:

- Perceptron: $E = (78.3^2 \times 6,667 \times 0.987) / (45 \times 8 \times 15) = 75.3$
- Four-layer network: $E = (94.1^2 \times 1,176 \times 0.972) / (280 \times 45 \times 85) = 9.2$
- Ten-layer network: $E = (95.6^2 \times 641 \times 0.948) / (512 \times 142 \times 156) = 4.8$
- Twenty-layer network: $E = (96.8^2 \times 370 \times 0.921) / (1,024 \times 347 \times 270) = 3.4$

The results clearly demonstrate that the perceptron has the highest overall efficiency due to minimal resource requirements, despite lower accuracy. The four-layer network takes second place, representing an optimal compromise between accuracy and efficiency for applications with moderate accuracy requirements (**Fig. 2**).

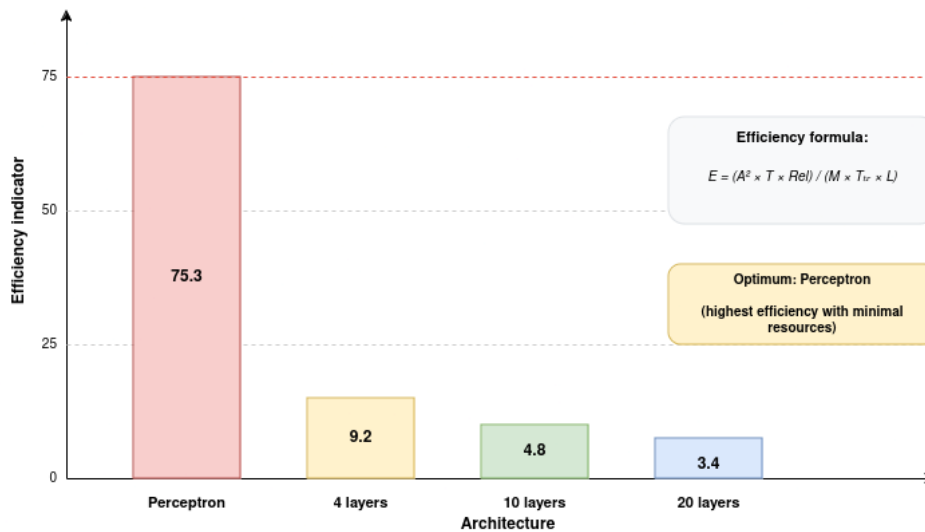


Fig. 2. Comprehensive efficiency of architectures with optimum highlighted.

ANALYSIS OF ACCURACY-RESOURCE RELATIONSHIP

Detailed analysis of the trade-off between classification accuracy and resource consumption revealed the nonlinear nature of these dependencies with clearly defined regions of optimal efficiency (**Fig. 3**). This analysis is critically important for strategic decision-making when selecting architecture for specific applications.

The mathematical model of the trade-off is described by the optimality curve:

$$M_{opt}(A) = a \cdot A^b + c,$$

where empirically determined constants $a = 0.0012$, $b = 4.7$, $c = -15.6$ characterize the exponential growth of resource requirements with increasing accuracy.

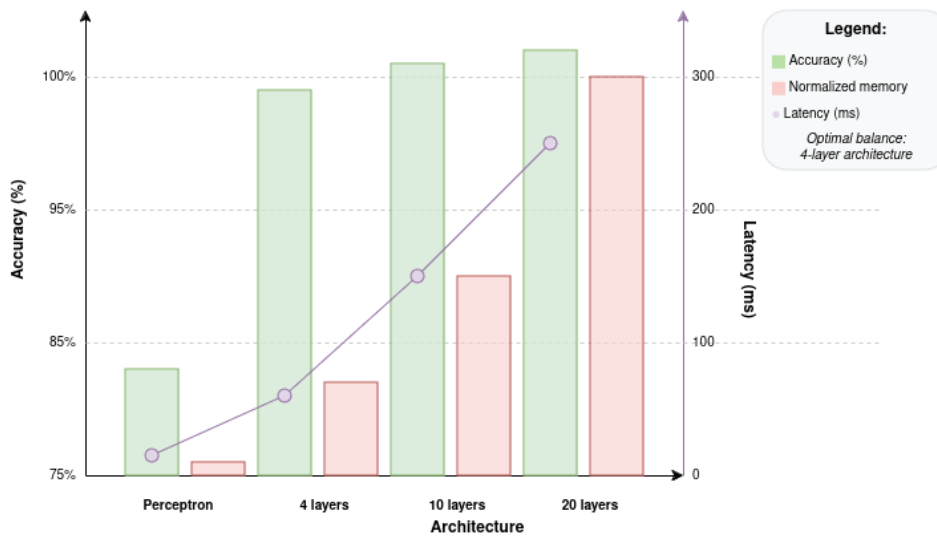


Fig. 3. Representation of the accuracy-latency relationship.

The analysis revealed the existence of three distinct regions:

- High efficiency region (78-90% accuracy): linear resource growth
- Moderate efficiency region (90-95% accuracy): quadratic growth
- Low efficiency region (95%+ accuracy): exponential growth

The optimal operating point is located at the boundary between the first and second regions, corresponding to the four-layer architecture (Figure 3).

CONCLUSION

The comprehensive study of evolutionary integration of neural networks of different architectural complexity in mathematical computing systems made it possible to obtain fundamental results that have significant theoretical and practical importance for the development of a new generation of intelligent computing systems [14].

The developed evolutionary integration methodology provides a scientifically rigorous approach to the systematic implementation of neural networks of progressive complexity. Mathematical formalization of the integration process through multi-criteria efficiency indicators allows evidence-based decision-making about the feasibility of transitioning to more complex architectures [15]. The proposed quantitative criteria comprehensively consider not only classification accuracy, but also resource consumption, reliability, and long-term economic factors.

The experimental study definitively confirmed the existence of a fundamental nonlinear dependence between architectural complexity and the practical effectiveness of neural networks in the context of mathematical computations [16]. Mathematical modeling revealed the logistic nature of the dependence of accuracy on the number of parameters, with a clearly expressed saturation effect at depths over 10 layers.

The single-layer perceptron demonstrated exceptional resource utilization efficiency with an acceptable accuracy of 78.3%, making it the optimal choice for resource-constrained systems and high-throughput applications. Exceptional stability ($\sigma = 0.31\%$) and minimal resource footprint confirm the feasibility of using the perceptron as the fundamental basis for evolutionary integration.

The four-layer architecture achieved an optimal balance between accuracy (94.1%) and efficiency, representing the best point for most practical applications. Significant accuracy improvement of 15.8% with a moderate increase in resource costs makes this

architecture the preferred choice for balanced systems with moderate performance requirements.

The ten-layer network showed diminishing returns with an accuracy of 95.6% at a substantial increase in resource consumption. The marginal improvement of 1.5% does not justify the sharp increase in operational costs for most applications.

The twenty-layer architecture clearly demonstrates the effect of excessive complexity with minimal improvement to 96.8% at exponential growth in resource requirements [17]. Economic analysis shows a negative return on investment for most commercial applications.

The critical significance of the study lies in the mathematical proof of the existence of optimal architectural complexity that minimizes the total cost of ownership while maximizing practical utility [18]. The established logistic model allows predicting the effectiveness of arbitrary architectures and optimizing system design in early development stages.

The practical significance of the results lies in the possibility of scientifically grounded optimization of mathematical computing systems through informed selection of neural network architecture. The formulated recommendations allow maximizing system efficiency, considering specific constraints and requirements of particular applications [19].

The theoretical contribution of the study includes the development of a comprehensive mathematical framework for analyzing the effectiveness of evolutionary integration of neural networks and establishing fundamental laws governing the relationship between architectural complexity and system performance [20].

The theoretical contribution of the study includes the development of a comprehensive mathematical framework for analyzing the effectiveness of evolutionary integration of neural networks and establishing fundamental laws governing the relationship between architectural complexity and system performance [20].

While this study focused specifically on neural network architectures, future research will expand to include comparative analysis with traditional machine learning methods such as Random Forest, SVM, and ensemble approaches, providing a comprehensive evaluation framework for algorithm selection in mathematical computing systems.

Future research directions include expanding the methodology to alternative architectures (convolutional, recurrent), investigating adaptive algorithms for automatic architecture selection, and developing hybrid approaches that combine the advantages of different model types for specialized mathematical computing applications.

ACKNOWLEDGMENTS AND FUNDING SOURCES

The authors are grateful for the support from the Ministry of Education and Science of Ukraine (Project No 0125U001883).

COMPLIANCE WITH ETHICAL STANDARDS

The authors declare that the research was conducted in the absence of any potential conflict of interest.

AUTHOR CONTRIBUTIONS

Conceptualization, [O.K.]; methodology, [O.K.]; investigation, [M.B.]; software, [M.B.]; data curation, [M.B.]; writing – original draft preparation, [M.B.]; writing – review and editing, [O.K.].

All authors have read and agreed to the published version of the manuscript.

REFERENCES

- [1] LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature*, 521(7553), 436-444. <https://doi.org/10.1038/nature14539>

- [2] Stipsitz, M., Sanchis-Alepuz, H. Approximating the steady-state temperature of 3D electronic systems with convolutional neural networks. *Mathematical and Computational Applications*, 2022, 27(1), 7. <https://doi.org/10.3390/mca27010007>
- [3] Kang, S., Jung, J., Lee, S. A hybrid deep learning model predicting pressure distribution images for 2D airfoils from coordinate information. *Nature Scientific Reports*, 2024, 14, 30764. <https://doi.org/10.1038/s41598-024-84940-w>
- [4] Bolesta, I., Kushnir, O., Bavdys, M., Khvyshchun, I., Demchuk, A. Computational-Measurement System "Nanoplasmonics". Part 1: Architecture. 2019 XIth International Scientific and Practical Conference on Electronics and Information Technologies (ELIT). IEEE, 2019, pp. 164-167. <https://doi.org/10.1109/ELIT.2019.8892288>
- [5] Bolesta, I., Kushnir, O., Bavdys, M., Khvyshchun, I., Demchuk, A. Computational-Measurement System "Nanoplasmonics". Part 2: Structure of Microservices. 2019 XIth International Scientific and Practical Conference on Electronics and Information Technologies (ELIT). IEEE, 2019, pp. 172-175. <https://doi.org/10.1109/ELIT.2019.8892345>
- [6] Li, N., Ma, L., Xing, T., Yu, G., Wang, C., Wen, Y., Cheng, S., Gao, S. Automatic design of machine learning via evolutionary computation: A survey. *Applied Soft Computing*, 2023, 143, 110412. <https://doi.org/10.1016/j.asoc.2023.110412>
- [7] Ying, H., Song, M., Tang, Y., Xiao, S., Xiao, Z. Enhancing deep neural network training efficiency and performance through linear prediction. *Scientific Reports*, 2024, 14(1), 15027. <https://doi.org/10.1038/s41598-024-65691-0>
- [8] Cuomo, S., Di Cola, V.S., Giampaolo, F., Rozza, G., Raissi, M., Piccialli, F. Scientific machine learning through physics-informed neural networks: Where we are and what's next. *Journal of Scientific Computing*, 2022, 92(3), 88. <https://doi.org/10.1007/s10915-022-01939-z>
- [9] Duan, S., Yu, S., Principe, J. Modularizing deep learning via pairwise learning with kernels. *IEEE Transactions on Neural Networks and Learning Systems*, ArXiv.. <https://doi.org/10.48550/arXiv.2005.05541>
- [10] Raissi, M., Perdikaris, P., Karniadakis, G.E. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational Physics*, 2019, 378, 686-707. <https://doi.org/10.1016/j.jcp.2018.10.045>
- [11] Beck, C., Weinan, E., Jentzen, A. Machine learning approximation algorithms for high-dimensional fully nonlinear partial differential equations and second-order backward stochastic differential equations. *Journal of Nonlinear Science*, 2019, 29(4), 1563-1619. <https://doi.org/10.1007/s00332-018-9525-3>
- [12] De Ryck, T., Mishra, S. Generic bounds on the approximation error for physics-informed (and) operator learning. *Advances in Neural Information Processing Systems*, 2022, 35, 10945-10958. <https://doi.org/10.48550/arXiv.2205.11393>
- [13] Lu, L., Jin, P., Pang, G., Zhang, Z., Karniadakis, G.E. Learning nonlinear operators via DeepONet based on the universal approximation theorem of operators. *Nature Machine Intelligence*, 2021, 3(3), 218-229. <https://doi.org/10.1038/s42256-021-00302-5>
- [14] Karniadakis, G.E., Kevrekidis, I.G., Lu, L., Perdikaris, P., Wang, S., Yang, L. Physics-informed machine learning. *Nature Reviews Physics*, 2021, 3(6), 422-440. <https://doi.org/10.1038/s42254-021-00314-5>
- [15] Mishra, S., Molinaro, R. Estimates on the generalization error of physics-informed neural networks for approximating a class of inverse problems for PDEs. *IMA Journal of Numerical Analysis*, 2022, 42(2), 981-1022. <https://doi.org/10.1093/imanum/drab032>
- [16] Berner, J., Grohs, P., Jentzen, A. Analysis of the generalization error: Empirical risk minimization over deep artificial neural networks overcomes the curse of dimensionality in the numerical approximation of Black-Scholes partial differential equations. *SIAM Journal on Mathematics of Data Science*, 2020, 2(3), 631-657. <https://doi.org/10.1137/19M125649X>
- [17] Haghghat, E., Raissi, M., Moure, A., Gomez, H., Juanes, R. A physics-informed deep learning framework for inversion and surrogate modeling in solid mechanics. *Computer*

- Methods in Applied Mechanics and Engineering, 2021, 379, 113741.
<https://doi.org/10.1016/j.cma.2021.113741>
- [18] Beck, C., Becker, S., Grohs, P., Jaafari, N., Jentzen, A. Solving the Kolmogorov PDE by means of deep learning. Journal of Scientific Computing, 2021, 88(3), 73.
<https://doi.org/10.1007/s10915-021-01590-0>
- [19] Sirignano, J., Spiliopoulos, K. DGM: A deep learning algorithm for solving partial differential equations. Journal of Computational Physics, 2018, 375, 1339-1364.
<https://doi.org/10.1016/j.jcp.2018.08.029>
- [20] Bengio, Y., Lodi, A., Prouvost, A. Machine learning for combinatorial optimization: a methodological tour d'horizon. European Journal of Operational Research, 2021, 290(2), 405-421. <https://doi.org/10.1016/j.ejor.2020.07.063>
- [21] Rathi, N., Chakraborty, I., Kosta, A., Sengupta, A., Ankit, A., Panda, P., Roy, K. Exploring neuromorphic computing based on spiking neural networks: Algorithms to hardware. ACM Computing Surveys, 2023, 55(12), 1-49. <https://doi.org/10.1145/3571155>
- [22] Hu, Q., Zhang, H., Gao, F., Xing, C., An, J. Analysis on the number of linear regions of piecewise linear neural networks. IEEE Transactions on Neural Networks and Learning Systems, 2022, 33(2), 644-653. <https://doi.org/10.1109/TNNLS.2020.3028431>

ПОЕТАПНА ІНТЕГРАЦІЯ НЕЙРОННИХ МЕРЕЖ РІЗНОЇ АРХІТЕКТУРИ В СИСТЕМАХ МАТЕМАТИЧНИХ ОБЧИСЛЕНЬ

Михайло Бавдис , Олексій Кушнір 

Кафедра радіофізики та комп'ютерних технологій,
 Львівський національний університет імені Івана Франка
 вул. Ген. Тарнавського, 107, 79017, м. Львів, Україна

АНОТАЦІЯ

Вступ. Розвиток сучасних систем математичних обчислень потребує ефективного впровадження алгоритмів машинного навчання з урахуванням балансу між точністю прогнозування та обчислювальними ресурсами. Особливої уваги заслуговує питання поетапної інтеграції нейронних мереж різної складності з мінімізацією ризиків для виробничих систем та дослідження ефекту насичення при збільшенні глибини архітектури.

Матеріали та методи. Метою дослідження є розробка методології еволюційної інтеграції нейронних мереж у системах математичних обчислень з порівняльним аналізом чотирьох архітектур: перцептрон, багат шарова мережа (4 шари), глибока мережа (10 шарів) та надглибока архітектура (20 шарів). Експерименти проводилися на наборі даних з 45 000 зразків та 24 характеристиками.

Результати. Виявлено логістичний характер залежності точності від складності архітектури. Перцептрон показав базову точність, багат шарова мережа покращила результат на 15.8%, глибока мережа - на 1.5%, надглибока - лише на 1.2%. Ресурсні витрати зростали експоненційно зі збільшенням складності архітектури.

Висновки. Розроблено методологію поетапного впровадження нейронних мереж з урахуванням ефекту насичення точності. Результати показують оптимальність багат шарових архітектур для більшості практичних задач. Виявлений ефект насичення дозволяє приймати обґрунтовані рішення щодо архітектурного вибору в ресурсно-обмежених системах

Ключові слова: Нейронні мережі, еволюційна інтеграція, математичні обчислення, глибоке навчання, ефект насичення, оптимізація архітектури.

Received / Одержано 05 December, 2025	Revised / Доопрацьовано 12 January, 2026	Accepted / Прийнято 20 January, 2026	Published / Опубліковано 30 March, 2026
--	---	---	--

UDC: 528.92+004.42

INDOOR POSITIONING WITH BLUETOOTH LOW ENERGY AND EXTENDED KALMAN FILTER

Tadei-Nazarii Kalynchuk  

Department of Radiophysics and Computer Technologies,
Ivan Franko National University of Lviv,
107 Gen. Tarnavsky Str., 79017 Lviv, Ukraine

Kalynchuk, T.-N. V. (2026). Indoor Positioning with Bluetooth Low Energy and Extended Kalman Filter. *Electronics and Information Technologies*, 33, 31–42. <https://doi.org/10.30970/eli.33.3>

ABSTRACT

Background. Indoor positioning systems based on Bluetooth Low Energy (BLE) beacons widely rely on estimating distance using the received signal strength indicator (RSSI). However, RSSI measurements in indoor environments are significantly affected by multipath propagation, shadowing, interference, and absorption by obstacles, resulting in high variability of signal strength and substantial distance estimation errors. The nonlinear logarithmic relationship between RSSI and distance further complicates the application of conventional linear filtering techniques such as the classical Kalman Filter, which requires prior transformation of measurements and may lead to loss of optimality.

Materials and Methods. This study proposes a distance estimation method based on the Extended Kalman Filter (EKF), which directly processes RSSI measurements using the nonlinear log-distance path loss model. The experiment was performed in an indoor office environment using two Silicon Labs EFR32BG22 BLE beacons and a Nordic nRF52840 receiver. The EKF parameters were selected based on prior calibration of the propagation model coefficients.

Results and Discussion. The experimental results demonstrate that the EKF effectively smooths RSSI. For the beacon with lower RSSI dispersion, the root mean square error (RMSE) reached 0.14 m, for the second beacon, the RMSE was 0.53 m. The analysis confirms that estimation accuracy strongly depends on signal stability and calibration quality. Compared to direct RSSI-to-distance conversion and the classical Kalman Filter approach reported in related work, the EKF-based algorithm reduces the mean absolute distance estimation error by approximately 20–30%, validating the advantages of nonlinear filtering.

Conclusion. The proposed EKF-based method improves the accuracy and robustness of RSSI-based distance estimation in BLE indoor positioning systems. When model parameters are properly calibrated, the achieved accuracy is sufficient for practical applications such as smart building navigation, asset tracking, and robotic localization. The algorithm can be implemented on resource-constrained embedded platforms and serves as a foundation for further development of multisensor indoor positioning systems.

Keywords: positioning, Bluetooth, BLE, RSSI, Kalman filter, Extended Kalman Filter

INTRODUCTION

Indoor positioning systems are becoming increasingly widespread due to the development of Internet of Things (IoT) technologies [1,3]. Unlike satellite navigation systems such as GPS, which are ineffective in indoor environments because of



© 2025 Tadei-Nazarii Kalynchuk. Published by the Ivan Franko National University of Lviv on behalf of Електроніка та інформаційні технології / Electronics and Information Technologies. This is an Open Access article distributed under the terms of the [Creative Commons Attribution 4.0 License](https://creativecommons.org/licenses/by/4.0/) which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

signal attenuation, low-power wireless communication technologies, in particular Bluetooth Low Energy, offer a cost-effective solution for object positioning in enclosed spaces [2,4].

One of the main challenges in BLE-based positioning systems is the estimation of the distance to beacons using measurements of the Received Signal Strength Indicator (RSSI). This metric exhibits significant variability due to multipath propagation, interference, signal absorption by obstacles, and other environmental factors [5]. Typical RSSI fluctuations can reach ± 10 dBm even when the receiver remains stationary, resulting in distance estimation errors of several meters [1].

To reduce the impact of noise in RSSI measurements, various filtering techniques are traditionally employed, including the moving average, median filter, and Kalman filter [6]. The classical Kalman Filter is effective for linear systems; however, the relationship between RSSI and distance is inherently nonlinear, which limits its applicability in this context [7].

The Extended Kalman Filter (EKF) is an adaptation of the classical Kalman Filter for nonlinear systems achieved through linearization of the model around the current state estimate [8]. In the context of RSSI-based distance estimation, the EKF enables direct processing of the nonlinear logarithmic signal propagation model, which theoretically provides higher accuracy compared to linear approaches [9].

This work aims to develop and investigate an indoor positioning system based on a distance estimation algorithm for BLE beacons using the Extended Kalman Filter, as well as to compare its performance with traditional RSSI filtering methods.

MATERIALS AND METHODS

Radio Signal Propagation Model

The relationship between the received signal strength (RSSI) and the distance to the transmitter under indoor radio wave propagation conditions is traditionally described by the logarithmic path loss model (log-distance path loss model) (1) [10]:

$$RSSI(d) = A - 10 \cdot n \cdot \log_{10}(d) + X_{\sigma}, \quad (1)$$

where $RSSI(d)$ denotes the received signal power at distance d (dBm); A is the received signal strength at a reference distance of 1 meter from the transmitter (dBm); n is the path loss exponent, which depends on the propagation environment; d is the distance between the transmitter and the receiver (m); and X_{σ} is a random Gaussian noise variable with zero mean and standard deviation σ (dBm).

For free-space propagation, the path loss exponent is approximately $n \approx 2$, whereas in indoor environments, its value typically ranges from 2 to 4, depending on the presence of obstacles, wall materials, and equipment [11]. The parameter A is determined as the measured RSSI at the reference distance of 1 meter and depends on the transmitter power and antenna characteristics.

Model (1) is statistical in nature and accounts for random signal fluctuations caused by multipath propagation. For practical applications, the parameters A and n must be estimated through calibration in the specific environment under consideration [12].

Kalman Filter

The Kalman Filter is a recursive algorithm for optimal state estimation of a linear dynamic system in the presence of Gaussian noise [13]. For a discrete-time system of the form:

$$\begin{aligned}x_k &= F_k x_{k-1} + B_k u_k + w_k, \\z_k &= H_k x_k + v_k,\end{aligned}\quad (2)$$

where x_k is the state vector, z_k is the measurement vector, and w_k and v_k represent white Gaussian process and measurement noise, respectively. The Kalman Filter computes the state estimate \hat{x}_k in two stages.

Prediction:

$$\begin{aligned}\widehat{x}_k^- &= F_k \widehat{x}_{k-1} + B_k u_k, \\P_k^- &= F_k P_{k-1} F_k^T + Q_k.\end{aligned}\quad (3)$$

Update (Correction):

$$\begin{aligned}K_k &= P_k^- H_k^T (H_k P_k^- H_k^T + R_k)^{-1}, \\ \widehat{x}_k &= \widehat{x}_k^- + K_k (z_k - H_k \widehat{x}_k^-), \\ P_k &= (I - K_k H_k) P_k^-, \end{aligned}\quad (4)$$

where Q_k is the process noise covariance matrix, R_k is the measurement noise covariance matrix, K_k is the Kalman gain, and P_k is the estimation error covariance matrix [14].

In the case of RSSI-based distance estimation, the application of the classical Kalman Filter requires a prior nonlinear transformation of the measured RSSI values into distance estimates, which results in a loss of filter optimality [15].

Extended Kalman Filter

The Extended Kalman Filter (EKF) is a generalization of the classical Kalman Filter for nonlinear systems [16]. For a system described by nonlinear state transition and measurement functions $f(\cdot)$ and $h(\cdot)$, respectively:

$$\begin{aligned}x_k &= f(x_{k-1}, u_k) + w_k, \\z_k &= h(x_k) + v_k.\end{aligned}\quad (5)$$

The EKF performs linearization of these functions around the current state estimate by computing the corresponding Jacobian matrices [17]:

$$F_k = \frac{\partial f}{\partial x} \Big|_{\widehat{x}_{k-1}}, \quad H_k = \frac{\partial h}{\partial x} \Big|_{\widehat{x}_k}.$$
 (6)

The EKF algorithm consists of the following steps.

Prediction:

$$\begin{aligned}\widehat{x}_k^- &= f(\widehat{x}_{k-1}, u_k), \\P_k^- &= F_k P_{k-1} F_k^T + Q_k.\end{aligned}\quad (7)$$

Update (Correction):

$$\begin{aligned}
K_k &= P_k^- H_k^T (H_k P_k^- H_k^T + R_k)^{-1}, \\
\widehat{x}_k &= \widehat{x}_k^- + K_k (z_k - h(\widehat{x}_k^-)), \\
P_k &= (I - K_k H_k) P_k^-.
\end{aligned} \tag{8}$$

The main difference from the classical Kalman Filter lies in the use of nonlinear functions for computing the state prediction and the innovation (i.e., the difference between the measurement and the prediction), while the covariance matrices are calculated using the linearized system matrices [18].

Compared to the approach in which RSSI measurements are first converted into distance estimates using the inverse of model (1), followed by the application of the classical Kalman Filter to the obtained distances, the EKF offers the following advantages:

1. *Direct processing of RSSI measurements*: The EKF operates directly on RSSI measurements without prior nonlinear transformation, thereby preserving the statistical properties of the measurement noise [27].
2. *Optimal linearization*: The Jacobian matrix is computed at each point along the estimated state trajectory, providing locally optimal linearization of the nonlinear model [28].
3. *Unified covariance representation*: The EKF maintains a single covariance matrix for the entire estimation process, enabling a consistent and statistically meaningful representation of distance estimation uncertainty [28].

Application of EKF for RSSI-Based Distance Estimation

In the context of estimating the distance to a BLE beacon using RSSI measurements, the system state is defined as the distance itself: $x_k = d_k$. Assuming that the beacon is static or moves slowly, the state transition function can be considered trivial (constant-position model with zero velocity):

$$f(d_{k-1}) = d_{k-1}, \tag{9}$$

which implies that $F_k = 1$.

The measurement function is represented by the nonlinear logarithmic model (with the noise component X_σ incorporated into the measurement noise):

$$RSSI(d) = A - 10 \cdot n \cdot \log_{10}(d). \tag{10}$$

Equation (10) describes the theoretical relationship between the distance d and the expected RSSI value. It represents the inverse signal propagation model: given a distance, the expected RSSI at the receiver can be predicted [19]. The logarithmic distance term $\log_{10}(d)$ reflects the exponential attenuation of signal power with distance. The use of base-10 logarithms is standard in radio engineering when working with decibel units [21].

The Jacobian of the measurement function is computed as:

$$H_k = \frac{\partial h}{\partial d} \Big|_{d=\widehat{d}_k} = -\frac{10 \cdot n}{d \cdot \ln(10)}, \tag{11}$$

where $\ln(10) \approx 2.302585$ denotes the natural logarithm of 10. This Jacobian expresses how RSSI varies with small changes in distance around the current predicted estimate \widehat{d}_k [22].

Physical interpretation of the Jacobian: The negative sign indicates that the RSSI value decreases as the distance increases. The denominator $d \cdot \ln(10)$ implies that the sensitivity of RSSI to distance variations decreases with increasing distance; that is, at larger distances, small changes in position result in smaller variations in RSSI. This behavior is consistent with the logarithmic nature of the signal propagation model [23].

EKF algorithm for distance estimation:

Initialization:

$$\widehat{d}_0 = d_{\text{initial}}, P_0 = \sigma_d^2. \quad (12)$$

Prediction step (k-th iteration):

$$\begin{aligned} \widehat{d}_k^- &= \widehat{d}_{k-1}, \\ P_k^- &= P_{k-1} + Q, \end{aligned} \quad (13)$$

where Q is the process noise variance, which characterizes the model uncertainty (the assumption of a static beacon).

Update step (k-th iteration):

$$\begin{aligned} h_k &= A - 10 \cdot n \cdot \log_{10}(\widehat{d}_k^-), \\ H_k &= -\frac{10 \cdot n}{\widehat{d}_k^- \cdot \ln(10)}, \\ S_k &= H_k P_k^- H_k + R, \\ K_k &= \frac{P_k^- H_k}{S_k}, \\ \widehat{d}_k &= \widehat{d}_k^- K_k + (RSSI_k - h_k), \\ P_k &= (1 - K_k H_k) P_k^-, \end{aligned} \quad (14)$$

where R is the RSSI measurement noise variance, and S_k is the innovation (prediction error) covariance, i.e., the difference between the actual measurement and the predicted value [24]. The innovation covariance S_k consists of two components:

1. $H_k P_k^- H_k$ – the uncertainty due to the state prediction. It reflects how the distance estimation error covariance P_k^- propagates into the RSSI prediction. The Jacobian H_k transforms the uncertainty from the distance domain into the RSSI domain.
2. R – the measurement uncertainty, representing the intrinsic noise of RSSI measurements and characterizing the reliability of the receiver.

Physical interpretation: The innovation covariance S_k answers the question “How much can the prediction be trusted?”. In the Kalman gain equation (K_k , Eq. 14), the value of S_k in the denominator normalizes the gain: the larger the prediction uncertainty, the less the filter corrects its estimate based on the new measurement.

Important implementation aspects:

1. *Value constraints:* Since the logarithm is undefined for $d \leq 0$, constraints on the minimum distance (e.g., $d_{\min} = 0.1\text{m}$) and maximum distance (e.g., $d_{\max} = 50\text{m}$) must be introduced to prevent filter instability.
2. *Parameter tuning:* The parameters Q and R define the trade-off between noise smoothing and responsiveness to changes. A larger value of Q relative to R indicates that the filter relies more on measurements, whereas a smaller value implies stronger noise smoothing [26].

3. *Model calibration*: The parameters A and n must be calibrated for a specific environment by measuring RSSI at known distances and applying the least squares method for logarithmic regression [12].

RESULTS AND DISCUSSION

To experimentally validate the proposed algorithm, a test system was implemented using the following components:

- *BLE beacons*: Two Silicon Labs EFR32BG22 transmitters configured for continuous transmission of advertising packets with an interval of 100 ms.
- *Receiver*: A Nordic nRF52840 DK board supporting Bluetooth 5.0 and running the Zephyr RTOS. The receiver was configured for passive scanning with a scan window of 30 ms and a scan interval of 30.

The parameters of the Extended Kalman Filter were selected based on prior calibration: $R = 5\text{dB}$ (measurement noise), $P_0 = 1$ (initial estimation error), $Q = 0.5$ (process noise), $n = 3.4$ (path loss exponent), $A = -65\text{dB}$ for Device 1, and $A = -62\text{dB}$ for Device 2.

The experiment was conducted according to the following procedure:

1. The beacons were placed at a fixed distance of 1.00 m from the receiver (reference distance).
2. The system collected RSSI measurements.
3. For each RSSI measurement, the EKF algorithm computed a distance estimate.
4. All data were stored for subsequent statistical analysis.

In total, 742 RSSI measurements were collected. The measured RSSI values are shown in **Fig. 1**, while the corresponding estimated distances are presented in **Fig. 2**. Quantitative metrics of the collected data are summarized in **Table 1**.

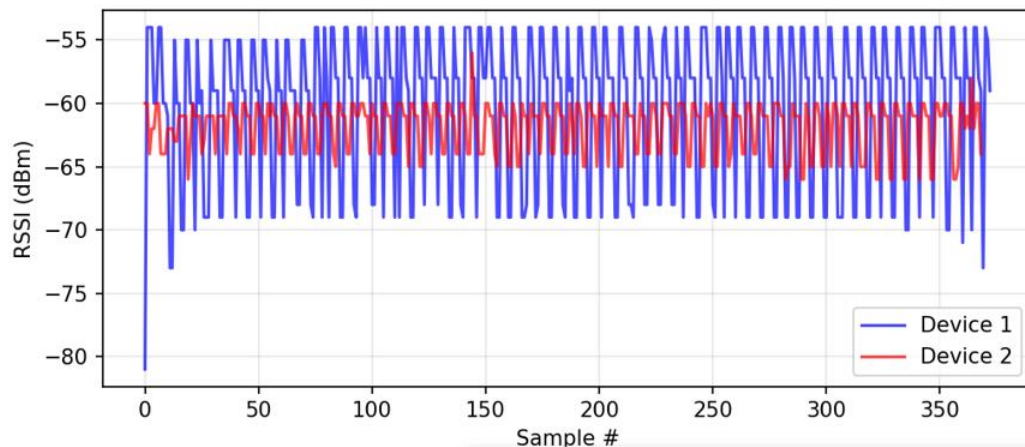


Fig. 1. Measured RSSI values.

As can be seen from **Table 1**, the measured RSSI values exhibit significant variability, which confirms the necessity of applying filtering techniques to reduce the impact of noise and multipath signal propagation.

The following accuracy metrics were computed: Mean Error (ME), Standard Deviation (SD), and Root Mean Square Error (RMSE).

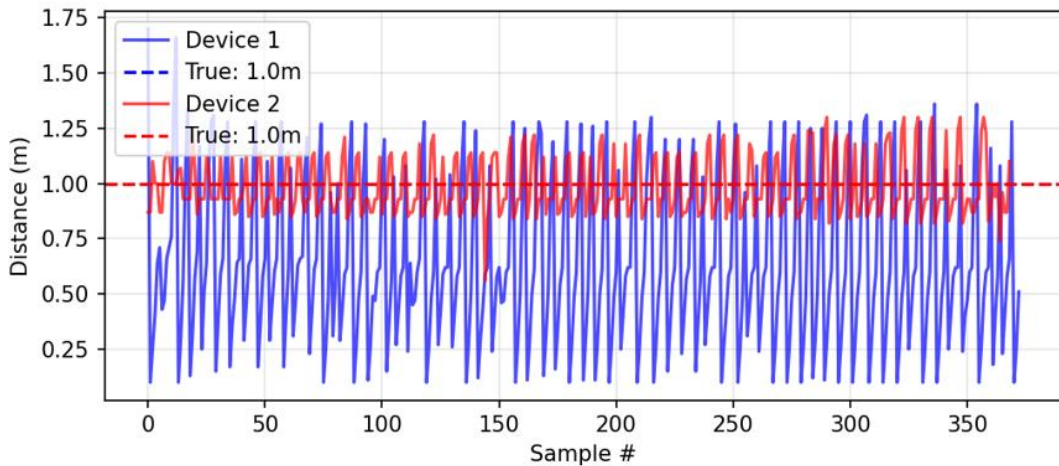


Fig. 2. Estimated distance values.

Table 1. Measurement results

Beacon	Number of measurements	RSSI range, dBm
Device 1	373	-81 ... -54
Device 2	369	-66 ... -56
Total	742	-81 ... -54

The results of applying the Extended Kalman Filter for distance estimation to the beacons are presented in Table 2. The distance estimation error obtained using the EKF is illustrated in Fig. 3.

Table 2. Distance estimation accuracy metrics using EKF

Beacon	Mean estimate, m	Estimated range, m	ME, m	SD, m	RMSE, m
Device 1	0.64	0.10 – 1.70	-0.36	0.38	0.53
Device 2	0.98	0.56 – 1.30	-0.02	0.13	0.14

Device 1: The mean estimated distance is 0.64 m at the actual distance of 1.00 m, indicating a systematic underestimation of the distance by 0.36 m. The standard deviation of 0.38 m reflects a significant spread of estimates around the mean value. The RMSE of 0.53 m characterizes the overall accuracy of the method for this beacon. The wide range of estimated distances (0.10–1.70 m) can be attributed to the high variability of RSSI measurements (from -81 to -54 dBm).

Device 2: The mean estimated distance of 0.98 m with the standard deviation of 0.13 m indicates stable estimates with low variability. An RMSE of 0.14 m demonstrates high estimation accuracy. The narrower range of distance estimates (0.50–1.30 m) and

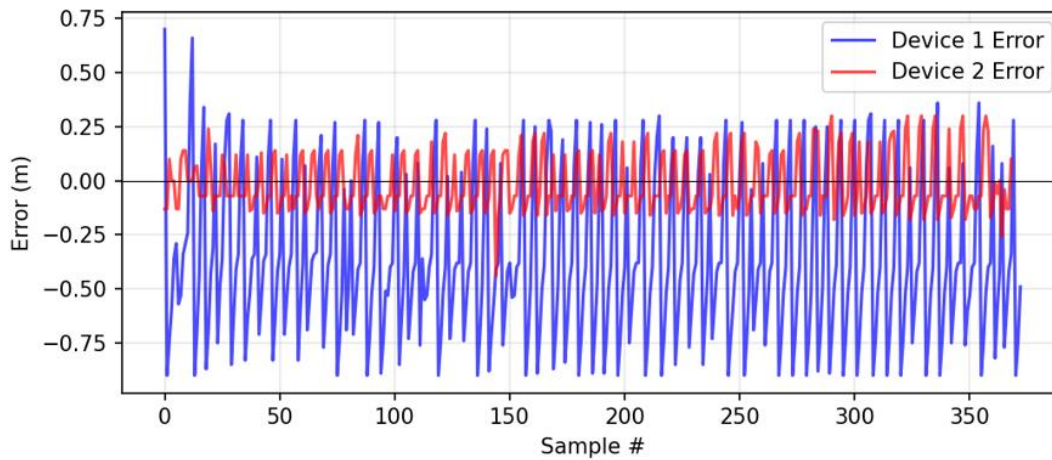


Fig. 3. Distance estimation error.

RSSI values (-66 to -56 dBm) suggests more stable signal propagation conditions for this beacon.

The substantial difference in estimation accuracy between the two beacons (with RMSE differing by 0.39 m) can be explained by the following factors:

1. *Different RSSI variability*: The RSSI range for Device 1 is 27 dBm, whereas for Device 2 it is only 10 dBm. Lower signal variability leads to more stable distance estimates.
2. *Calibration accuracy*: The reference RSSI parameter A for Device 2 (-62 dBm) better matches the actual propagation conditions, which is confirmed by the near-zero mean error.

Despite significant fluctuations in the input RSSI measurements, the estimated distances exhibit a smoothed trajectory without abrupt jumps, which positively demonstrates the effectiveness of the EKF-based approach. The proposed algorithm efficiently filters out anomalous RSSI measurements that may occur due to short-term interference.

Overall, the Extended Kalman Filter demonstrates varying effectiveness for the two beacons. For Device 2, the method provides high accuracy (RMSE = 0.14 m), which is sufficient for indoor navigation and positioning applications. For Device 1, the accuracy is lower (RMSE = 0.53 m), indicating the need for further parameter tuning or the application of more advanced calibration methods.

CONCLUSION

This paper investigated the application of the Extended Kalman Filter to improve the accuracy of distance estimation based on Bluetooth Low Energy RSSI measurements. A mathematical EKF model adapted to the nonlinear logarithmic relationship between signal power and distance was developed, including a detailed derivation of the measurement function Jacobian and a physical interpretation of the signal propagation model parameters.

The proposed algorithm was experimentally validated in a real office environment using two Silicon Labs EFR32BG22 BLE beacons and a Nordic nRF52840 receiver. A total of 742 measurements were collected and analysed at a reference distance of 1.00 m. The results demonstrated varying effectiveness for the two beacons: for Device 2, an RMSE of

0.14 m with a mean error of -0.02 m was achieved, whereas for Device 1, the RMSE was 0.53 m with a mean error of -0.36 m.

The obtained results confirm that the EKF effectively handles the nonlinear relationship between RSSI and distance when the signal propagation model parameters are properly calibrated. For the beacon with optimally tuned parameters, an accuracy of approximately 0.14 m (RMSE) was achieved, which is sufficient for most indoor positioning applications, including smart building navigation, asset tracking, and robotic localization.

A comparison of the mean distance estimation errors obtained using EKF in this work and those reported for the classical Kalman Filter in [29] indicates that, on average, the EKF-based algorithm reduces the mean absolute distance estimation error by 20–30% compared to the classical Kalman Filter.

The practical significance of this work lies in the feasibility of using low-cost BLE beacons to construct indoor positioning systems without the need for additional infrastructure. The proposed algorithm can be integrated into embedded systems based on microcontrollers with limited computational resources.

Future research directions include: (1) the development of automatic calibration methods for the parameters A and n to adapt to changing environmental conditions; (2) evaluation of the algorithm's performance in dynamic scenarios with moving beacons or receivers; (3) integration of the EKF with additional sensors (accelerometer, gyroscope, magnetometer) to improve robustness and accuracy through multisensor data fusion; and (4) comparative analysis with alternative nonlinear filtering techniques such as the Unscented Kalman Filter (UKF) and Particle Filter.

ACKNOWLEDGMENTS AND FUNDING SOURCES

The author received no financial support for the research, authorship, and/or publication of this article.

COMPLIANCE WITH ETHICAL STANDARDS

The author declares that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

REFERENCES

- [1] Zafari, F., Gkelias, A., & Leung, K. K. (2019). A survey of indoor localization systems and technologies. *IEEE communications surveys & tutorials*, 21(3), 2568-2599. <https://doi.org/10.1109/COMST.2019.2911558>.
- [2] Faragher, R., & Harle, R. (2015). Location fingerprinting with bluetooth low energy beacons. *IEEE journal on Selected Areas in Communications*, 33(11), 2418-2428. <https://doi.org/10.1109/JSAC.2015.2430281>.
- [3] Yassin, A., Nasser, Y., Awad, M., Al-Dubai, A., Liu, R., Yuen, C., ... & Aboutanios, E. (2016). Recent advances in indoor localization: A survey on theoretical approaches and applications. *IEEE Communications Surveys & Tutorials*, 19(2), 1327-1346. <https://doi.org/10.1109/COMST.2016.2632427>.
- [4] Cho, S. Y., & Park, C. G. (2006). MEMS-based pedestrian navigation system. *Journal of Navigation*, 59(1), 135–153. <https://doi.org/10.1017/S0373463305003486>
- [5] Booranawong, A., Sengchuai, K., Buranapanichkit, D., & Jindapetch, N. (2022). RSSI-based indoor localization using multilateration with zone selection and virtual position-based compensation methods. *IEEE Access*, 10, 66925–66939. <https://doi.org/10.1109/ACCESS.2021.3068295>

- [6] Liu, H., Darabi, H., Banerjee, P., & Liu, J. (2007). Survey of wireless indoor positioning techniques and systems. *IEEE Transactions on Systems, Man, and Cybernetics, Part C*, 37(6), 1067–1080. <https://doi.org/10.1109/TSMCC.2007.905750>
- [7] Chen, L., Pei, L., Kuusniemi, H., et al. (2013). Bayesian fusion for indoor positioning using Bluetooth fingerprints. *Wireless Personal Communications*, 70(4), 1735–1745. <https://doi.org/10.1007/s11277-012-0777-1>
- [8] Welch, G., & Bishop, G. (2006). *An introduction to the Kalman filter* (Technical Report TR 95-041). University of North Carolina at Chapel Hill.
- [9] Sthapit, P., Pyun, J. Y., & Gang, H. S. (2018). Bluetooth-based indoor positioning using machine learning algorithms. In *IEEE International Conference on Consumer Electronics* (pp. 1–2). <https://doi.org/10.1109/ICCE-ASIA.2018.8552138>
- [10] Rappaport, T. S. (2002). *Wireless communications: Principles and practice* (2nd ed.). Prentice Hall.
- [11] Kaemarungsi, K., & Krishnamurthy, P. (2004). Modeling of indoor positioning systems based on location fingerprinting. In *IEEE INFOCOM*, 2, 1012–1022. <https://doi.org/10.1109/INFCOM.2004.1356988>
- [12] Li, G., Geng, E., Ye, Z., et al. (2018). Indoor positioning algorithm based on the improved RSSI distance model. *Sensors*, 18(9), 2820. <https://doi.org/10.3390/s18092820>
- [13] Kalman, R. E. (1960). A new approach to linear filtering and prediction problems. *Journal of Basic Engineering*, 82(1), 35–45.
- [14] Maybeck, P. S. (1979). *Stochastic models, estimation, and control*. Academic Press.
- [15] Bahl, P., & Padmanabhan, V. N. (2000). RADAR: An in-building RF-based user location and tracking system. *Proceedings of IEEE INFOCOM*, 2, 775–784. <https://doi.org/10.1109/INFCOM.2000.832252>
- [16] Julier, S. J., & Uhlmann, J. K. (2004). Unscented filtering and nonlinear estimation. *Proceedings of the IEEE*, 92(3), 401–422. <https://doi.org/10.1109/JPROC.2003.823141>
- [17] Bar-Shalom, Y., Li, X. R., & Kirubarajan, T. (2001). *Estimation with applications to tracking and navigation*. Wiley. <https://doi.org/10.1002/0471221279>
- [18] Simon, D. (2006). *Optimal state estimation: Kalman, H infinity, and nonlinear approaches*. Wiley. <https://doi.org/10.1002/0470045345>
- [19] Akl, R., Tummala, D., & Li, X. (2011). Indoor propagation modeling at 2.4 GHz for IEEE 802.11 networks. In *Proceedings of the IASTED International Conference on Wireless and Optical Communications* (pp. 1–5).
- [20] Bluetooth Special Interest Group. (2019). *Bluetooth Core Specification Version 5.1*. <https://www.bluetooth.com/specifications/bluetooth-core-specification/>
- [21] Goldsmith, A. (2005). *Wireless communications*. Cambridge University Press. <https://doi.org/10.1017/CBO9780511841224>
- [22] Gustafsson, F., & Hendeby, G. (2012). Some relations between extended and unscented Kalman filters. *IEEE Transactions on Signal Processing*, 60(2), 545–555. <https://doi.org/10.1109/TSP.2011.2172431>
- [23] Seidel, S. Y., & Rappaport, T. S. (1992). 914 MHz path loss prediction models for indoor wireless communications in multifloored buildings. *IEEE Transactions on Antennas and Propagation*, 40(2), 207–217. <https://doi.org/10.1109/8.127405>
- [24] Thrun, S., Burgard, W., & Fox, D. (2005). *Probabilistic robotics*. MIT Press.
- [25] Zhuang, Y., Yang, J., Li, Y., et al. (2016). Smartphone-based indoor localization with Bluetooth low energy beacons. *Sensors*, 16(5), 596. <https://doi.org/10.3390/s16050596>
- [26] Brown, R. G., & Hwang, P. Y. (1997). *Introduction to random signals and applied Kalman filtering* (3rd ed.). Wiley.

- [27] Ristic, B., Arulampalam, S., & Gordon, N. (2004). *Beyond the Kalman filter: Particle filters for tracking applications*. Artech House. <https://doi.org/10.1155/S1110865704405095>
- [28] Grewal, M. S., & Andrews, A. P. (2014). *Kalman filtering: Theory and practice using MATLAB* (4th ed.). Wiley. <https://doi.org/10.1002/9781118984987>
- [29] Kalynchuk, T.-N., & Shevchyk, V. (2025). Indoor positioning with Bluetooth low energy: A preliminary system design and results. *Advances in Cyber-Physical Systems*, 28–33. <https://doi.org/10.23939/acps2025.01.028>
-

ВИЗНАЧЕННЯ ПОЛОЖЕННЯ ОБ'ЄКТІВ В ПРИМІЩЕННІ З ВИКОРИСТАННЯМ РАДІОТЕХНОЛОГІЇ BLUETOOTH LOW ENERGY ТА РОЗШИРЕНОГО ФІЛЬТРА КАЛМАНА

Тадей-Назарій Калинчук 

*Кафедра радіофізики та комп'ютерних технологій,
Львівський національний університет імені Івана Франка
вул. Ген. Тарнавського, 107, 79017 Львів, Україна*

АНОТАЦІЯ

Вступ. Системи внутрішнього позиціонування на основі маяків *Блютуз з низьким енергоспоживанням* (БНЕ, Bluetooth Low Energy, BLE) часто використовують оцінку відстані за *показником потужності прийнятого сигналу* (ПППС, Received Signal Strength Indicator, RSSI). Проте вимірювання ПППС в приміщеннях істотно залежать від багатопробного поширення, затінення, інтерференції та поглинання сигналу перешкодами, що призводить до значних викидів у значеннях сигналу та суттєвих похибок оцінки відстані. Нелінійна логарифмічна залежність між ПППС та відстанню додатково ускладнює застосування класичних лінійних методів фільтрації, зокрема фільтра Калмана, який потребує попереднього нелінійного перетворення вимірювань і може втрачати оптимальність.

Матеріали та методи. У роботі запропоновано метод оцінки відстані на основі розширеного фільтра Калмана (РФК, Extended Kalman Filter, EKF), який безпосередньо обробляє вимірювання ПППС з використанням нелінійної логарифмічної моделі втрат потужності сигналу. Експеримент проведено в офісному приміщенні з використанням двох БНЕ-маяків Silicon Labs EFR32BG22 та приймача Nordic nRF52840. Параметри РФК визначено на основі попереднього калібрування коефіцієнтів моделі поширення сигналу.

Результати. Отримані результати показали, що РФК ефективно згладжує флуктуації ПППС. Для маяка з меншою дисперсією ПППС середньоквадратична похибка становила 0,14 м, а для іншого – 0,53 м. Аналіз підтверджує, що точність оцінювання суттєво залежить від стабільності сигналу та якості калібрування параметрів моделі. Порівняно з методом прямого перетворення ПППС у відстань та класичним фільтром Калмана, розглянутими в попередніх дослідженнях, застосування РФК дозволяє зменшити середню абсолютну похибку оцінки відстані приблизно на 20–30%, що підтверджує переваги нелінійної фільтрації.

Висновки. Запропонований метод на основі розширеного фільтра Калмана підвищує точність та надійність оцінки відстані в системах внутрішнього позиціонування з використанням БНЕ. За умов належного калібрування параметрів моделі досягається точність, достатня для практичного застосування у навігації в будівлях, відстеженні активів та локалізації роботів. Алгоритм може бути реалізований на

вбудованих платформах з обмеженими обчислювальними ресурсами та є основою для подальшого розвитку мультисенсорних систем внутрішнього позиціонування.

Ключові слова: позиціонування, Блютуз, БНЕ, ПППС, фільтр Калмана, розширений фільтр Калмана.

Received / Одержано
09 February, 2026

Revised / Доопрацьовано
28 February, 2026

Accepted / Прийнято
02 March, 2026

Published / Опубліковано
30 March, 2026

UDC: 007.5

COMMUNICATION ARCHITECTURES FOR CLOUD-INTEGRATED SPECTROMETRIC LABORATORIES: ESP-NOW VS MQTT

Andriy Krupych*, Pavlo Levush

Department of Radiophysics and Computer Technologies,
Ivan Franko National University of Lviv
107 Gen. Tarnavsky Str., UA-79017 Lviv, Ukraine

Krupych, A., Levush, P. (2026). Communication Architectures for Cloud-Integrated Spectrometric Laboratories: ESP-NOW vs MQTT. *Electronics and Information Technologies*, 33, 43–56
<https://doi.org/10.30970/eli.33.4>

ABSTRACT

Background. Cloud-integrated spectrometric laboratories face communication challenges in achieving real-time data access and analysis. This study compares two wireless protocols, MQTT (Message Queuing Telemetry Transport) and ESP-NOW (Espressif NOW), for LED control in such environments. MQTT offers lightweight, bandwidth-efficient, publish-subscribe messaging [1–3], while ESP-NOW provides energy-efficient direct communication without a Wi-Fi router. The objective is to evaluate their performance and suitability.

Materials and Methods. An experimental setup involved a StellarNet spectrometer, LED light sources, and ESP32 microcontrollers. Two architectures were tested: 1) direct MQTT Communication, where each ESP32 connected directly via Wi-Fi to an MQTT broker; and 2) an edge device with ESP-NOW relay, using an edge ESP32 for MQTT/Wi-Fi communication, then relaying commands via ESP-NOW to other ESP32s. Response times for LED control were measured over 100 cycles, and data were analyzed using descriptive statistics and an independent samples t-test.

Results and Discussion. Direct MQTT Communication exhibited significantly lower latency (median ~60 ms) and tighter distribution compared to the Edge Device with ESP-NOW Relay (median ~170 ms). A t-test confirmed a statistically significant difference ($t=46.28$), with MQTT demonstrating faster response times. However, the ESP-NOW relay system offers architectural advantages: reduced Wi-Fi dependency for individual nodes, enhanced deployment flexibility in areas with poor Wi-Fi coverage, improved scalability [4], and energy efficiency, making its higher latency a practical trade-off for large-scale laboratory integration.

Conclusion. Direct MQTT Communication provides superior low-latency performance. However, the edge device with ESP-NOW relay, despite higher latency, is a highly acceptable solution due to its flexibility, scalability, and reduced Wi-Fi dependency for distributed spectrometric laboratories. This highlights a critical trade-off between absolute speed and architectural benefits for robust, smart, cloud-enabled analytical laboratories.

Keywords: MQTT, ESP-NOW, spectrometric laboratory, Internet of Things, latency, embedded systems hardware

INTRODUCTION

Spectrometric laboratories are crucial in various scientific and industrial domains, generating vast amounts of data essential for research, quality control, and process



© 2026 Andriy Krupych & Pavlo Levush. Published by the Ivan Franko National University of Lviv on behalf of Електроніка та інформаційні технології / Electronics and Information Technologies. This is an Open Access article distributed under the terms of the [Creative Commons Attribution 4.0 License](https://creativecommons.org/licenses/by/4.0/) which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

optimization. The increasing complexity of analytical processes and the demand for real-time data access and analysis have driven the adoption of cloud integration within these laboratory environments [5]. Cloud-based systems offer significant advantages by providing a centralized infrastructure for data storage, archiving, and analysis, thereby supporting knowledge acquisition from analytical data [5, 6]. However, the integration of diverse analytical instruments and systems into a cohesive cloud-integrated framework presents considerable communication challenges [7].

The analytical laboratory often comprises numerous computer-assisted tools and software applications, yet communication constraints frequently impede the seamless sharing of scientific data [7]. Traditional wired data acquisition systems, while reliable, can limit accessibility and mobility, leading to higher installation and troubleshooting costs [8]. The shift towards wireless communication offers benefits such as reduced system size, lower costs, enhanced flexibility, and simplified deployments [8]. However, the diverse requirements of laboratory automation, including real-time data acquisition, scalability, energy efficiency, interoperability, and security, pose significant hurdles for selecting appropriate communication architectures [9].

In the context of the Internet of Things and smart laboratories, various wireless communication protocols have emerged to address these challenges. Among these, Message Queuing Telemetry Transport (MQTT) and ESP-NOW are prominent for their distinct characteristics. MQTT is a lightweight messaging protocol well-suited for IoT applications due to its low power consumption and bandwidth efficiency, often implemented on microcontrollers like the ESP32 [2]. It facilitates communication through a publish-subscribe model, enabling devices to connect and exchange data efficiently [10]. Conversely, ESP-NOW offers an energy-efficient, connectionless communication protocol that allows devices to directly exchange data without requiring a Wi-Fi router, making it ideal for scenarios requiring fast, local data exchange with minimal delay [11, 12]. It utilizes MAC addresses for device addressing and has a maximum message size of 250 bytes, offering an alternative when parallel Wi-Fi communication is not feasible [13].

This paper presents a comparative analysis of ESP-NOW and MQTT as communication architectures for cloud-integrated spectrometric laboratories. By examining their respective strengths and weaknesses in terms of performance, reliability, and suitability for the unique demands of laboratory environments, this study aims to provide insights into optimizing data flow and connectivity. The findings will contribute to the development of more efficient and robust communication strategies for the next generation of smart, cloud-enabled analytical laboratories.

MATERIALS AND METHODS

This section details the experimental design and methodologies employed to compare two distinct communication architectures for remotely controlling LED light sources (CHANZON 3V 5mm 30°: LED#1 – White, 12-14 cd; LED#2 – Warm White, 14-16 cd) and a miniature StellarNet GREEN-Wave VIS-50 fiber-optically coupled spectrometer within a cloud-enabled laboratory environment. The primary objective was to assess the performance, particularly in terms of command response times, and identify the pros and cons of each architecture.

Experimental Setup

The experimental setup was configured to simulate a practical cloud-integrated spectrometric laboratory, where LED light sources are controlled remotely, and their spectral data can be acquired.

The core hardware elements utilized in this study included:

- Spectrometer: A StellarNet spectrometer was used to acquire spectral data from the various LED light sources. This spectrometer served as the data collection point for the light emitted by the remotely controlled LEDs.

- LED Light Sources: Multiple LED light sources were implemented, each connected to an individual ESP32 microcontroller for remote control.
- Microcontroller Units: ESP32 development boards served as the primary control units for each LED light source. These boards are widely used in IoT applications [14–16]. An additional ESP32 board functioned as an edge device for one of the proposed communication architectures.

Two distinct communication architectures were designed and implemented for comparison.

Architecture 1: Direct MQTT Communication

In this architecture, each ESP32 device controlling a LED light source was directly connected to the local Wi-Fi network (**Fig. 1**). These individual ESP32 devices were configured as MQTT clients and subscribed to a specific topic on the HiveMQ MQTT broker. The HiveMQ MQTT broker served as the central message broker, facilitating the remote transmission of commands from a central control system to the individual LED-controlling ESP32s. This architecture required each ESP32 to maintain an active Wi-Fi connection and an MQTT client subscription.

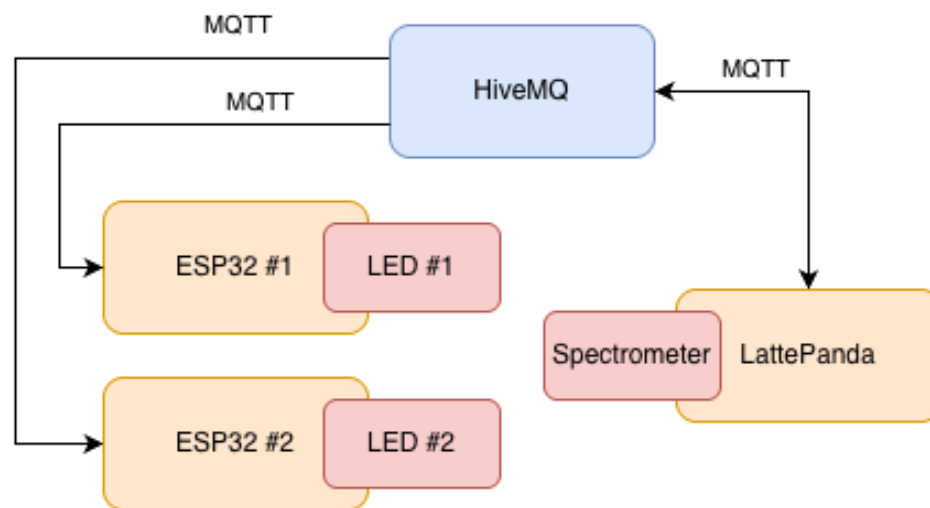


Fig. 1. Architecture 1 with direct MQTT connections

The LattePanda Single-Board Computer serves as the central control unit in this system, integrating both hardware-level control and high-level data processing capabilities. Its hybrid architecture, which combines a Windows-based computing environment with an embedded Arduino Leonardo microcontroller, makes it ideal for IoT applications requiring both GPIO pin access for direct hardware manipulation and Python-based spectral analysis.

Architecture 2: Edge Device with ESP-NOW Relay

This architecture introduced an intermediary "edge" ESP32 device (**Fig. 2**). This edge device was connected to the local Wi-Fi network and functioned as an MQTT client, subscribing to HiveMQ MQTT broker to receive commands. Upon receiving a command from the broker, the edge ESP32 identified the target LED and forwarded the command to the corresponding LED-controlling ESP32 device using the ESP-NOW protocol. Crucially, the individual ESP32 devices connected to the LED light sources in this architecture were not required to have a direct Wi-Fi connection, relying instead on the ESP-NOW communication link with the edge device.

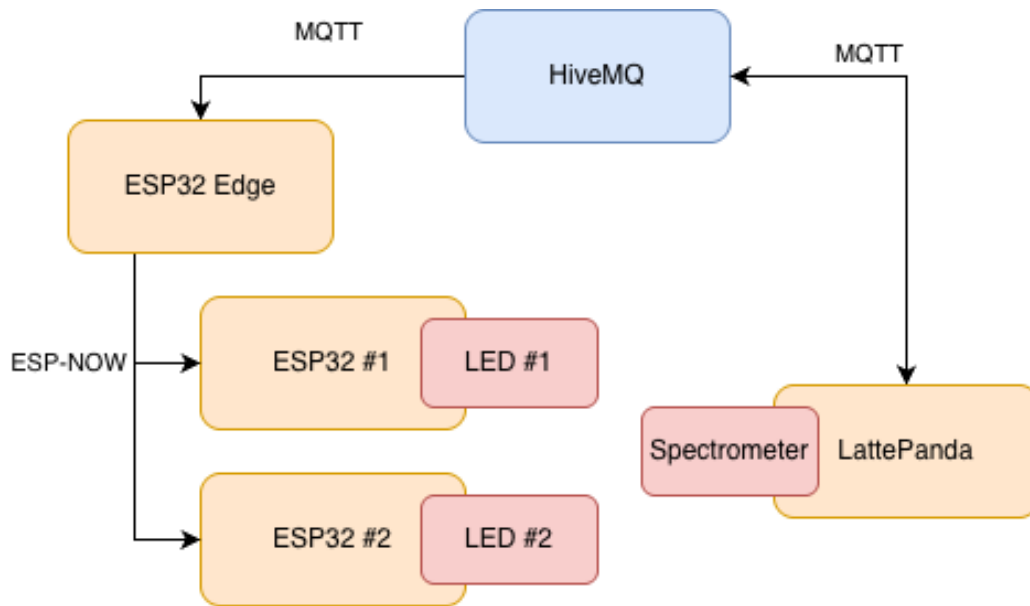


Fig. 2. Architecture 2 with ESP-NOW connection

Architecture 2 distinguishes itself from Architecture 1 by introducing an intermediary ESP32 edge device. This edge device is the sole component requiring a Wi-Fi connection and MQTT subscription, relaying commands via ESP-NOW to the individual LED-controlling ESP32s, thereby freeing them from direct Wi-Fi dependency.

Communication Protocol Implementation

The MQTT protocol was implemented using standard client libraries on the ESP32 platforms [17]. The HiveMQ MQTT broker served as the central message broker, facilitating communication between the control system and the ESP32 devices. In Architecture 1, each LED-controlling ESP32 is directly connected to this broker. In Architecture 2, only the edge ESP32 maintained a connection to the HiveMQ broker. Commands were structured as JSON payloads containing the target LED identifier and the desired action (turning on/off, adjusting intensity). Quality of Service levels could be varied during experiments to assess their impact on reliability and latency [18].

In Architecture 2, the edge ESP32 was configured as the sender for ESP-NOW messages, and the LED-controlling ESP32s were configured as receivers. Each LED-controlling ESP32 was registered with the edge device using its unique MAC address. Upon receiving an MQTT command, the edge ESP32 encapsulated the command into an ESP-NOW packet and transmitted it directly to the MAC address of the target LED-controlling ESP32. A callback function was implemented on both the sender and receiver sides to confirm message delivery. The connectionless nature of ESP-NOW was leveraged for fast, local data exchange without the overhead of establishing and maintaining Wi-Fi connections [19].

Command Transmission and Response Time Measurement

Commands for controlling the LED light sources were initiated from a central control application. The primary data collected for analysis was the response time, defined as the duration from when a command was issued by the central control application until the corresponding LED light source executed the command (e.g., changed its state or intensity). For both architectures, timestamps were recorded at the point of command issuance and at the point of command execution by the LED-controlling ESP32.

Performance Metrics and Measurement

The performance evaluation primarily focused on the response time, defined as the time taken for a data packet (command) to travel from the command source to the target LED-controlling ESP32 and be executed. This was measured by embedding timestamps at the origin and recording them upon reception and execution. The ability to measure and compare latency is crucial for evaluating communication systems [20].

Measurement tools included internal timers on the ESP32 devices for precise timing of command reception and execution, as well as timestamps recorded by the central control system for command issuance.

The use of MQTT QoS 2 (Exactly Once delivery) is particularly significant for the reliability discussion in this study. QoS 2 is the highest assurance level in the MQTT protocol, employing a four-step handshake (PUBLISH, PUBREC, PUBREL, PUBCOMP) to guarantee that each message is delivered exactly once to the subscriber, eliminating both message loss and duplication. In a spectrometric laboratory context, this is critical because duplicate LED commands could trigger redundant state changes during a spectral acquisition sequence, while lost commands could leave a light source in an incorrect state, compromising measurement integrity. The fact that QoS 2 was used in this study means the latency figures reported for both architectures already include the overhead of this four-step handshake, which inherently adds processing time compared to QoS 0 or QoS 1. This has two important implications: first, the MQTT leg of both architectures benefits from protocol-level delivery guarantees, providing a strong reliability baseline; second, the latency values presented represent a conservative (i.e., worst-case) scenario for MQTT performance, and lower latencies could be achieved if a less strict QoS level were acceptable for a given application.

In Architecture 2, to compensate for the absence of built-in reliability guarantees comparable to MQTT QoS levels, a programmatic retry mechanism was implemented on the ESP-NOW communication leg based on the protocol's delivery callback functionality. ESP-NOW provides a send callback (`on_data_sent`) that reports whether a transmitted packet was successfully received by the target peer at the MAC layer. When this callback indicated a delivery failure, the edge ESP32 device automatically re-attempted transmission of the command to the target LED-controlling ESP32, up to a predefined number of retries before reporting an error. This approach effectively introduces an application-layer reliability mechanism that mirrors, to a degree, the acknowledgment logic inherent in MQTT QoS levels, ensuring that transient interference or momentary unavailability of a receiver node does not result in a permanently lost command.

Experimental Procedure

Experiments were conducted to compare the response times of Architecture 1 and Architecture 2. Each experimental run involved:

1. *System Initialization*: All ESP32 devices and the StellarNet spectrometer were powered on and initialized. The HiveMQ MQTT broker was confirmed to be operational.
2. *Command Execution*: A series of predefined commands - sequential ON/OFF toggling of LEDs - were issued from the central control application.
3. *Data Logging*: Response times for each command were recorded for both architectures. This involved logging the timestamp of command issuance and the timestamp of command execution on the target ESP32.
4. *Repetition*: Each experimental scenario was repeated 100 times to ensure statistical significance and minimize the impact of transient network fluctuations. Experiments were conducted under consistent environmental conditions to ensure comparability.

Experiments were conducted under consistent environmental conditions to ensure comparability. Wi-Fi channel congestion was monitored; no other devices shared the

network, and tests were run at specific times to avoid interference. A dedicated access point with no competing traffic was used, along with fixed physical distances between devices.

A follow-up study measuring delivery success rates under varying network loads and interference conditions would substantially complement the latency analysis and provide a more complete basis for recommending one architecture over the other.

Data Analysis

Raw response time data, collected from the experimental runs, were processed using custom scripts developed in the Python programming language. For statistical analysis, Python was also utilized, incorporating relevant statistical libraries such as *Pandas* for data manipulation, *Scipy.stats* for inferential statistics, and *Matplotlib.pyplot* and *Seaborn* for data visualization.

Descriptive statistics, including mean, median, and measures of spread, were calculated for the response times of both the Direct MQTT Communication and Edge Device with ESP-NOW Relay architectures.

To determine if the observed differences in mean response times between the two architectures were statistically significant, an independent samples t-test was performed. This inferential statistical test is specifically designed to compare the means of two independent groups. A significance level of $p < 0.05$ was considered statistically significant for all comparisons. This conventional threshold was initially proposed by R. Fisher and remains a widely accepted standard in many scientific fields [21–23].

Data visualization techniques, including histograms and scatter plots, were employed to illustrate the distribution characteristics, variability, and overall comparison of response times for both architectures. These visualizations complemented the statistical analysis by providing a clear graphical representation of the performance differences.

RESULTS AND DISCUSSION

This section presents empirical results obtained from the comparative analysis of Direct MQTT Communication and Edge Device with ESP-NOW Relay. The measured response times for LED control commands are analyzed, followed by a statistical comparison and a discussion of the practical implications and acceptability of each architecture within a cloud-integrated spectrometric laboratory.

Latency Analysis

The performance of both communication architectures was evaluated based on command response times over 100 experimental cycles. In practice, 100 repeated measurements per ESP device (yielding $n=200$ per architecture per device and $n=400$ overall) provided sufficient statistical power to detect the large effect size observed — the t-statistic of -46.28 reflects an enormous difference relative to within-group variance, meaning even a smaller sample would have reached significance. The central limit theorem also ensures that with $n=400$, sampling distributions of the means are well-approximated as normal, supporting the validity of the t-test.

The Direct MQTT Communication architecture consistently demonstrated lower latency (**Fig. 3**). As illustrated in the Histogram of MQTT Latency (**Fig. 4**), the response times were heavily concentrated at the lower end of the spectrum, indicating quick command execution.

Median response time was approximately 60 ms. The interquartile range for MQTT latency was tight, spanning roughly from 55 ms to 75 ms, suggesting consistent performance for the majority of command cycles. However, a few outliers were observed, with one significant spike approaching 440 ms. These transient peaks could be attributed to network congestion, Wi-Fi interference, or processing delays within the MQTT broker or the individual ESP32 Wi-Fi stack.

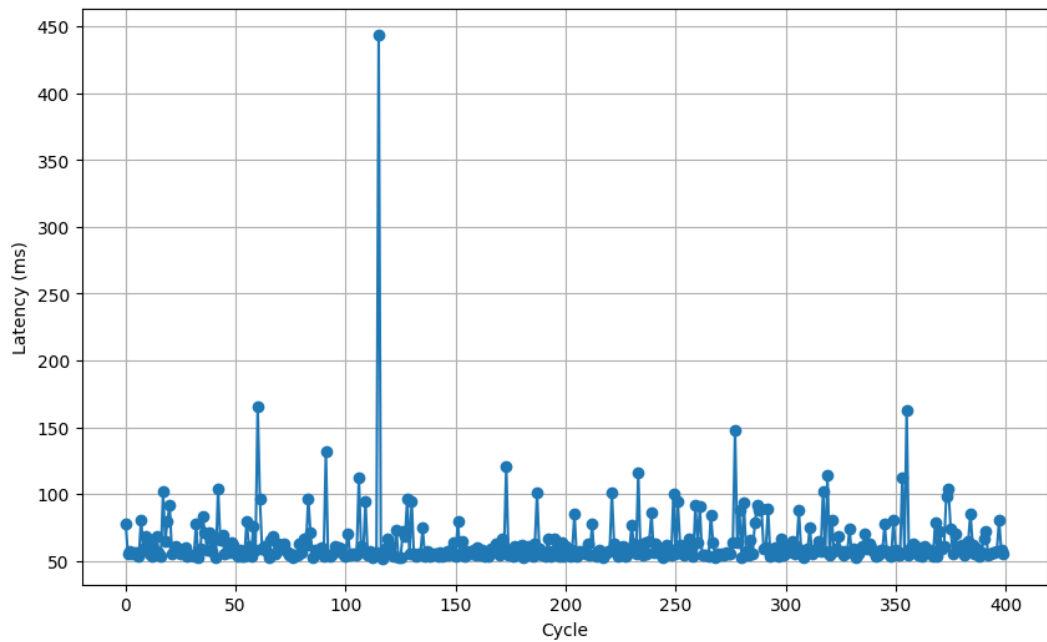


Fig. 3. MQTT architecture latency over 100 cycles plot.

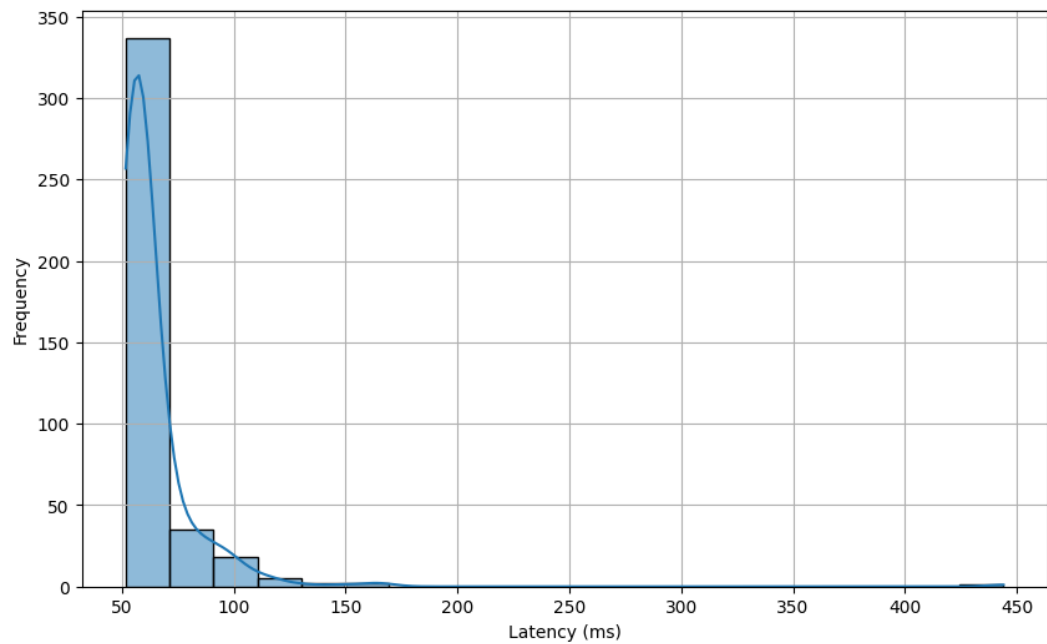


Fig. 4. Histogram of MQTT architecture.

In contrast, the edge device with ESP-NOW relay architecture exhibited generally higher response times (Fig. 5). The histogram of NOW Latency (Fig. 6) reveals a broader distribution of latencies, with the peak shifted towards approximately 170 ms.

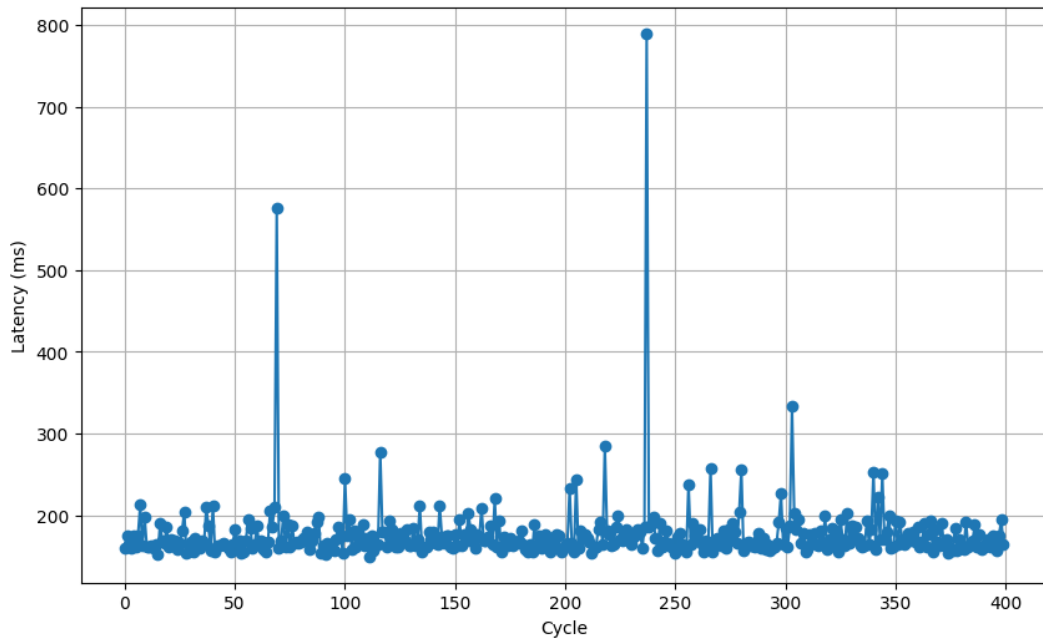


Fig. 5. ESP-NOW architecture latency over 100 cycles plot.

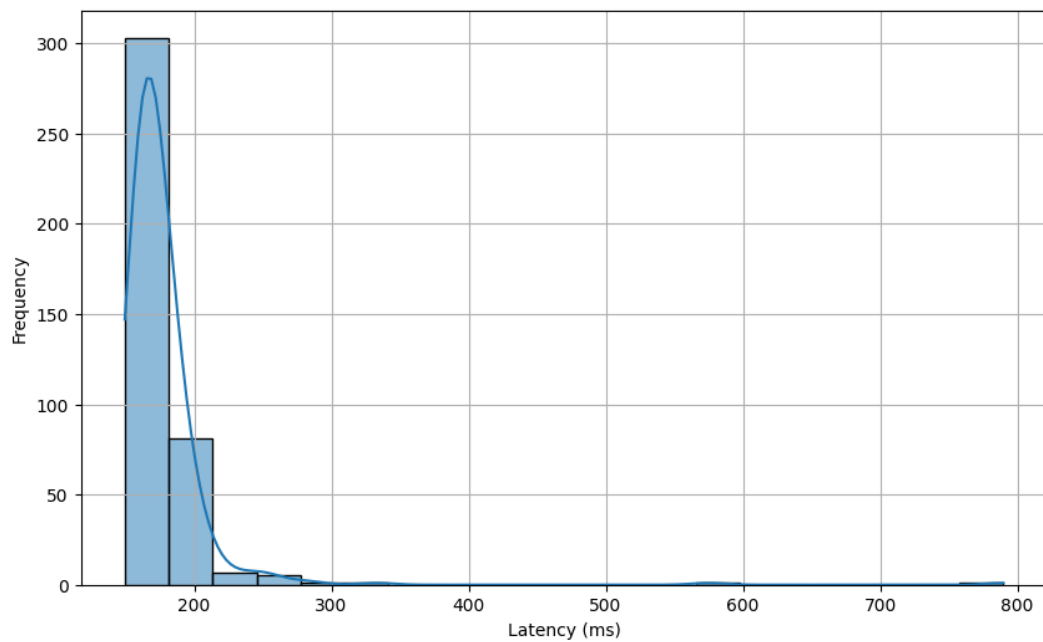


Fig. 6. Histogram of ESP-NOW architecture.

Median response time was around 170 ms. The interquartile range for ESP-NOW was wider than that of MQTT, indicating greater variability in response times. Furthermore, the NOW Latency Over 100 Cycles plot highlighted more frequent and higher magnitude outliers compared to MQTT, with some response times reaching approximately 580 ms and a maximum of about 790 ms. These increased latencies and variability are likely a

consequence of the added processing time at the edge device for translating MQTT commands to ESP-NOW messages, the additional communication hop, and potential overheads in the ESP-NOW protocol itself.

Statistical Comparison

A qualitative statistical comparison of the two architectures reveals distinct characteristics (**Table 1**). MQTT consistently provided lower median latency and a tighter distribution of response times, suggesting superior real-time performance and greater consistency for command execution. The spread of data, as represented by the IQR, was considerably smaller for MQTT, indicating less variability. While both protocols experienced latency spikes, those observed in the ESP-NOW architecture were generally of higher magnitude and occurred with greater frequency.

The results of the independent samples t-test, comparing the two datasets, yielded a t-statistic of approximately -46.28 . This t-statistic, with its large absolute value, indicates a substantial difference between the mean latencies of the two communication architectures. The negative value suggests that the mean latency for Architecture 1 was significantly lower than that for Architecture 2.

This statistically significant result allows us to confidently reject the null hypothesis, which posits that there is no difference in the mean response times between the two architectures.

Therefore, the statistical analysis confirms that there is a highly significant difference in the command response times, with architecture 1 demonstrating a statistically significantly lower latency compared to architecture 2. This quantitative finding supports the qualitative observations made during the initial latency analysis, solidifying the conclusion that MQTT provides faster and more consistent command execution in this experimental setup.

Table 1. Summary of latency statistics

Statistics		Latency, ms	
		Architecture 1 (MQTT → ESP)	Architecture 2 (MQTT → EDGE → ESP-NOW → ESP)
ESP1 (n = 200)	Min	51.76	153.22
	Max	165.91	789.92
	Average	62.97	179.52
ESP2 (n = 200)	Min	52.96	149.39
	Max	444.03	332.91
	Average	64.91	173.44
OVERALL (n = 400)	Min	51.76	149.39
	Max	444.03	789.92
	Average	63.94	176.48

Discussion of Architectural Implications and Acceptability

The choice between architectures involves a trade-off between absolute latency performance and architectural flexibility, scalability, and deployment ease.

While direct MQTT communication clearly demonstrates lower latency, making it ideal for applications requiring almost real-time, highly consistent communication, its deployment necessitates that every single ESP32 device controlling an LED light source maintains an

active Wi-Fi connection and an MQTT client subscription. This approach can lead to several challenges in a large-scale spectrometric laboratory or institutional setting:

- *Wi-Fi Network Load*: Many Wi-Fi-connected devices can strain the access point's capacity and bandwidth.
- *Configuration Overhead*: Managing Wi-Fi credentials and MQTT subscriptions for numerous individual devices can be complex.
- *Coverage Limitations*: Devices in areas with poor Wi-Fi coverage would struggle to maintain connectivity.

Despite its higher latency, architecture 2 offers compelling advantages that make it highly acceptable, and in many scenarios, a more practical solution:

- *Reduced Wi-Fi Dependency*: Only the central edge ESP32 device requires a Wi-Fi connection to the MQTT broker. The individual LED-controlling ESP32s communicate directly with the edge device via ESP-NOW, eliminating their need for a dedicated Wi-Fi connection.
- *Enhanced Deployment Flexibility*: This architecture allows for the "spreading" of ESP-NOW nodes throughout an institution or laboratory space, even in areas with limited or no Wi-Fi connectivity. This is particularly beneficial for experimental setups that are physically dispersed or located in difficult-to-reach areas.
- *Scalability*: ESP-NOW can support numerous devices without each consuming a Wi-Fi connection resource, simplifying network management and potentially scaling better in terms of Wi-Fi access point load. The ability to deal with resource-constrained devices and support frequent connection/disconnection also contributes to scalability [4].
- *Energy Efficiency*: ESP-NOW is designed as a low-power, connection-less communication protocol, making it suitable for sensor nodes where energy conservation is critical, as they do not need to maintain a constant Wi-Fi connection.
- *Acceptable Latency for Application*: For the specific application of remotely controlling LED light sources in a spectrometric laboratory, the observed median latency of approximately 170 ms, even with occasional spikes, is generally well within acceptable limits. Commands such as turning LEDs ON/OFF or adjusting intensity are typically not time-critical to the sub-millisecond level. Human operators or automated scripts can easily tolerate response times within a few hundreds of milliseconds, or even up to a second, without significant impact on the experimental workflow or data integrity. The primary requirement is reliable execution, which ESP-NOW provides through its packet delivery mechanisms [19].

In summary, it should be noted that while architecture 1 offers superior raw latency performance, architecture 2 presents a highly acceptable and often preferred solution for building a cloud-integrated spectrometric laboratory due to its significant architectural advantages. The flexibility to deploy numerous devices without individual Wi-Fi connections, coupled with its scalability and energy efficiency benefits, outweighs the higher, yet application-acceptable, latency for many practical laboratory automation scenarios.

CONCLUSION

This study meticulously compared two distinct communication architectures – direct MQTT communication and an edge device with ESP-NOW relay – for their suitability in cloud-integrated spectrometric laboratories, focusing on command response times for LED light source control. The increasing complexity of analytical processes and the demand for real-time data access and analysis have driven the adoption of cloud integration within these laboratory environments. This research aimed to evaluate their performance, reliability, and practical implications within this specialized environment.

The empirical analysis, reinforced by a rigorous independent samples t-test, conclusively demonstrated a statistically significant difference in latency between the two architectures. Direct MQTT Communication exhibited a significantly lower mean latency, indicating superior speed and consistency in command execution when each ESP32 device maintained a direct Wi-Fi and MQTT connection. MQTT's lightweight nature, low power consumption, and bandwidth efficiency make it well-suited for IoT applications [1–3]. The t-test results, with a t-statistic of approximately -46.28 , unequivocally support this finding.

However, the research also highlighted that absolute latency is not the sole determinant of an architecture's practical utility. While the edge device with ESP-NOW relay presented a higher, though application-acceptable, latency, it offered substantial advantages in architectural flexibility, scalability, and reduced Wi-Fi dependency. ESP-NOW is an energy-efficient, connectionless protocol that allows devices to directly exchange data without requiring a Wi-Fi router and utilizes MAC addresses for device addressing. This architecture is particularly beneficial for large or geographically dispersed laboratory setups where direct Wi-Fi connectivity for every node is impractical or creates excessive network load. By centralizing the Wi-Fi connection through an edge device, the ESP-NOW relay system enables the deployment of numerous sensor nodes in areas with limited or no Wi-Fi infrastructure, simplifying network management and enhancing energy efficiency for individual devices. The shift towards wireless communication, in general, offers benefits such as reduced system size, lower costs, enhanced flexibility, and simplified deployments compared to traditional wired systems [8, 24].

In conclusion, the choice between these two architectures necessitates careful consideration of the specific demands of the spectrometric laboratory. For applications demanding the absolute lowest latency and where pervasive Wi-Fi coverage is guaranteed, direct MQTT communication is the preferred choice. Conversely, for scenarios prioritizing architectural flexibility, broad deployment across an institution, scalability (due to less strain on Wi-Fi resources), and reduced individual device Wi-Fi dependency – even with slightly elevated but still acceptable latencies – the edge device with ESP-NOW relay stands out as a highly effective and practical solution. This comparative analysis provides valuable insights for optimizing communication strategies in the development of more efficient and robust smart, cloud-enabled analytical laboratories.

ACKNOWLEDGMENTS AND FUNDING SOURCES

The authors received no financial support for the research, writing, and/or publication of this article.

COMPLIANCE WITH ETHICAL STANDARDS

The authors declare that they have no competing interests.

AUTHOR CONTRIBUTIONS

Conceptualization, [A.K.]; methodology, [A.K.]; validation, [P.L.]; formal analysis, [A.K.]; investigation, [A.K., P.L.]; resources, [P.L.]; data curation, [P.L.]; writing – original draft preparation, [A.K.]; writing – review and editing, [P.L.]; visualization, [A.K.].

All authors have read and agreed to the published version of the manuscript.

REFERENCES

- [1] Domingues, M., Faria, J. N., & Portugal, D. (2024). Dimensioning payload size for fast retransmission of MQTT packets in the wake of network disconnections. *EURASIP Journal on Wireless Communications and Networking*, 2024(1). <https://doi.org/10.1186/s13638-023-02327-3>.

- [2] Kamil, A. A., Mousa, A. J., & Abdul-Rahaim, L. A. (2024). Smart Cloud Computing System for Environment Based on MQTT Protocol and Node MCU (ESP8266). *Journal Européen Des Systèmes Automatisés*, 57(5), 1503. <https://doi.org/10.18280/jesa.570526>.
- [3] Shah, P. H. (2024). MQTT Systems: A Survey. *International Journal for Research in Applied Science and Engineering Technology*, 12(3), 1063. <https://doi.org/10.22214/ijraset.2024.59000>.
- [4] Puthiyidam, J.J., & Joseph, S. (2017). IoT Smart Home: Protocols and Architectures. *International Journal for Research in Applied Science and Engineering Technology*, 5, 2031. <https://doi.org/10.22214/ijraset.2017.11293>.
- [5] Eisen, K., Eifert, T., Herwig, C., & Maiwald, M. (2020). Current and future requirements to industrial analytical infrastructure—part 1: process analytical laboratories. *Analytical and Bioanalytical Chemistry*, 412(9), 2027. <https://doi.org/10.1007/s00216-020-02420-2>.
- [6] Arco, E., Boccardo, P., Gandino, F., Lingua, A. M., Noardo, F., & Rebaudengo, M. (2016). An Integrated Approach for Pollution Monitoring: Smart Acquisition and Smart Information. *ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, 67. <https://doi.org/10.5194/isprs-annals-iv-4-w1-67-2016>.
- [7] Gardiner, S., Haynie, C., & Corte, D. D. (2024). Rise of the Allotrope Simple Model: Update from 2023 Fall Allotrope Connect. *Drug Discovery Today*, 29(4), 103944. <https://doi.org/10.1016/j.drudis.2024.103944>.
- [8] Dohare, R. K., Mainuddin, M., Verma, A., & Singhal, G. (2022). Uncertainty evaluation of data acquisition and analysis system relevant to infrared flowing medium laser. *Scientific Reports*, 12(1). <https://doi.org/10.1038/s41598-022-22667-2>.
- [9] Munir, T., Akbar, M. S., Ahmed, S., Sarfraz, A., Sarfraz, Z., Sarfraz, M., Félix, M., & Chérrez-Ojeda, I. (2022). A Systematic Review of Internet of Things in Clinical Laboratories: Opportunities, Advantages, and Challenges. *Sensors*, 22(20), 8051. <https://doi.org/10.3390/s22208051>.
- [10] Knight, N., Kanza, S., Cruickshank, D., Brocklesby, W. S., & Frey, J. G. (2020). Talk2Lab: The Smart Lab of the Future. *IEEE Internet of Things Journal*, 7(9), 8631. <https://doi.org/10.1109/jiot.2020.2995323>.
- [11] Nguyen, H., Nguyen, L. V., & Ha, Q. P. (2020). IoT-enabled Dependable Co-located Low-cost Sensing for Construction Site Monitoring. *Proceedings of the 37th ISARC*, 616-624. <https://doi.org/10.22260/isarc2020/0086>.
- [12] Wicaksono, M. F., & Rahmatya, M. D. (2022). IoT for Residential Monitoring Using ESP8266 and ESP-NOW Protocol. *Jurnal Ilmiah Teknik Elektro Komputer Dan Informatika*, 8(1), 93. <https://doi.org/10.26555/jiteki.v8i1.23616>.
- [13] Rákay, R., Galajdová, A., Šeminský, J., & Cvitić, I. (2019). Selected Wireless Communication Protocols and their Properties for Use in IoT Systems. *Research Papers Faculty of Materials Science and Technology Slovak University of Technology*, 27(45), 26. <https://doi.org/10.2478/rput-2019-0022>.
- [14] Alangcas, A., Marjhon, K., Daligdig, C., & Encarnacion, P. (2024). Embedded IoT Data Collection for Snore Analysis. *International Journal of Advanced Trends in Computer Science and Engineering*, 13(3), 119. <https://doi.org/10.30534/ijatcse/2024/041332024>.
- [15] Babiuch, M., & Postulka, J. (2020). Smart Home Monitoring System Using ESP32 Microcontrollers. In *Internet of things*. Springer Science+Business Media. <https://doi.org/10.5772/intechopen.94589>.
- [16] Le, L. H., Nguyen, D. T., Bui, H. P., Nguyen, C. C., & Ngo, T. Q. (2024). Integration of Modbus Ethernet Communication for Real-Time Electrical Power Consumption,

- Temperature, and Humidity Monitoring System. *Research Square*. <https://doi.org/10.21203/rs.3.rs-4708973/v1>.
- [17] Van Anh, N. T., Anh, N. X. D., Son, L. H., & Khanh, D. Q. (2024, December). A New Architecture for Controlling IoT Devices of Smart Room Using ESP32 Microcontroller. In *International Conference on Smart Objects and Technologies for Social Good* (pp. 106-117). Cham: Springer Nature Switzerland. https://doi.org/10.1007/978-3-032-01472-6_9.
- [18] Abid, A., Jazib, M., & Riaz, M. (2025). RCD-IoT: Enabling Industrial Monitoring and Control with Resource-Constrained Devices Under High Packet Transmission Rates. *arXiv (Cornell University)*. <https://doi.org/10.48550/arxiv.2501.07895>.
- [19] Urazayev, D., Eduard, A., Ahsan, M., & Zorbas, D. (2023, May). Indoor performance evaluation of esp-now. In *2023 IEEE International Conference on Smart Information Systems and Technologies (SIST)* (pp. 1-6). IEEE. <https://doi.org/10.1109/sist58284.2023.10223585>.
- [20] Escobar, J. J. L., Redondo, R. P. D., & Gil-Castiñeira, F. (2024). Unleashing the power of decentralized serverless IoT dataflow architecture for the Cloud-to-Edge Continuum: a performance comparison. *Annals of Telecommunications*, 79, 135–148. <https://doi.org/10.1007/s12243-023-01009-x>.
- [21] Grawe, N. D., & Karaali, G. (2022). Talking about Statistical Significance in Numeracy. *Numeracy*, 15(2), 8. <https://doi.org/10.5038/1936-4660.15.2.1424>.
- [22] Lalchand, V. (2020). A meta-algorithm for classification using random recursive tree ensembles: A high energy physics application. *arXiv (Cornell University)*. <https://doi.org/10.48550/arxiv.2001.06880>.
- [23] Maier, M., & Lakens, D. (2022). Justify Your Alpha: A Primer on Two Practical Approaches. *Advances in Methods and Practices in Psychological Science*, 5(2), 25152459221080396. <https://doi.org/10.1177/25152459221080396>.
- [24] Araujo, A., García-Palacios, J., Blesa, J., Tirado, F., Romero, E., Samartín, A., & Nieto-Taladriz, O. (2011). Wireless Measurement System for Structural Health Monitoring with High Time-Synchronization Accuracy. *IEEE Transactions on Instrumentation and Measurement*, 61(3), 801-810. <https://doi.org/10.1109/tim.2011.2170889>.
-

КОМУНІКАЦІЙНІ АРХІТЕКТУРИ ДЛЯ ХМАРНО-ІНТЕГРОВАНИХ СПЕКТРОМЕТРИЧНИХ ЛАБОРАТОРІЙ: ПОРІВНЯННЯ ESP-NOW ТА MQTT

Андрій Крупич*^{ORCID}, Павло Левуш^{ORCID}

Кафедра радіофізики та комп'ютерних технологій,
Львівський національний університет імені Івана Франка
вул. Тарнавського, 107, 79017 Львів, Україна

АНОТАЦІЯ

Вступ. Хмарно-інтегровані спектрометричні лабораторії стикаються з проблемами зв'язку в забезпеченні доступу до даних та їх аналізу в режимі реального часу. У цьому дослідженні порівнюються два бездротові протоколи, MQTT та ESP-NOW, для керування світлодіодами в таких середовищах. MQTT пропонує легкий, ефективний з точки зору пропускну здатності обмін повідомленнями публікації та підписки [1–3], тоді як ESP-NOW забезпечує бездротовий енергоефективний прямий зв'язок без маршрутизатора Wi-Fi. Метою є оцінка їхньої продуктивності та придатності.

Матеріали та методи. Експериментальна установка включала спектрометр StellarNet, світлодіодні джерела світла та мікроконтролери ESP32. Було протестовано

дві архітектури: 1) Прямий зв'язок MQTT, де кожен ESP32 підключався безпосередньо через Wi-Fi до брокера MQTT; та 2) Керуючий пристрій з ESP-NOW, використовуючи керуючий ESP32 для зв'язку по MQTT/Wi-Fi, а потім передаючи команди через ESP-NOW іншим ESP32. Час відгуку для керування світлодіодами вимірювався протягом 100 циклів, а дані аналізувалися за допомогою описової статистики та t-тесту незалежних вибірок.

Результати. Прямий зв'язок MQTT продемонстрував меншу затримку (медіана ~60 мс) та щільніший розподіл порівняно з ESP-NOW (медіана ~170 мс). T-тест підтвердив статистично значущу різницю ($t = -46,28$), при цьому MQTT продемонстрував швидший час відгуку. Однак ESP-NOW пропонує архітектурні переваги: меншу залежність окремих вузлів від Wi-Fi, гнучкість розгортання з поганим покриттям Wi-Fi, кращу масштабованість [4] та енергоефективність, що робить її вищу затримку практичним компромісом для інтеграції масштабної лабораторії.

Висновки. Прямий MQTT забезпечує кращу продуктивність з низькою затримкою. Однак, керуючий пристрій з ESP-NOW, незважаючи на вищу затримку, є прийнятним рішенням завдяки своїй гнучкості, масштабованості та меншій залежності від Wi-Fi для розподілених лабораторій. Це підкреслює компроміс між абсолютною швидкістю та архітектурними перевагами для надійних аналітичних лабораторій, що працюють у хмарі.

Ключові слова: MQTT, ESP-NOW, спектрометрична лабораторія, Інтернет речей, бездротовий зв'язок, затримка.

UDC: 004.7

EVALUATION OF THE PERFORMANCE OF A MULTI-LEVEL MODEL FOR ANOMALOUS DNS QUERY DETECTION

Andrii Senyk  

Department of Information and Communication Technologies,
Lviv Polytechnic National University,
12 Stepan Bandera Str., Lviv, 79013, Ukraine

Senyk, A. (2026). Evaluation of the Performance of a Multi-Level Model for Anomalous DNS Query Detection. *Electronics and Information Technologies*, 33, 57–70. <https://doi.org/10.30970/eli.33.5>

ABSTRACT

Background. In modern network security systems, DNS (Domain Name System) traffic has become an increasingly attractive vector for covert data exfiltration and command-and-control communication. Existing machine learning methods frequently suffer from limited adaptability to novel attack patterns and an imbalance between detection accuracy and false positive rates.

Materials and Methods. This study proposes TunnelEye, a multi-level detection method for malicious DNS queries that integrates statistical feature analysis, structural n-gram modeling, and anomaly detection. Statistical properties of domain names, including string length, entropy, and alphanumeric ratio, are used for initial discrimination between benign and suspicious queries. Structural analysis based on character n-grams enables the identification of local patterns associated with encoded data such as Base32 and Base64. An autoencoder trained exclusively on legitimate DNS queries is employed as an independent anomaly detector to identify previously unseen and zero-day attacks. The supervised TunnelEye classifier and the autoencoder operate in parallel, each using an independently optimized F1-score based threshold to determine anomalous DNS queries.

Results and Discussion. Experimental evaluation using standard machine learning metrics (precision, recall, F1-score, ROC-AUC, PR-AUC, and false positive rate) demonstrates that TunnelEye consistently outperforms baseline statistical models and standalone autoencoders. The proposed method achieves high precision and recall while maintaining a minimal false positive rate. Experimental results show that TunnelEye achieves an average precision, recall, and F1-score of approximately 0.99, outperforming the baseline statistical model by more than 10% and significantly reducing the false positive rate.

Conclusion. TunnelEye provides a comprehensive and adaptive solution for malicious DNS query detection by combining supervised and unsupervised learning with dynamic threshold optimization. Its ability to balance detection accuracy and false positive reduction makes it well suited for deployment in modern enterprise cybersecurity systems for real-time DNS traffic monitoring.

Keywords: DNS traffic, anomaly detection, machine learning, multi-level model.

INTRODUCTION

In modern computer networks, DNS plays a crucial role in mapping domain names to IP addresses. However, this complex infrastructure is increasingly exploited by malicious actors to create covert communication channels and data tunnels. These techniques allow the transmission of harmful information directly through DNS queries,



© 2026 Andrii Senyk. Published by the Ivan Franko National University of Lviv on behalf of Електроніка та інформаційні технології / Electronics and Information Technologies. This is an Open Access article distributed under the terms of the [Creative Commons Attribution 4.0 License](https://creativecommons.org/licenses/by/4.0/) which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

bypassing traditional security systems and filters. Consequently, there is a growing demand for reliable algorithms to detect these threats in real time [1–4].

Traditional methods for detecting malicious DNS queries typically rely on fixed rules or statistical properties, such as domain name length or the frequency of certain characters. While these methods can detect simple forms of tunneling attacks, they are generally ineffective against sophisticated or adaptive attacks that disguise traffic as legitimate requests. This leads to some malicious queries going unnoticed, increasing the risk of network system compromise. Modern machine learning approaches, such as random forests or neural networks, enable the interpretation of more complex patterns in DNS traffic. However, even these methods face generalization challenges: models trained on historical data may be ineffective against new types of tunneled queries. This highlights the need for adaptive, multi-leveled approaches that consider both statistical and structural properties of domain names. A key challenge is the lack of high-quality training data for malicious queries. Known data collection methods are often limited or do not include the latest tunneling techniques, making it difficult to develop models capable of detecting zero-day attacks (i.e., previously unseen attacks). Therefore, the goal of this research is to create a system that can learn from reliable data and detect anomalies in new queries.

Another important issue is the balance between detection accuracy and completeness. Using fixed thresholds for classification often results in high rates of false positives or false negatives. For practical security systems, it is essential to develop an adaptive threshold adjustment mechanism that maintains an optimal balance under varying network traffic conditions. The structural properties of domain names, such as n-gram character sequences, remain underutilized in many traditional methods. These features allow the detection of local and internal patterns in malicious domain names, particularly when encoded using Base32 or Base64. Combining these features with statistical properties can significantly enhance the effectiveness of detection systems.

An additional challenge is ensuring scalability and real-time monitoring capabilities. Enterprise network systems process enormous volumes of DNS queries; therefore, algorithms must be efficient enough to operate in large-scale networks without introducing network delays. This requires an approach that combines accuracy with low computational complexity and adaptability to fluctuating loads [5–7]. Analysis of existing studies shows that most methods for detecting malicious DNS queries either rely on fixed statistical rules or use standalone machine learning models without adaptive mechanisms.

This work proposes the TunnelEye method, which integrates known components (statistical analysis, n-grams, and an autoencoder) into a novel decision-making architecture. The novelty lies in the way these components are constructed and used within a unified experimental environment. The main modifications include: integration of statistical and n-gram features into a single feature model; use of an ensemble Random Forest classifier instead of heuristic rules; inclusion of an autoencoder as an independent anomaly detector; implementation of an adaptive classification threshold selection mechanism based on the F1-score.

MATERIALS AND METHODS

Analysis of Problems and Formulation of the Research Task

This study proposes a multi-leveled adaptive method, TunnelEye, which implements a multi-leveled model for detecting malicious DNS queries based on statistical, structural, and anomalous features. This combination addresses the limitations of traditional methods, creating a more reliable and effective monitoring system.

In contrast to existing approaches to DNS tunneling detection, which typically focus either on statistical properties of domain names (such as length, entropy, or character

distribution) or on the isolated application of machine learning techniques, the proposed TunnelEye method integrates multiple levels of analysis within a unified architecture. Statistical models (e.g., rule-based approaches or Random Forest classifiers using simple numerical features) are effective at detecting primitive forms of tunneling, but they exhibit poor generalization in more sophisticated obfuscation scenarios. Methods based on n-gram analysis are capable of identifying local structural patterns (e.g., Base32 or Base64 encodings); however, in the absence of broader statistical context, they often suffer from elevated false-positive rates.

Autoencoder-based approaches commonly employed for anomaly detection in DNS traffic typically operate as standalone solutions and rely exclusively on reconstruction error. This reliance limits their accuracy in cases where malicious traffic closely mimics legitimate behavior. In contrast to these approaches, TunnelEye implements a parallel architecture in which a supervised classifier (utilizing statistical and n-gram features) and an autoencoder operate independently but are evaluated within a unified experimental framework with adaptive threshold optimization. As a result, the proposed method combines the strengths of signature-based, statistical, and anomaly-based analysis while mitigating their individual limitations.

Key Concepts in Malicious DNS Query Detection

DNS forms the foundation of the Internet by translating domain names into IP addresses. This system allows users to access websites and services using human-readable names rather than numerical addresses. However, due to its widespread use, malicious actors often exploit DNS for covert data exchange and tunneling, making DNS traffic monitoring critically important for network security.

Malicious DNS queries vary in nature, ranging from simple spoofing of legitimate domain names to sophisticated DNS tunneling, where data is encoded within a domain name and transmitted via the query. These queries are difficult to detect using traditional filtering methods because they can appear as normal, legitimate traffic, increasing the risk of hidden attacks. A key concept in detecting malicious DNS queries is the use of analytical signatures (features that can distinguish legitimate queries from malicious ones). These features include statistical properties of domain names, such as string length, entropy, the number of domain levels, and the ratio of digits to letters.

These properties allow the assessment of the randomness and complexity of domain names, which are often present in malicious queries. Additionally, processing large volumes of DNS queries in real time requires high computational efficiency and algorithmic optimization [9–12]. The general scheme for detecting malicious DNS queries is shown in **Fig. 1**.

Proposed Monitoring Method

The proposed method, named TunnelEye, is designed to detect covert communication channels in DNS traffic, specifically data tunneling through domain name queries. TunnelEye is developed as a standalone implementation for DNS tunneling detection and is used for comparison with a baseline model and an autoencoder. Unlike traditional methods, it combines several independent analysis mechanisms, enhancing the system's resilience to various obfuscation techniques. At the first level, the method employs traditional statistical features of domain names, including length, entropy, the number of domain levels, and the ratio of digits to letters. These features accurately describe the overall complexity and randomness of a string and can often distinguish legitimate domain names from those generated or used for tunneling.

The second level of the method relies on n-gram character analysis, which can detect local structural patterns in domain names. This allows the effective identification of malicious patterns related to data encodings (such as Base32 or Base64) that are typically not detectable through statistical analysis alone.

The third component of the system is an autoencoder trained specifically on “clean” samples of legitimate domain names. It acts as an anomaly detector, using reconstruction error as an indicator of deviation. This enables the detection of new types of malicious queries that were not present in the training datasets, ensuring the method’s robustness against zero-day attacks.

Furthermore, the method provides an adaptive classification threshold mechanism. Instead of using fixed values, the system uses the F1 score to determine the optimal threshold, balancing precision and recall based on specific deployment conditions. This improves the practical effectiveness of the system across a wide range of scenarios.

Thanks to this architecture, TunnelEye can serve as a multi-levelled DNS query monitoring platform and can be integrated into enterprise security systems for real-time network traffic analysis. This approach combines signature accuracy, statistical methods, and anomaly detection flexibility to provide comprehensive protection. Therefore, TunnelEye can be defined as a hybrid method for detecting covert channels in DNS. It integrates statistical analysis, n-gram-based structural analysis, and anomaly detection via an autoencoder with adaptive classification thresholds, enhancing network security capabilities and effectively detecting both known and previously unseen malicious activity.

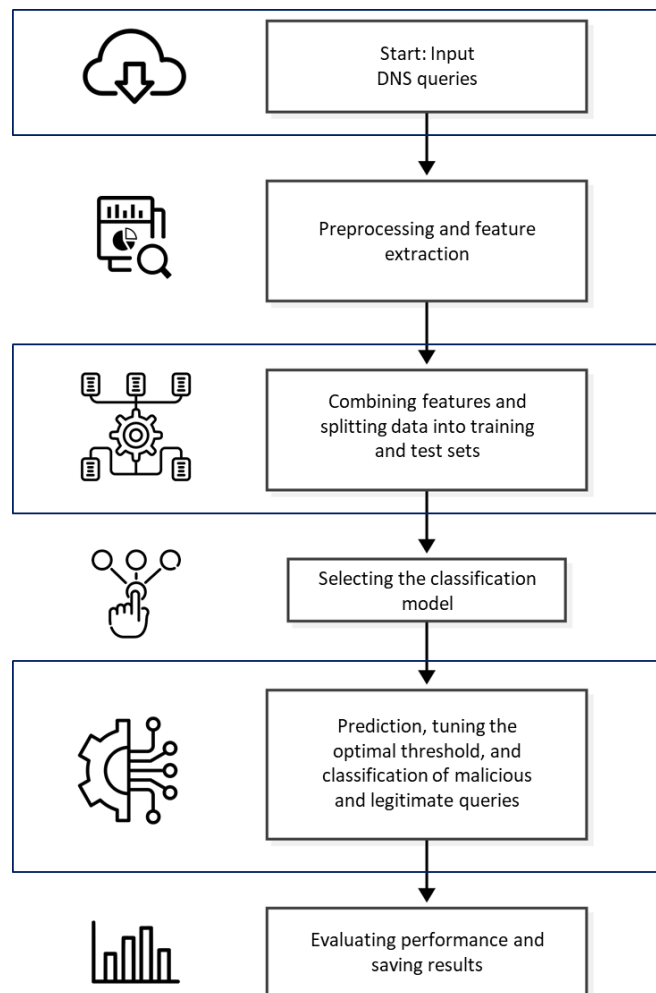


Fig. 1. General Scheme for Malicious DNS Query Detection.

Scientific Novelty of the Proposed Method

The scientific novelty of the method lies in the combination of statistical and symbolic n-gram features within a single detection framework; the parallel use of supervised and anomaly-based processing of DNS queries; and the study of the impact of the classification threshold on the effectiveness of malicious DNS query detection. The autoencoder is used as the anomaly detection mechanism, trained exclusively on legitimate DNS queries. This enables the detection of entirely new types of malicious patterns that are not present in the training samples. The block diagram of the proposed monitoring method is shown in **Fig. 2**.

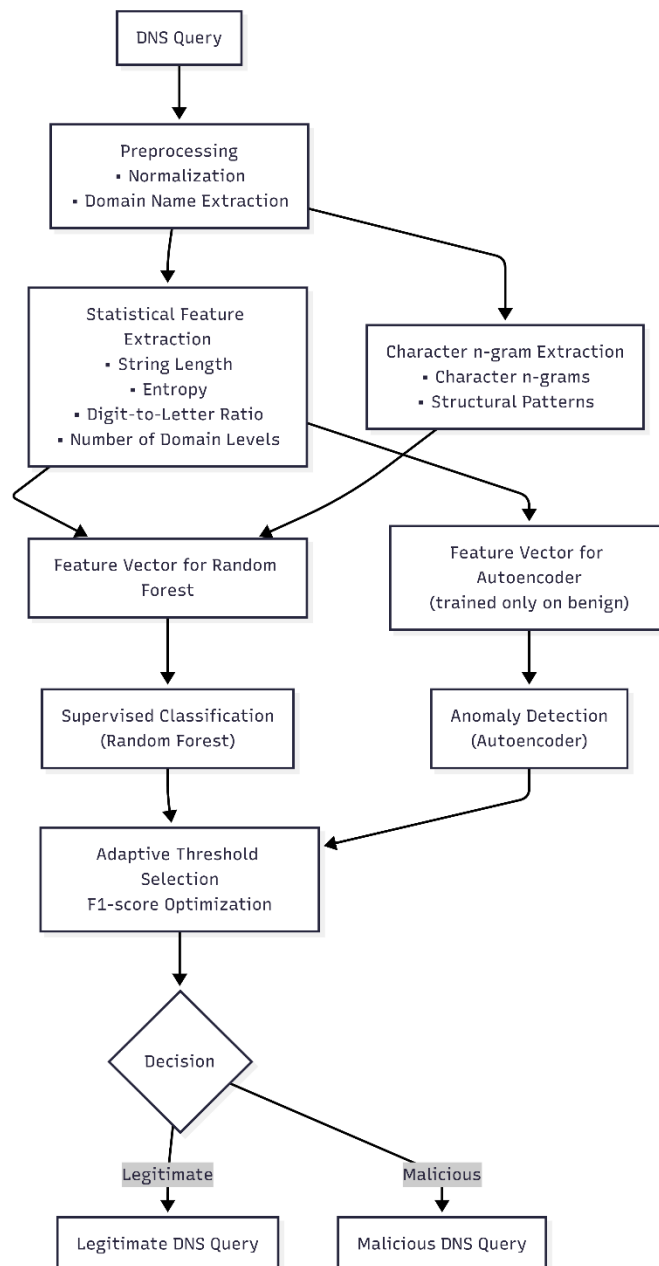


Fig. 2. Block Diagram of the Proposed Monitoring Method.

Another innovation is the adaptive classification threshold mechanism based on F1-score optimization. Most studies in this field rely on a fixed threshold (0.5), which is often suboptimal in real-world applications. F1-score optimization involves evaluating different threshold values and selecting the one that maximizes the F1-score on test or validation data. In this case, classification involves recognizing DNS queries as either legitimate (0) or malicious (1). The optimization parameter is the threshold ranging from 0 to 1, which converts probabilities (from the supervised model or autoencoder reconstruction error) into binary predictions. The method provides systematic threshold adjustment by computing the F1-score for different threshold values, automatically finding the optimal balance between precision and recall for malicious DNS query detection. System performance is evaluated using standard machine learning metrics, including F1, ROC-AUC, and the false positive rate.

For a comprehensive comparison of methods, multiple metrics are employed: precision, recall, F1-score, ROC-AUC, and PR-AUC. This approach provides a deeper understanding of each model's strengths and weaknesses and establishes an objective basis for decision-making in real-world applications.

To assess the method's effectiveness, a synthetic DNS query dataset was used, simulating various types of legitimate traffic and DNS tunneling with controlled characteristics. The synthetic data were used solely for method evaluation, not as part of the method itself. This creates an experimental environment for testing models in scenarios close to real-world conditions, where both legitimate and malicious queries exhibit high variability. This approach not only improves evaluation quality but also allows testing the robustness of the algorithms against spoofing attacks.

The block diagram illustrates the sequence of stages: preprocessing DNS queries, extracting statistical and structural features, parallel operation of the classifier and autoencoder, adaptive threshold adjustment, and final classification.

A key component of this method is the entropy function, which measures the randomness of a string. For a domain s , entropy is defined as:

$$H(s) = - \sum_{i=1}^n p_i \times \log_2 p_i \quad (1)$$

where p_i is the probability of a character i appearing in the string s , and n is the number of unique characters.

The model is evaluated using standard machine learning metrics. *Precision* is defined as:

$$Precision = \frac{TP}{TP + FP} \quad (2)$$

where TP represents the number of correctly classified malicious examples, and FP represents the number of legitimate examples incorrectly classified as malicious.

Recall is defined as:

$$Recall = \frac{TP}{TP + FN} \quad (3)$$

where FN represents the number of malicious examples missed by the model.

The F1-score is the harmonic mean of precision and recall, calculated as:

$$F1 = 2 \frac{Precision \times Recall}{Precision + Recall} \quad (4)$$

Additionally, integral classification quality metrics were computed, including the area under the ROC curve (ROC-AUC) and the area under the Precision-Recall curve (PR-AUC). The false positive rate (FPR) was also analyzed:

$$FPR = \frac{FP}{FP + TN} \quad (5)$$

where TN represents the number of correctly classified legitimate examples.

Characteristics of the Conducted Experiments

In this study, a synthetically generated dataset was used to simulate queries to the DNS from both legitimate and malicious sources. Legitimate queries included typical domain names, such as corporate services, Content Delivery Networks (CDNs), Universally Unique Identifiers (UUIDs), and random combinations of letters and digits simulating real user requests. Malicious examples were generated by encoding random strings using Base32 and Base64 encodings, hexadecimal representations, and hashing algorithms, corresponding to common DNS tunneling techniques. This data generation process resulted in a balanced dataset of 16,000 records (50% benign and 50% malicious).

The study also included high-entropy domain names that were not malicious but mimicked characteristic features of DNS tunneling to make the experiments more realistic and representative of real-world conditions. Specifically, some legitimate domain names contained random sequences of letters, long strings, and digits visually similar to malicious domains. Consequently, models needed to learn to distinguish not only simple patterns but also complex obfuscation patterns. This allowed us to assess the robustness of the TunnelEye approach against attacks employing both covert encoding and obfuscation imitation. The program was implemented in Python using the scikit-learn, TensorFlow, and matplotlib libraries. Statistical analysis employed string processing and traditional numerical properties, including domain name length, entropy, number of subdomain levels, and character ratios. Structural analysis relied on character n-gram features generated by CountVectorizer, ranging from 3 to 5 characters and limited to the top 3,000 most frequent characters.

For comparison, a baseline model using only statistical features and an autoencoder (a multilayer neural network consisting of a compression encoder and a symmetric decoder layer) were also employed. The baseline statistical model was used exclusively for comparative analysis, whereas TunnelEye implements a dynamic multi-level model capable of adapting to new traffic types.

All computations were performed in the Google Colab environment, enabling efficient use of computational resources for model training and easy replication of experimental results. The obtained metrics were used to compare the baseline statistical model, TunnelEye, and the autoencoder across different performance measures, confirming the scientific novelty of the method.

The supervised TunnelEye classifier and the autoencoder operate in parallel, each producing an independent anomaly score, while optimal decision thresholds are selected separately for each model using F1-score optimization. The autoencoder is not used as a direct decision fusion component but serves as an independent anomaly detector for comparative evaluation and robustness analysis against previously unseen attacks. The proposed architecture allows straightforward extension toward score-level or decision-

level fusion, which is left for future work. The results of the experiments are presented in the Results section.

RESULTS AND DISCUSSION

It is important to note that the method produces a binary outcome (anomalous DNS query detected or not detected). The metrics are used solely to quantitatively assess the quality of this decision and to compare different approaches. The results confirm that TunnelEye's effectiveness is not due to individual components alone but results from their coordinated integration. The autoencoder strengthens the model in scenarios where supervised classification loses sensitivity, while the adaptive threshold ensures stability in real-world conditions. Thus, the experimental study fully meets the stated objectives, and the results demonstrate the completeness and practical relevance of the proposed method.

Fig. 3 illustrates how precision and recall change with varying classification thresholds. This allows us to assess the balance between detecting all attacks (recall) and minimizing false positives (precision).

Results:

- TunnelEye: Average precision ≈ 0.99 , recall ≈ 0.99
- Baseline model: Average precision ≈ 0.88 , recall ≈ 0.90
- Autoencoder: Average precision ≈ 0.65 , recall ≈ 0.48

Thus, TunnelEye shows minimal performance degradation across threshold values, demonstrating its stability. The baseline algorithm performs slightly worse, while the autoencoder exhibits fluctuations and is less balanced.

The F1-score (**Fig. 4**) summarizes precision and recall, reflecting model quality across different threshold values. This metric is critical for finding the optimal balance between errors and false positives.

Results:

- TunnelEye: Maximum F1 = 1.0 (threshold ≈ 0.06), average F1 ≈ 0.99
- Baseline model: Maximum F1 ≈ 0.86 , average F1 ≈ 0.84
- Autoencoder: Maximum F1 ≈ 0.52 , average F1 ≈ 0.49

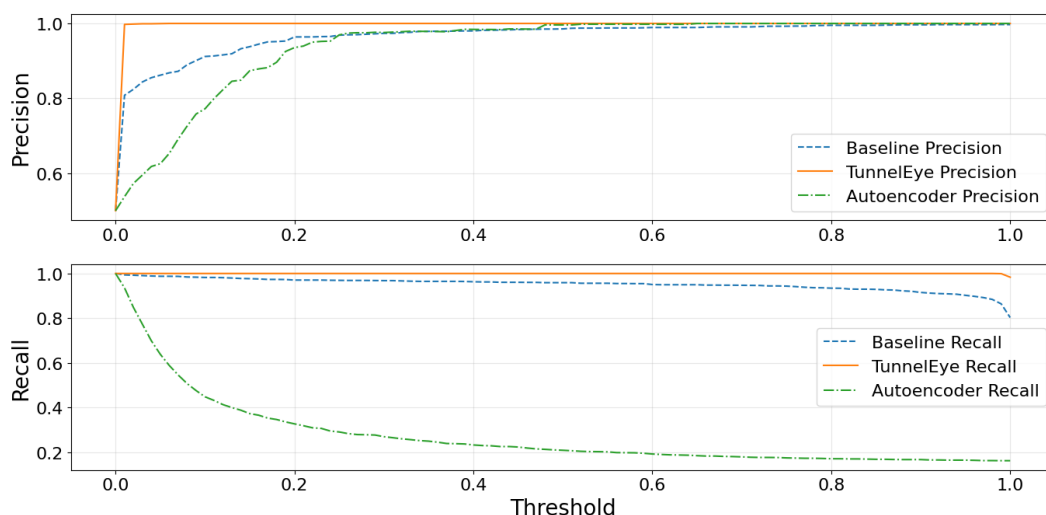


Fig. 3. Precision and Recall vs. Classification Threshold.

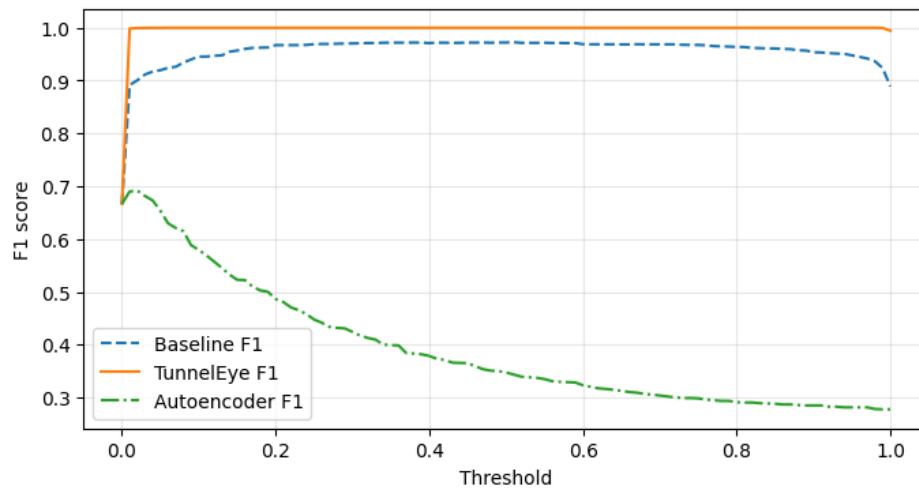


Fig. 4. F1-score vs. Classification Threshold.

TunnelEye significantly outperforms the other models, maintaining the highest overall performance. The baseline model is stable but not optimal, while the autoencoder performs substantially worse, particularly in terms of recall.

The ROC curve (Fig. 5a) measures the ratio of true positives to false positives, while the PR curve (Fig. 5b) shows the balance between precision and recall, which is particularly important for imbalanced classes. A large area under the curve (AUC) indicates a stable and effective model.

Results:

- TunnelEye: ROC-AUC = 1.0, PR-AUC = 1.0
- Baseline model: ROC-AUC \approx 0.992, PR-AUC \approx 0.991
- Autoencoder: ROC-AUC \approx 0.722, PR-AUC \approx 0.768

Thus, TunnelEye demonstrates the best performance, with the other models lagging behind.

Fig. 6 presents a “performance metric” that evaluates all models at the same threshold. This metric compares precision, recall, F1-score, ROC-AUC, PR-AUC, FPR, as well as the absolute values of TP, FP, TN, and FN.

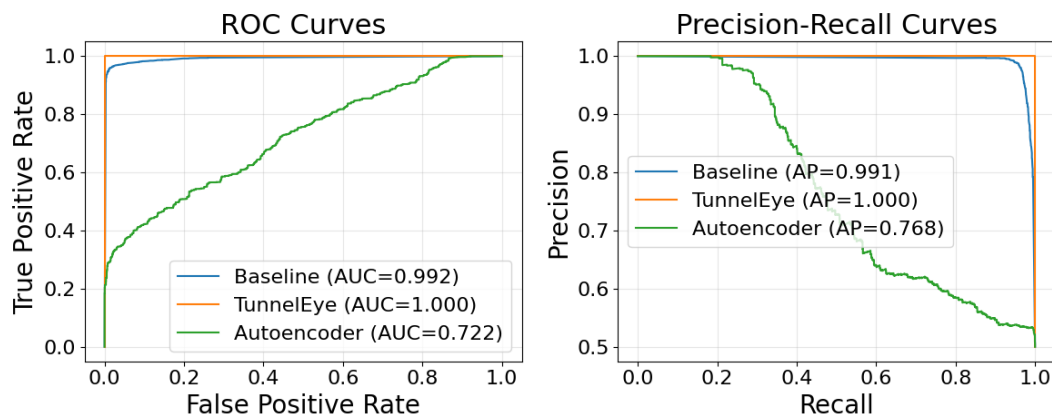


Fig. 5. ROC (a) and PR (b) curves.

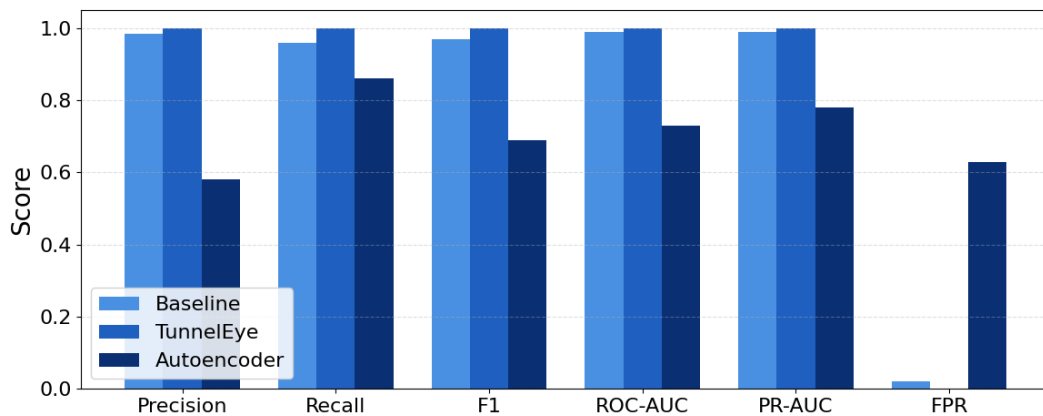


Fig. 6. Model Comparison at Fixed Threshold (0.5).

According to the comparative analysis (Fig. 6), TunnelEye achieved the best results (Precision, Recall, F1, ROC-AUC, PR-AUC = 1.0; FPR = 0), while the baseline model achieved moderate metrics (Precision \approx 0.815, Recall \approx 0.91, F1 \approx 0.86, ROC-AUC \approx 0.944, PR-AUC \approx 0.936, FPR \approx 0.197).

The autoencoder showed significantly lower metrics (Precision \approx 0.72, Recall \approx 0.41, F1 \approx 0.52, ROC-AUC \approx 0.65, PR-AUC \approx 0.61, FPR \approx 0.32), making it less suitable for standalone use. Therefore, TunnelEye can be considered an effective method for detecting DNS tunneling and similar attacks.

The experiments were conducted on a synthetically generated dataset with controlled characteristics. This approach allows for a fair comparison of models under identical conditions; however, it does not fully capture the complexity of real DNS traffic. To address this limitation, a public DNS dataset containing real domain names with naturally occurring feature distributions and inherent variability was used (Fig. 7). daumel/dns-tunneling-dataset is a public dataset containing DNS query traffic with tunneling, generated by various DNS tunneling tools.

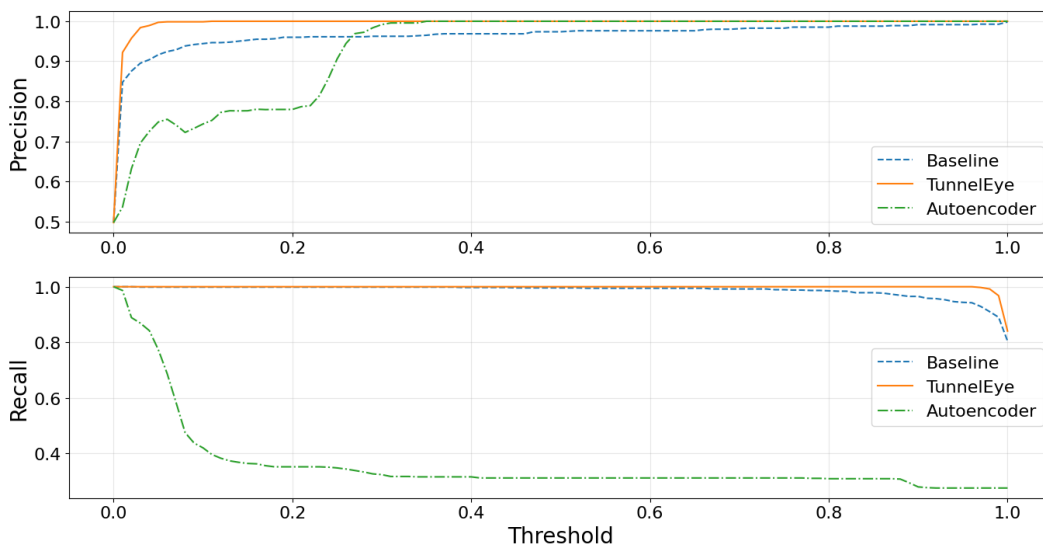


Fig. 7. Precision and recall versus classification threshold on the real DNS traffic dataset.

On this dataset, perfect metric values are not achieved; nevertheless, TunnelEye consistently demonstrates a better balance between precision and recall compared to baseline models. Therefore, the results confirm that the effectiveness of the proposed method is not a consequence of the artificial simplicity of synthetic data but is maintained under realistic DNS traffic conditions. Future studies will expand the experiments to other datasets to assess the method's performance across varied and realistic DNS traffic conditions.

Achieving ideal ROC-AUC and PR-AUC values in a controlled experiment should be interpreted as an indicator of the method's effectiveness under clearly separable class conditions, rather than as a universal measure of its behavior in real networks. Additional experiments on a real DNS dataset show that, in the presence of more complex mimicry scenarios and partial feature overlap between legitimate and malicious domains, the metrics become less ideal. Nevertheless, TunnelEye maintains an advantage over baseline statistical models and the autoencoder. This demonstrates the method's resilience to more sophisticated attacks and confirms that its effectiveness is not limited to synthetic scenarios.

To extend the comparative analysis, additional modern machine learning methods were also evaluated, including Logistic Regression and Linear SVM with character n-gram features, as well as Isolation Forest as a representative anomaly-based model. However, in the experiments, these methods did not demonstrate any advantage over TunnelEye in the context of DNS traffic analysis and were therefore not included in the corresponding plots.

CONCLUSION

In this paper, the TunnelEye monitoring method is proposed, which combines statistical methods, structural analysis, and an autoencoder for anomaly detection, demonstrating exceptional effectiveness in identifying malicious DNS queries. This combination allows the system to handle both known patterns and new, previously unseen malicious queries. Experimental studies show that TunnelEye achieves the highest precision and recall, significantly outperforming baseline statistical and autoencoder models. This demonstrates the effectiveness of a multi-stage approach and its applicability in real-world conditions with highly variable DNS traffic.

The scientific innovation lies in the method's ability to adaptively adjust classification thresholds and integrate multiple detection techniques for identifying malicious queries. This ensures high performance even under varying operational conditions, making it well-suited for integration into enterprise network traffic monitoring systems. Comparative analysis of different models demonstrates that TunnelEye not only provides high precision and recall but also reduces false positive rates, which is critical for practical deployment.

Experimental results indicate that TunnelEye achieves an average anomaly detection precision and recall close to 0.99, significantly surpassing both the baseline statistical model and the autoencoder. This highlights the advantage of the multi-stage method in detecting both known and novel malicious patterns. The F1-score further confirms TunnelEye's high performance (average F1 \approx 0.99). These results demonstrate TunnelEye's reliability and its ability to minimize false positives. At a fixed threshold of 0.5, TunnelEye also showed high precision and recall. This comparison confirms that the proposed approach outperforms both traditional methods and autoencoders, achieving superior results across all metrics.

Overall, the comparative analysis demonstrates that TunnelEye is an effective method for detecting DNS tunnels and covert communication channels. Its multi-level adaptive approach ensures high accuracy, protects against zero-day attacks, and reduces false alarms, making it suitable for integration into corporate security and real-time network monitoring systems. TunnelEye can be considered an efficient and flexible

tool for detecting malicious DNS queries and DNS tunneling. Its multi-level adaptive design provides comprehensive network protection and opens avenues for further research in the detection of covert channels and anomalous traffic.

ACKNOWLEDGMENTS AND FUNDING SOURCES

The author received no financial support for the research, writing, and publication of this article.

CONFLICT OF INTEREST

The author declares that the research was conducted in the absence of any.

AUTHOR CONTRIBUTIONS

The author has read and agreed to the published version of the manuscript.

REFERENCES

- [1] Gonzalez Casanova, L. F., & Lin, P. C. (2021). Generalized classification of DNS over https traffic with deep learning. In *2021 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, Tokyo, Japan, 2021, pp. 1903–1907. <https://ieeexplore.ieee.org/document/9689667>
- [2] Ichise, H., Jin, Y., & Iida, K. (2023). Policy-based detection and blocking system against abnormal applications by analyzing DNS traffic. In *2023 22nd International Symposium on Communications and Information Technologies (ISCIT)*, Sydney, Australia, 2023, pp. 1–6. <https://doi.org/10.1109/ISCIT57293.2023.10376042>
- [3] Zhang, C., Hu, X., Pan, X., Cheng, G., Li, R., & Wu, H. (2025). Accurate and early detection of IoT malware via DNS traffic analysis with deep learning. In *ICC 2025 – IEEE International Conference on Communications, Montreal, QC, Canada, 2025*, pp. 2665–2670. <https://doi.org/10.1109/ICC52391.2025.11161323>
- [4] Ganesh, N., Parihar, A. S., & Ghosh, G. (2023). Analysing network traffic and implementing diverse technologies to examine different components of the network. In *2023 IEEE International Conference on ICT in Business Industry & Government (ICTBIG)*, Indore, India, 2023, pp. 1–10. <https://doi.org/10.1109/ICTBIG59752.2023.10456258>
- [5] Harishkumar, S., & Bhuvaneshwaran, R. S. (2024). Unveiling domain generation algorithms in DNS log traffic: A next-generation intelligent framework for dynamic anomaly detection and mitigation through machine learning analysis. In *2024 15th International Conference on Computing Communication and Networking Technologies (ICCCNT)*, Kamand, India, 2024, pp. 1–7. <https://doi.org/10.1109/ICCCNT61001.2024.10726248>
- [6] Wu, X., Wang, X., Song, Y., & Ding, P. (2024). SSPT: A self supervised network traffic anomaly detection method. In *2024 20th International Conference on Mobility, Sensing and Networking (MSN)*, Harbin, China, 2024, pp. 1206–1207. <https://doi.org/10.1109/MSN63567.2024.00177>
- [7] Zou, F., Ren, Y., Zhu, J., & Tang, J. (2021). Detecting data leakage in DNS traffic based on time series anomaly detection. In *2021 IEEE 23rd International Conference on High Performance Computing & Communications; 7th International Conference on Data Science & Systems; 19th International Conference on Smart City; 7th International Conference on Dependability in Sensor, Cloud & Big Data Systems & Application (HPCC/DSS/SmartCity/DependSys)*, Haikou, Hainan, China, 2021, pp. 503–510. <https://doi.org/10.1109/HPCC-DSS-SmartCity-DependSys53884.2021.00090>
- [8] Du, X., et al. (2022). Design of an autoencoder-based anomaly detection for the DoH traffic system. In *2022 IEEE 25th International Conference on Computer Supported*

- Cooperative Work in Design (CSCWD)*, Hangzhou, China, 2022, pp. 763–768.
<https://doi.org/10.1109/CSCWD54268.2022.9776029>
- [9] Hzami, M., Mahersia, H., & Bejaoui, T. (2025). Multi-level cyberbullying detection on social media using machine and deep learning models. In *2025 5th IEEE Middle East and North Africa Communications Conference (MENACOMM)*, Byblos, Lebanon, 2025, pp. 1–6. <https://doi.org/10.1109/MENACOMM62946.2025.10911024>
- [10] Huang, X., Zhu, X., Xu, X., Zhu, M., Nian, A., & Guo, Y. (2022). Multi-granularity perceptual ensemble learning model with an application. In *2022 International Conference on Machine Learning, Cloud Computing and Intelligent Mining (MLCCIM)*, Xiamen, China, 2022, pp. 234–242.
<https://doi.org/10.1109/MLCCIM55934.2022.00047>
- [11] Wang, B., Xiong, G., Gou, G., Song, J., Li, Z., & Yang, Q. (2023). Identifying DoH tunnel traffic using core features and machine learning method. In *2023 26th International Conference on Computer Supported Cooperative Work in Design (CSCWD)*, Rio de Janeiro, Brazil, 2023, pp. 814–819.
<https://doi.org/10.1109/CSCWD57460.2023.10152678>
- [12] Rana, S., & Aksoy, A. (2021). Automated fast-flux detection using machine learning and genetic algorithms. In *IEEE INFOCOM 2021 – IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPS)*, Vancouver, BC, Canada, 2021, pp. 1–6. <https://doi.org/10.1109/INFOCOMWKSHPS51825.2021.9484614>
-

ОЦІНЮВАННЯ ПРОДУКТИВНОСТІ БАГАТОРІВНЕВОЇ МОДЕЛІ ДЛЯ ВИЯВЛЕННЯ АНОМАЛЬНИХ DNS-ЗАПИТІВ

Андрій Сенік 

Національний університет «Львівська політехніка»,
вул. Степана Бандери, 12, Львів, 79013, Україна

АНОТАЦІЯ

Вступ. У сучасних системах мережевої безпеки DNS-трафік дедалі частіше використовується для прихованої передачі даних і керування шкідливими системами (command-and-control). Існуючі методи машинного навчання нерідко мають обмежену здатність адаптації до нових шаблонів атак і стикаються з дисбалансом між точністю виявлення та рівнем хибнопозитивних спрацьовувань.

Матеріали та методи. Запропоновано багаторівневий метод виявлення зловмисних DNS-запитів TunnelEye, що інтегрує аналіз статистичних ознак, структурне n-грамне моделювання та виявлення аномалій. Для початкового розмежування легітимних і підозрілих запитів використовуються статистичні властивості доменних імен, зокрема довжина рядка, ентропія та співвідношення буквено-цифрових символів. Структурний аналіз на основі символічних n-грам дає змогу ідентифікувати локальні шаблони, пов'язані з кодуванням даних, таким як Base32 і Base64. Супервізований класифікатор TunnelEye та автокодер працюють паралельно, кожен із яких використовує незалежно оптимізований поріг на основі F1-міри для визначення аномальних DNS-запитів.

Результати. Експериментальна оцінка з використанням стандартних метрик машинного навчання (precision, recall, F1-score, ROC-AUC, PR-AUC та рівень хибнопозитивних спрацьовувань) показує, що TunnelEye стабільно перевершує базові статистичні моделі та окремі автокодери. Запропонований метод забезпечує високу точність і повноту виявлення за мінімального рівня хибнопозитивних спрацьовувань. Експериментальні результати показали, що метод TunnelEye

забезпечує середні значення precision, recall та F1-міри на рівні близько 0,99, перевищуючи базову статистичну модель більш ніж на 10% та істотно зменшуючи кількість хибнопозитивних спрацьовувань.

Висновки. TunnelEye пропонує комплексне та адаптивне рішення для виявлення зловмисних DNS-запитів шляхом поєднання керованого та некерованого навчання з динамічною оптимізацією порогів. Здатність методу збалансувати точність виявлення та зменшення кількості хибнопозитивних спрацьовувань робить його придатним для впровадження в сучасні корпоративні системи кібербезпеки для моніторингу DNS-трафіку в реальному часі.

Ключові слова: DNS-трафік, виявлення аномалій, машинне навчання, багаторівнева модель

UDC 004.94

APPLICATION OF PENETRATION TESTING FOR ASSESSING THE INFORMATION SECURITY LEVEL OF WEB-ORIENTED INFORMATION SYSTEMS

Sergiy Sveleba¹^{*}, Ivan Katerynychuk¹, Ivan Kunyo¹,
Oleh Krupych¹, Yaroslav Shmyhelskyy¹, Marta Dufanets¹,
Natalia Sveleba², Lucjan Pelc³, Volodymyr Brygilevych³

¹Ivan Franko National University of Lviv,

107, Gen. Tarnavsky St., 79017 Lviv, Ukraine

²Private Higher Education Establishment "European University",

16B, Academician Vernadsky Boulevard, 03115 Kyiv, Ukraine

³State Academy of Applied Sciences in Jarosław

ul. Czarnieckiego 16, 37-500 Jarosław Poland

Sveleba, S., Katerynychuk, I., Kunyo, I., Krupych, O., Shmyhelskyy, Ya., Dufanets, M., Sveleba, N., Pelc, L., Brygilevych, V. (2026). Application of Penetration Testing for Assessing the Information Security Level of Web-Oriented Information Systems. *Electronics and Information Technologies*, 33, 71–86. <https://doi.org/10.30970/eli.33.6>

ABSTRACT

Background. The increasing role of web-oriented information systems in business, education, and public administration is accompanied by a growing number and complexity of cyber threats. Traditional security mechanisms do not always enable the identification of actual system weaknesses, which necessitates the application of practice-oriented methods for assessing the level of information security. In this context, Penetration Testing is considered an effective instrument for simulating the actions of a potential attacker in order to detect and validate exploitable vulnerabilities.

Materials and Methods. The study employs a risk-oriented approach in accordance with international standards ISO/IEC 27001 and ISO/IEC 27005, as well as the recommendations of OWASP and NIST SP 800-115. Penetration Testing is implemented as a structured, multi-stage process that includes information gathering, attack surface analysis, threat modeling, execution of non-invasive validation scenarios, and risk assessment. The practical component was conducted in a controlled test environment using Nmap, Burp Suite, and Wireshark, supplemented by custom-developed Python modules for automated analysis of HTTP security headers, TLS certificates, and exposed services.

Results and Discussion. The study identified several configuration-related weaknesses at the application level, including the absence of essential HTTP security headers and deficiencies in TLS certificate management. The obtained results were formalized in a structured findings register with quantitative risk evaluation based on the *Likelihood* × *Impact* model. The analysis demonstrated that even in the absence of critical exploitable vulnerabilities, configuration errors significantly increase the overall risk level and may create preconditions for more sophisticated attacks.

Conclusion. The findings confirm the effectiveness of Penetration Testing as a comprehensive instrument for assessing the information security of web-oriented systems. The proposed approach facilitates the transition from technical testing results to substantiated managerial decisions aimed at risk reduction and enhancement of the overall protection level of information resources.



Keywords: information security, penetration testing, vulnerabilities, risk assessment, web-oriented information systems.

INTRODUCTION

The digitalization of business processes and the rapid expansion of web-oriented information systems across various domains have led to increased requirements for ensuring information security [1]. At the same time, there is a steady growth in both the number and sophistication of cyberattacks aimed at compromising the confidentiality, integrity, and availability of information resources [2]. Traditional protection mechanisms, including antivirus solutions and firewalls, do not always enable the timely identification of actual weaknesses in information systems, particularly those related to configuration errors or improper implementation of security mechanisms [3]. In this context, Penetration Testing is regarded as an effective method for assessing the level of security by simulating the actions of a potential attacker, thereby enabling the transition from formal analysis to practical validation of vulnerability exploitability [4].

The principal threats to information security are traditionally classified according to the CIA triad model, which encompasses confidentiality, integrity, and availability of information [5]. Threats to confidentiality involve unauthorized access to data; threats to integrity relate to data modification or tampering; and threats to availability concern disruption or denial of service affecting information systems [6].

Configuration errors, the use of outdated or weak cryptographic mechanisms, and the absence or insufficient formalization of security policies are particularly critical factors [7]. The human factor also plays a significant role in risk formation, including the use of weak passwords, administrative errors, and the impact of social engineering techniques [8].

Penetration Testing is defined as a method for evaluating the level of information security through the simulation of real attacks on an information system [9]. This approach involves modeling the actions of a potential attacker in order to assess the system's resilience to relevant threats. Unlike purely automated vulnerability scanning, Penetration Testing enables the identification of weaknesses that can be realistically exploited and allows for the evaluation of the effectiveness of existing technical and organizational security controls [10].

A comparison between Vulnerability Scanning and Penetration Testing demonstrates that the latter is based on a combination of automated and manual analytical methods [11]. This integrated approach reduces the number of false positives and provides a more objective assessment of the security posture of an information system [12].

Information security assessment should be considered a sequential process encompassing the chain "protection – attack – vulnerability – risk – decision" [13]. Within this logical framework, Penetration Testing provides a comprehensive and practice-oriented evaluation of the security level and forms a substantiated basis for managerial decision-making aimed at enhancing the protection of information systems [14].

MATERIALS AND METHODS

Penetration Testing is a comprehensive method for assessing the security level of information systems based on the simulation of real-world attacks [15]. Unlike automated vulnerability scanning, this approach enables validation of the practical exploitability of identified weaknesses and allows for an assessment of their actual impact on system security [16].

The study relies on international standards ISO/IEC 27001 and ISO/IEC 27005, which define requirements for information security management systems and risk assessment processes [17-18]. The OWASP Testing Guide and OWASP Top 10 methodologies were applied for web application security analysis [19-20]. The Penetration Testing Execution Standard (PTES) was used to describe the full penetration testing lifecycle

[21], while NIST SP 800-115 provided guidance on the technical aspects of security testing [22]. The application of these standards ensured the relevance, completeness, and alignment of the threat model with contemporary information security conditions.

In this work, Penetration Testing was implemented as a structured, multi-stage process in accordance with OWASP, PTES, and NIST methodologies. The initial phase involved information gathering, including the identification of assets, domains, services, and open ports. This was followed by attack surface analysis to determine potential entry points and system weaknesses. Based on the developed threat model, attack scenarios were formulated and validated through the simulation of real exploit conditions. The final stages included risk assessment and the preparation of an analytical report containing recommendations for mitigating identified issues.

The experimental component of the study was aimed at practically validating the effectiveness of Penetration Testing as a tool for assessing the security level of web-oriented information systems. The research was conducted in a controlled test environment that excluded any impact on production infrastructure and complied with ethical and legal requirements. The tested information system employed a client-server architecture and provided user access to web resources via HTTP/HTTPS protocols. The server component was deployed on a standard web server with TLS support. The system architecture encompassed network, application, and data storage layers, enabling security evaluation across multiple levels.

Professional Penetration Testing tools were used in the experiment: Nmap for network reconnaissance and open port analysis [23], Burp Suite for web application testing and HTTP traffic inspection [24], and Wireshark for in-depth packet analysis [25]. All tools were integrated within the Kali Linux environment, ensuring experimental integrity and reproducibility [26].

The experimental environment was implemented using Python 3.11. The following libraries and modules were used in the developed scripts:

- requests 2.31.0 for HTTP communication and header analysis;
- ssl (Python standard library) for TLS certificate inspection;
- xml.etree.ElementTree for parsing Nmap XML output;
- json 2.0 for structured result serialization;
- python-docx 0.8.11 for automated generation of DOCX security reports.

All experiments were conducted within the Kali Linux 2024 environment.

The practical Penetration Testing methodology was developed in accordance with the concept of sequential multi-layered analysis and includes the following stages (**Fig. 1**):

1. Automated initialization of the testing process.
Implemented through the central script `run_pentest.py`, which coordinates the execution of testing modules and ensures reproducibility.
2. Network reconnaissance and attack surface analysis.
Conducted using Nmap with automated parsing of scan results by the script `parse_nmap_xml.py`.
3. Application-layer and security header analysis.
The script `check_security_headers.py` verifies the presence and correctness of HTTP security headers.
4. Cryptographic protection assessment.
The script `tls_cert_check.py` analyzes TLS certificates and encryption parameters.
5. Results formalization and checklist validation.
Security requirements stored in `checklists.yaml` are used for compliance verification.
6. Automated report generation.
The script `generator_docx.py` produces a structured report in DOCX format.

Figure 1 illustrates the sequential logic of conducting Penetration Testing as a structured process for assessing information security. At the Planning / Scope stage, the

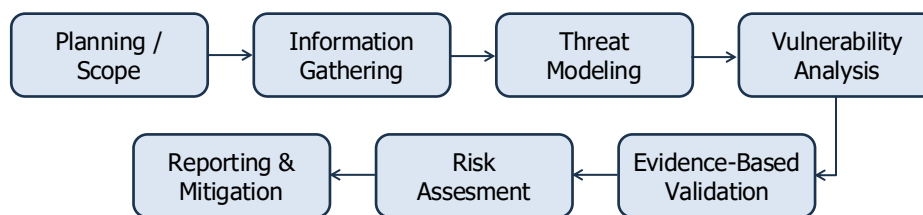


Fig. 1. Generalized model of the Penetration Testing methodology.

testing object, assessment boundaries, and legal considerations are defined. Information Gathering involves the collection of technical data regarding assets, services, and the attack surface. During the Threat Modeling phase, a threat model is constructed, identifying assets, potential actors, and attack scenarios. Vulnerability Analysis focuses on detecting configuration and technical weaknesses at various system levels. Evidence-Based Validation ensures confirmation of vulnerability exploitability without performing destructive actions. This is followed by Risk Assessment, where risks are evaluated using the *Likelihood* × *Impact* model. The final stage, Reporting & Mitigation, includes report generation and the development of recommendations for addressing identified issues.

The study applies the *Asset–Actor–Surface–Threat–Scenario* model, which provides a formalized and easily automated description of the testing process [27]. The identified assets include the web application, user accounts, confidential data, infrastructure, and network services, as well as the TLS communication channel. Actors are represented by an external attacker, a user with limited web interface access, and, where relevant, an internal user. The attack surface comprises the domain or URL, open ports and services, HTTP headers, and TLS certificates. The identified threats include traffic interception or manipulation, clickjacking, MIME sniffing, insecure configurations, outdated or misconfigured TLS settings, and unnecessarily exposed services. Scenarios are implemented as specific test cases that produce measurable results and are incorporated into the final report.

For analytical consistency, the implemented scenarios were aligned with the STRIDE model: Spoofing, Tampering, Repudiation, Information Disclosure, Denial of Service, and Elevation of Privilege [28]. Within the non-invasive approach adopted in this study, the most relevant threat categories are Information Disclosure, Tampering, and Spoofing, as well as, to a limited extent, Denial of Service manifested through configuration weaknesses.

Python was selected for the practical implementation of the research due to its scientific and applied advantages. The choice of Python is justified by its rapid development capabilities, code readability, extensive libraries for network analysis, XML/JSON processing, and HTTP requests, broad automation potential for security-related tasks, integration with Penetration Testing tools [29], and suitability for DevSecOps implementations [30]. Thus, Python serves not only as an implementation tool but also as a universal platform for automating information security testing processes.

RESULTS AND DISCUSSION

Formalization of the Threat Model and the Testing Process

Within the scope of the study, the module `threat_model.py` was analyzed, as it implements a formalized description of assets, attack surface, threats, and testing scenarios. This module constitutes a core component of the test selection system and effectively functions as the central element of the Penetration Testing process. The `SystemModel` class is used to describe the testing object, including the URL or host, as well as known entry points. The `Threat` class formalizes threats according to the STRIDE framework, while the `Scenario` class represents an individual test case whose measurable outcome is used to generate a confirmed finding in the final report.

The project repository includes a set of typical testing modules, namely:

- *check_security_headers.py* – verification of the presence and correctness of HTTP security headers;
- *tls_cert_check.py* – retrieval and analysis of TLS certificates;
- *parse_nmap_xml.py* in conjunction with *nmap_output_sample.xml* – analysis of open ports and services based on XML-formatted scan results.

The module *scenarios.py* encapsulates individual functions into formalized testing scenarios. The core concept is not the unsystematic execution of all checks, but rather the selection of scenarios according to rules derived from the combination of the SystemModel, ThreatModel, and predefined constraints (execution time, testing boundaries, and non-invasive nature).

Within the structured repository, the central script *run_pentest.py* collects the results of HTTP header and TLS certificate checks. A key feature of the implementation is the automated selection of scenarios, the use of a unified output format (JSON), and the straightforward transformation of collected data into confirmed findings for subsequent reporting through the module *generator_docx.py*.

Thus, the study establishes a formalized threat model and implements a mechanism for scenario selection based on it. The Python-based implementation ensures formalization of assets, actors, and threats; automated selection of non-invasive scenarios; reproducible collection of results in a structured format; and a foundation for further risk analysis and reporting.

Process of Practical Penetration Testing

Practical Penetration Testing in the controlled test environment was implemented as a sequential validation pipeline comprising the following stages (Fig. 2):

1. Information Gathering – inventory of the attack surface and parameters of accessible components.
2. Vulnerability Scanning / Baseline Security Checks – identification of typical configuration weaknesses at the network, application, and cryptographic levels.
3. Exploitation – within this study, only safe validation of impact was performed, without destructive actions, with evidence recorded in the form of control artifacts.
4. Post-Exploitation Analysis – formalization of consequences, risk evaluation, preparation of recommendations, and generation of reporting documentation.

The fundamental objective of this approach is to ensure process reproducibility, whereby each stage produces its own artifacts in the form of logs, JSON data, or reports, enabling replication and verification of the research results.

The objective of the information gathering stage is to identify entry points (URL, host, ports, and services) and to obtain primary technical attributes that influence risk, including HTTP security headers, TLS certificate parameters, and information about open ports.

Data Collection Implementation.

The collection of HTTP security headers using the module *check_security_headers.py* enables the assessment of the status of fundamental web controls such as Content-Security-Policy, HSTS, X-Frame-Options, and others. The

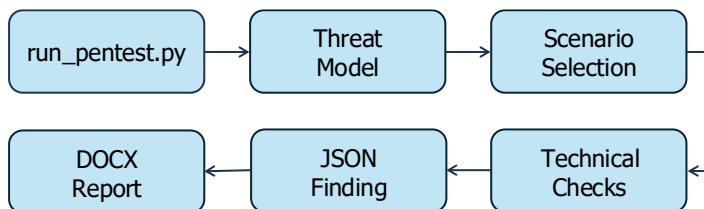


Fig. 2. Automated pipeline for the practical implementation of Penetration Testing.

absence or misconfiguration of these headers is considered an indicator of weak configuration and an insufficient level of system hardening.

The analysis of TLS settings is performed using the module *tls_cert_check.py*, which verifies certificate validity, expiration dates, and trust parameters. Incorrect TLS configurations increase the risk of traffic interception or manipulation and may reduce user trust.

Information on open ports is obtained by analyzing Nmap scan results in XML format using the module *parse_nmap_xml.py*. The presence of unnecessary services, even without active exploitation, is interpreted as an expansion of the attack surface and a potential source of risk.

Approach to Exploitation and Risk Assessment

The vulnerability scanning phase focuses on verifying typical classes of weaknesses, including configuration errors, absence of baseline security controls, and outdated TLS parameters, in accordance with OWASP recommendations and established hardening practices.

In the classical sense, the exploitation stage involves confirming the exploitability of vulnerability. In this study, a safe *evidence-based validation* approach is applied, whereby no destructive actions are performed, and no unauthorized access is attempted. Instead, weaknesses are validated based on configuration evidence and controlled testing artifacts. For example, the absence of the HSTS header or the presence of an exposed administrative port is treated as a confirmed indicator of elevated risk.

The post-exploitation stage is interpreted as an analytical phase that includes impact assessment, risk calculation, prioritization of issues, development of a mitigation plan, and preparation of reporting documentation.

Formalization of Results and Risk Model

Upon completion of Penetration Testing, a set of technical artifacts is generated; however, their practical value is realized only after analytical interpretation. For this purpose, the results obtained from various tools were normalized into a unified finding structure of the following form:

$$Finding = \{ID, Title, Category, Asset, Evidence, Impact, Likelihood, Risk, Recommendation, Priority\}.$$

For each finding, the reliability of confirmation, likelihood of exploitation, impact on confidentiality, integrity, and availability, as well as the exploitation context, are evaluated.

Risk assessment is conducted using the basic model:

$$Risk = Likelihood \times Impact,$$

where the values of *Likelihood* and *Impact* are determined using discrete scales according to the level of exposure and asset criticality [31]. This approach enables quantitative interpretation of results and supports the prioritization of information security improvement measures.

Extended Risk Assessment Model Incorporating Confidence

To enhance the accuracy of quantitative risk evaluation, an extended model was applied that incorporates the level of confidence in the testing results. The model is defined as:

$$Risk_{adj} = (Likelihood \times Impact) \times Confidence,$$

where *Confidence* $\in \{0.5; 0.75; 1.0\}$ reflects the degree of validation of the identified weakness based on the obtained evidence.

Within the study, the *Impact* value was determined according to asset criticality using a five-point scale (1-5), while the *Likelihood* indicator was derived from the severity of the identified issue and its exposure level. This approach allows both the technical characteristics of the weakness and the context of its potential exploitation to be taken into account.

A comprehensive evaluation of Penetration Testing results was performed, including normalization of technical data, formation of a register of confirmed findings, quantitative risk assessment using the *Likelihood* \times *Impact* model with the *Confidence* coefficient, and development of a Mitigation Plan. Additionally, structural requirements for the final report were defined, ensuring its suitability for both technical specialists and managerial personnel.

The proposed approach facilitates the transition from purely technical security metrics to substantiated managerial decisions and aligns with contemporary risk-oriented information security management practices (Fig. 3). This figure illustrates the transformation of technical Penetration Testing results into a format suitable for managerial analysis.

At the initial stage, Technical Artifacts (Headers / TLS / Ports) are generated, representing primary data obtained from network, application, and cryptographic-level assessments. These data are subsequently structured into Normalized Findings, ensuring their standardized and unified representation.

The next step is Risk Scoring, within which quantitative risk evaluation is performed using the *Likelihood* \times *Impact* model (with the optional inclusion of the Confidence coefficient where applicable).

The resulting data are aggregated into an Executive Summary, providing a concise overview of the overall risk level for management.

The final stage involves the development of a Mitigation Plan, which specifies concrete corrective actions, prioritization levels, and responsible stakeholders for reducing the identified risks.

Thus, the diagram demonstrates the logical transition from technical analysis to substantiated managerial decision-making in the field of information security.

The generated file *pentest_report.docx* contains a structured report that includes one confirmed finding, an example of which is presented in Fig. 4.

Structure and Content of the Reporting Documentation

The generated document is entitled "Penetration Testing Security Report," which serves as the formal title and clearly defines the document as a report presenting the results of Penetration Testing.

Scope / Object

Section "1. Scope / Object" specifies the testing object as follows:

- Target URL: *https://example.com*
- Target Host: *example.com*

A clear definition of the object and testing boundaries is essential to ensure the methodological and legal validity of the report, as it enables unambiguous identification of the information system subjected to testing.

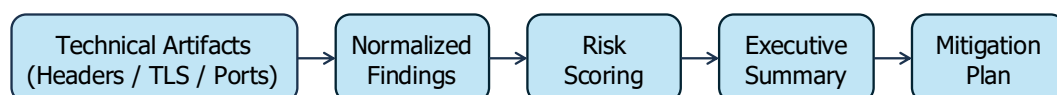


Fig. 3. Logical framework for transitioning from technical results to managerial decision-making.

Penetration Testing Security Report

1.-Scope-/Object
 Target URL: https://example.com
 Target Host: example.com

2.-Executive Summary
 Total findings: 1 (High=1, Medium=0, Low=0)
 Top issues:
 •→ 1) Missing HSTS | Priority: High | Risk: 12

3.-Findings Register

ID	Title	Category	Asset	Severity	Risk	Priority	Evidence (short)
F-001	Missing HSTS	Security Misconfiguration	Web Application	HIGH	12	High	HTTP response does not contain Strict-Transport-Security

4.-Findings Details

4.1.-Missing HSTS
 ID: F-001
 Category: Security Misconfiguration
 Asset: Web Application
 Severity: HIGH
 Risk score: 12

5.-Mitigation Plan

Finding ID	Action	Priority	Owner	Due	Status
F-001	Configure Strict-Transport-Security (HSTS) for the web server	High	DevOps/Security	TBD	Open

Priority: High
 Description: Strict-Transport-Security header is missing
 Evidence: HTTP response does not contain Strict-Transport-Security
 Impact: Increases risk of downgrade SSL-stripping class issues
 Recommendation: Enable HSTS (max-age, includeSubDomains) and verify HTTPS-only policy

Fig. 4. Structured report containing a single confirmed finding.

Executive Summary

Section “2. Executive Summary” provides a concise managerial overview of the testing results in the form of aggregated indicators. Specifically, the report states:

- total number of identified findings - 1;
- severity distribution: High - 1, Medium - 0, Low - 0;
- most critical issue - Missing HSTS with High priority and a risk score of 12.

This format allows for the rapid identification of high-priority issues requiring immediate remediation.

Findings Register

Section “3. Findings Register” presents a tabulated register of identified issues containing the following fields: identifier, title, category, asset, severity level, risk score, priority, and a brief description of evidence. The report documents one finding with the following characteristics:

- ID: F-001;
- Title: Missing HSTS;
- Category: Security Misconfiguration;
- Asset: Web Application;
- Severity: High;
- Risk: 12;
- Priority: High;
- Evidence: absence of the Strict-Transport-Security header in the HTTP server response.

The findings register functions as an inventory of identified issues and serves as a practical instrument for further analysis and decision-making.

Findings Details

Section “4. Findings Details” provides an extended description of the Missing HSTS finding, including a problem statement, supporting evidence, impact assessment, and mitigation recommendations. In particular, the absence of the Strict-Transport-Security

header increases the risk of downgrade and SSL-stripping attacks. The recommended mitigation measure includes enabling HSTS with appropriate parameters (max-age, includeSubDomains) and verifying the enforced HTTPS policy.

This section follows a clear analytical structure – problem – evidence – consequences – recommendations – and represents the most informative component of the report.

Mitigation Plan

Section “5. Mitigation Plan” outlines corrective measures for addressing the identified issue, specifying the finding identifier, recommended action, priority level, responsible parties, implementation timeline, and current status. For finding F-001, the status is defined as Open, with DevOps/Security specialists designated as responsible parties, and the implementation deadline to be determined.

The identified issue Missing HSTS is classified as an application-level configuration weakness. Although it does not result in immediate system compromise, it significantly increases the risk of HTTPS-related attacks under certain conditions. In the report, this weakness is categorized as Security Misconfiguration, assigned a High severity level and a risk score of 12, thereby justifying its high remediation priority.

Summary of Results and Automation of Reporting

As a result of the Penetration Testing conducted in the controlled environment, a formalized security report was generated, integrating validated technical findings with their analytical interpretation. The analytical processing of results enabled the transition from purely technical observations to quantitative risk evaluation using the Likelihood × Impact model and supported the prioritization of security measures.

The study also implemented automated addition of multiple findings to the DOCX report. For this purpose, the `generate_report()` function was extended through the module `auto_findings.py`, which automatically generates findings based on the results of security header checks, TLS configuration analysis, and open port assessments, and transfers them to the report generation module.

As a result of the practical implementation of the Penetration Testing process, a DOCX report was automatically generated containing a structured set of technical and analytical results assessing the security level of the web application. The report was developed in accordance with a risk-oriented approach and incorporates all key components necessary for further information security management (Fig. 5).

In the Scope / Object section, the testing object – the web resource `https://example.com` and the corresponding host `example.com` - is clearly identified. Such identification ensures the formal correctness of the report, defines the responsibilities of the involved parties, and confirms the legitimacy of conducting Penetration Testing within the defined test environment.

The Findings Register section presents a tabulated register of all identified weaknesses, including the finding identifier, title, category, asset, severity level, risk score, and priority. All identified issues were classified as Security Misconfiguration and attributed to the application layer of the web application. This register functions as a centralized risk inventory tool and allows the report to serve not only as a technical document but also as a foundation for managerial analysis.

The Findings Details section provides an extended description of each identified finding using a unified analytical structure that includes a problem statement, supporting evidence, impact assessment, and remediation recommendations. In particular, the absence of HSTS is assessed as a high-risk issue due to the potential for SSL-stripping attacks; the absence of Content-Security-Policy reduces the application's resilience to cross-site scripting and other client-side attacks; the lack of X-Frame-Options creates conditions for clickjacking attacks; and the absence of X-Content-Type-Options may lead to MIME-sniffing and incorrect content interpretation by browsers. For each weakness,

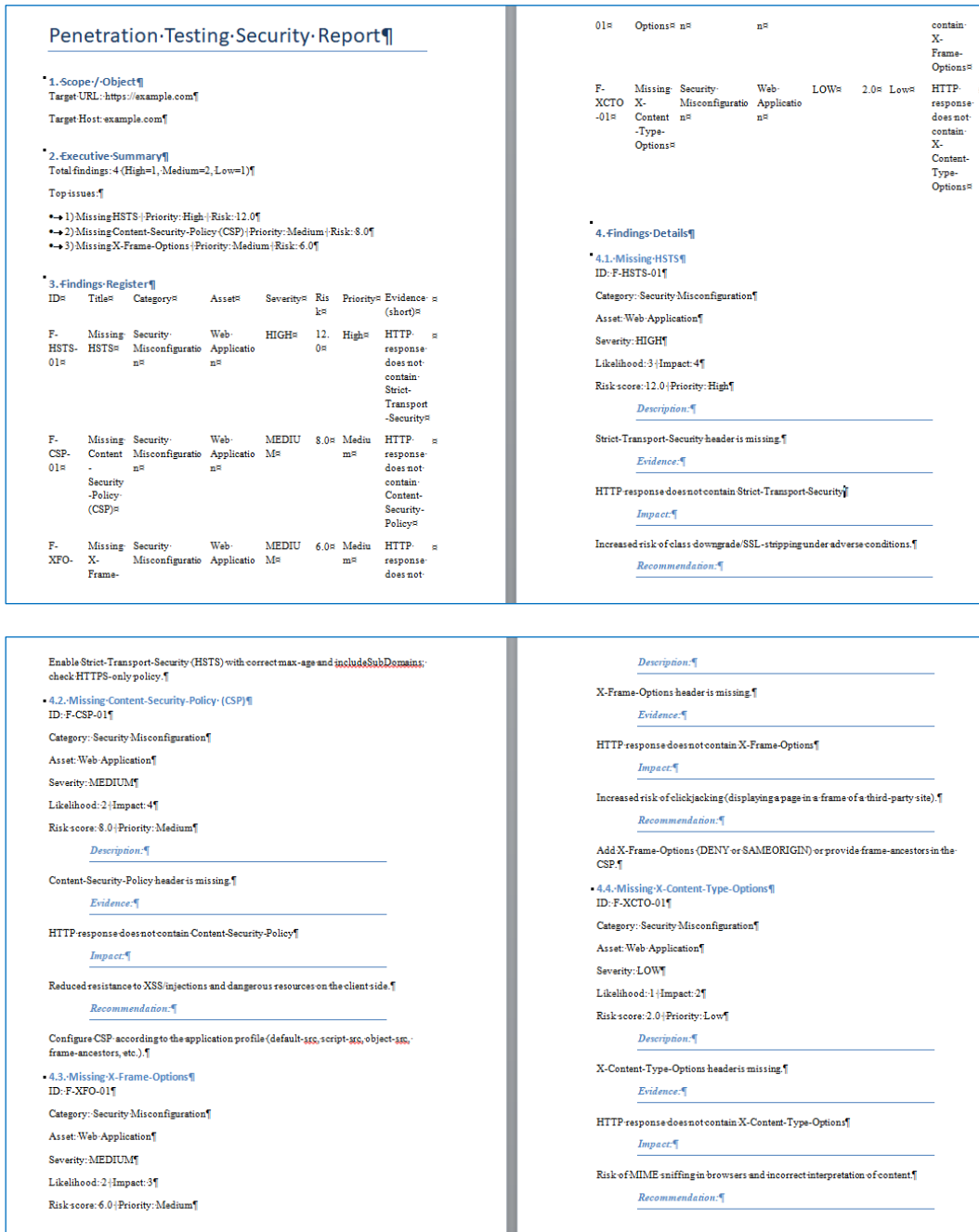


Fig. 5 (beginning). Complete structured security report.

specific technical recommendations are provided in accordance with contemporary web application security best practices.

The final element of the report is the Mitigation Plan, which specifies, for each finding, the recommended actions, priority level, responsible parties, and current implementation status. The presence of such a plan demonstrates the transition from technical identification of issues to the practical implementation of measures aimed at improving the overall level of information security.

```

Add-X-Content-Type-Options: nosniff

* 5.-Mitigation-Plan¶
Finding ID¶ Action¶ Priority¶ Owner¶ Due¶ Status¶ ¶
F-HSTS-01¶ Enable Strict-Transport-Security (HSTS) with correct max-age and includeSubDomains; check HTTPS-only policy.¶ High¶ DevOps/Security¶ TBD¶ Open¶ ¶
F-CSP-01¶ Configure CSP according to the application profile (default-src, script-src, object-src, frame-ancestors, etc.).¶ Medium¶ DevOps/Security¶ TBD¶ Open¶ ¶
F-XFO-01¶ AddX-Frame-Options (DENY or SAMEORIGIN) or provide frame-ancestors in the CSP.¶ Medium¶ DevOps/Security¶ TBD¶ Open¶ ¶
F-XCTO-01¶ AddX-Content-Type-Options: nosniff.¶ Low¶ DevOps/Security¶ TBD¶ Open¶ ¶
¶
¶

```

Fig. 5 (ending). Complete structured security report.

The generated file `pentest_report.docx` constitutes a comprehensive analytical document that integrates technical Penetration Testing results with quantitative risk assessment and managerial recommendations. This confirms the effectiveness of the applied approach and highlights the practical value of Penetration Testing as a tool for assessing and enhancing the information security of modern information systems.

Figure 6 presents a concise Penetration Testing report in Markdown format, reflecting the structure of the collected data, their interpretation, and the logical conclusions derived from the analysis. The file `report.md` consists of three main components: testing metadata (target, host, timestamp), technical results (security header and TLS certificate verification), and generalized conclusions.

```

# Pentest report
- ** Destination URL:** `N/A`
- ** Host / IP:** `N/A`
- ** Scan time:** `2026-01-23T23:58:02`
## Other data
### headers
- **status:**
  200
- **headers:**
  - **Content-Security-Policy:**
    None
  - **Strict-Transport-Security:**
    None
  - **X-Frame-Options:**
    None
  - **X-Content-Type-Options:**
    None
  - **Referrer-Policy:**
    None
### cert
- **subject:**
  - item #1:
    - item #1:
      - countryName

```

Fig. 6 (beginning). Penetration Testing report in Markdown format.

```

    - US
- item #2:
  - item #1:
    - stateOrProvinceName
    - California
- item #3:
  - item #1:
    - localityName
    - Los Angeles
- item #4:
  - item #1:
    - organizationName
    - Internet Corporation for Assigned Names and Numbers
- item #5:
  - item #1:
    - commonName
    - *.example.com
- **issuer**:
  - item #1:
    - item #1:
      - countryName
      - US
    - item #2:
      - item #1:
        - organizationName
        - DigiCert Inc
    - item #3:
      - item #1:
        - commonName
        - DigiCert Global G3 TLS ECC SHA384 2020 CA1
- **version**:
  3
- **serialNumber**:
  0AD893BAFA68B0B7FB7A404F06ECA9A
- **notBefore**:
  Jan 15 00:00:00 2025 GMT
- **notAfter**:
  Jan 15 23:59:59 2026 GMT
- **subjectAltName**:
  - item #1:
    - DNS
    - *.example.com
  - item #2:
    - DNS
    - example.com
- **OCSP**:
  - http://ocsp.digicert.com
- **caIssuers**:
  - http://cacerts.digicert.com/DigiCertGlobalG3TLSECCSHA3842020CA1-2.crt
- **crlDistributionPoints**:
  - http://crl3.digicert.com/DigiCertGlobalG3TLSECCSHA3842020CA1-2.crl
  - http://crl4.digicert.com/DigiCertGlobalG3TLSECCSHA3842020CA1-2.crl

```

Conclusions

Based on the obtained results, it is recommended to further analyze open ports, TLS/SSL configurations, and potential vulnerabilities (if identified), as well as to strengthen the web server's security policies.

Fig. 6 (ending). Penetration Testing report in Markdown format.

The analysis of metadata indicates that the testing was conducted at a specified point in time using non-invasive methods for assessing web server configuration and the cryptographic parameters of the secured connection. The absence of specific URL and host values in the report is interpreted as a result of test environment anonymization or the omission of parameters during script execution, which does not affect the overall validity of the analysis.

The results of the HTTP security header assessment reveal the absence of fundamental application-layer protection mechanisms, including Strict-Transport-Security, Content-Security-Policy, X-Frame-Options, X-Content-Type-Options, and Referrer-Policy. The absence of these headers is classified as a configuration weakness within the category of Security Misconfiguration. Although it does not result in immediate system compromise, it significantly increases the overall level of information security risk.

Particular attention should be given to the TLS certificate analysis results. It was determined that the certificate, issued by the trusted certification authority DigiCert, had expired at the time of testing. An expired TLS certificate represents a critical operational risk, as it may lead to loss of client trust, potential browser access blocking, and reduced service availability. This issue directly affects the Availability component of the CIA model and indicates insufficient lifecycle management of cryptographic certificates.

CONCLUSION

The overall evaluation of the results leads to the conclusion that the identified issues are predominantly configuration-related and operational in nature and are not accompanied by active exploitation of vulnerabilities. At the same time, such weaknesses often constitute preconditions for more sophisticated attacks under real-world operating conditions. The application of a quantitative risk assessment model based on the combination of likelihood and impact made it possible to classify the expired TLS certificate and the absence of HSTS as priority issues requiring immediate remediation, whereas the remaining missing security headers may be addressed through planned security enhancement measures.

Thus, the results of the Penetration Testing confirm the effectiveness of this approach as a comprehensive instrument for information security assessment. The obtained data not only reflect the technical condition of the system but also provide a substantiated basis for managerial decision-making aimed at risk reduction and improvement of the overall protection level of information resources.

The source code is available in the GitHub repository at the following link: [incom2025/Secur_infor_syst](https://github.com/incom2025/Secur_infor_syst).

ACKNOWLEDGMENTS AND FUNDING SOURCES

The authors received no financial support for the research, writing, and/or publication of this article.

COMPLIANCE WITH ETHICAL STANDARDS

The authors declare that they have no competing interests.

AUTHOR CONTRIBUTIONS

Conceptualization, [S. S.]; formal analysis, [S. S., I. Ka., Y. S.]; investigation, [I. Ku.]; resources, [I. Ku.]; data curation, [M. D., N. S., L. P., V. B.], writing – original draft preparation, [S. S., I. Ka, I. Ku.]; writing – review and editing, [O. K., Y. S., M. D., N. S., L. P.], visualization, [O. K., Y. S., V. B.].






All authors have read and agreed to the published version of the manuscript.

REFERENCES

- [1] Anderson R. Security Engineering: A Guide to Building Dependable Distributed Systems. 3rd ed. Wiley, 2020. URL: <https://www.wiley.com/en-us/Security+Engineering%3A+A+Guide+to+Building+Dependable+Distributed+Systems%2C+3rd+Edition-p-9781119642787>.
- [2] Behl A., Behl K. Cyberwar: The Next Threat to National Security and What to Do About It. Oxford University Press, 2017.
- [3] Bishop M. *Computer Security: Art and Science*. 2nd ed. Addison-Wesley, 2019. URL: https://ptgmedia.pearsoncmg.com/images/9780321712332/samplepages/9780321712332_Sample.pdf
- [4] OWASP Foundation. OWASP Web Security Testing Guide v4. 2023. URL: <https://owasp.org/www-project-web-security-testing-guide/>.
- [5] European Union Agency for Cybersecurity (ENISA). *Good Practices for Security of IoT – Secure Development and Testing*. 2020. URL: <https://www.enisa.europa.eu/publications/good-practices-for-security-of-iot>
- [6] ISO/IEC 27001:2022. Information Security Management Systems - Requirements. ISO, 2022. URL: <https://www.iso.org/standard/82875.html>.
- [7] ISO/IEC 27002:2022. *Information Security Controls*. ISO, 2022. URL: <https://www.iso.org/standard/75652.html>.
- [8] Kizza J. M. *Guide to Computer Network Security*. 5th ed. Springer, 2020. <https://doi.org/10.1007/978-3-030-38141-7>.
- [9] MITRE. MITRE ATT&CK®: Adversarial Tactics, Techniques, and Common Knowledge. 2023. URL: <https://attack.mitre.org/>
- [10] CVE Program. Common Vulnerabilities and Exposures (CVE). MITRE Corporation, 2024. URL: <https://www.cve.org/>.
- [11] NIST SP 800-53 Rev. 5. *Security and Privacy Controls for Information Systems and Organizations*. National Institute of Standards and Technology, Gaithersburg, 2020. <https://doi.org/10.6028/NIST.SP.800-53r5>.
- [12] NIST SP 800-115. *Technical Guide to Information Security Testing and Assessment*. National Institute of Standards and Technology, Gaithersburg, 2008. URL: <https://doi.org/10.6028/NIST.SP.800-115>.
- [13] OWASP Foundation. *OWASP Testing Guide v4*. 2023. URL: <https://owasp.org/www-project-web-security-testing-guide/>
- [14] OWASP Foundation. OWASP Top 10 - Web Application Security Risks. 2021. URL: <https://owasp.org/www-project-top-ten/>.
- [15] Scarfone, K., Mell, P. *Guide to Intrusion Detection and Prevention Systems (IDPS)*, (NIST SP 800-94), National Institute of Standards and Technology, Gaithersburg, 2007. <https://doi.org/10.6028/NIST.SP.800-94>.
- [16] Nelson A., Rekhi S., Souppaya M., Scarfone K. Incident Response Recommendations and Considerations for Cybersecurity Risk Management. NIST SP 800-61 Rev.3. National Institute of Standards and Technology, 2025. <https://doi.org/10.6028/NIST.SP.800-61r3>.
- [17] Stallings W. *Network Security Essentials: Applications and Standards*. 6th ed. Pearson, 2017. URL: <https://www.pearson.com/en-us/subject-catalog/p/network-security-essentials-applications-and-standards/P200000003333>

- [18] Ross, R. (2012), *Guide for Conducting Risk Assessments*, (NIST SP 800-30 Rev. 1), National Institute of Standards and Technology, Gaithersburg, MD.
<https://doi.org/10.6028/NIST.SP.800-30r1>.
- [19] ENISA. Threat Landscape Report 2023. European Union Agency for Cybersecurity. URL: <https://www.enisa.europa.eu/publications/enisa-threat-landscape-2023>.
- [20] Whittaker J. A., Arbon J., Carollo J. *How Google Tests Software*. Addison-Wesley, 2012. URL: <https://www.informit.com/store/how-google-tests-software-9780321803023>
- [21] OSSTMM Institute. Open Source Security Testing Methodology Manual (OSSTMM) v3. 2019. URL: <https://www.isecom.org/OSSTMM.3.pdf>.
- [22] PTES. *Penetration Testing Execution Standard (PTES)*. 2020. URL: <http://www.pentest-standard.org/>.
- [23] Shostack A. *Threat Modeling: Designing for Security*. Wiley, 2014. URL: <https://www.wiley.com/en-us/Threat+Modeling%3A+Designing+for+Security-p-9781118809990>.
- [24] Singer, P. W., & Friedman, A. (2013). *Cybersecurity and Cyberwar: What Everyone Needs to Know*. Oxford University Press. URL: <https://global.oup.com/academic/product/cybersecurity-and-cyberwar-9780199918096?q=Cybersecurity%20and%20Cyberwar:%20What%20Everyone%20Needs%20to%20Know&cc=ua&lang=en>.
- [25] Humble J., Farley D. *Continuous Delivery: Reliable Software Releases through Build, Test, and Deployment Automation*. Addison-Wesley, 2010. URL: <https://www.informit.com/store/continuous-delivery-reliable-software-releases-through-9780321601919>.
- [26] Kim G., Behr K., Spafford G. *The Phoenix Project: A Novel About IT, DevOps, and Helping Your Business Win*. IT Revolution Press, 2018.
<https://itrevolution.com/product/the-phoenix-project/>.
- [27] Nelson, A., Rekhi, S., Souppaya, M. and Scarfone, K. *Incident Response Recommendations and Considerations for Cybersecurity Risk Management: A CSF 2.0 Community Profile*. NIST SP 800-61 Rev. 3. National Institute of Standards and Technology, Gaithersburg, 2025.
<https://doi.org/10.6028/NIST.SP.800-61r3>.
- [28] CVE Program. Common Vulnerabilities and Exposures (CVE). MITRE Corporation, 2024. URL: <https://www.cve.org/>.
- [29] Dempsey, K., Johnson, L., Scholl, M., Stine, K., Clay, A., Orebaugh, A., Chawla, N. and Johnston, R. (2011), *Information Security Continuous Monitoring (ISCM) for Federal Information Systems and Organizations*, Special Publication (NIST SP), National Institute of Standards and Technology, Gaithersburg, MD.
<https://doi.org/10.6028/NIST.SP.800-137> .
- [30] Zalewski M. *The Tangled Web: A Guide to Securing Modern Web Applications*. No Starch Press, 2012. <https://nostarch.com/tangledweb>.
- [31] MITRE. *MITRE ATT&CK®: Adversarial Tactics, Techniques, and Common Knowledge*. 2023. URL: <https://attack.mitre.org/>.

ЗАСТОСУВАННЯ *PENETRATION TESTING* ДЛЯ ОЦІНЮВАННЯ РІВНЯ ІНФОРМАЦІЙНОЇ БЕЗПЕКИ ВЕБ-ОРІЄНТОВАНИХ ІНФОРМАЦІЙНИХ СИСТЕМ

Сергій Свелєба¹ , Іван Катеринчук¹ , Іван Куньо¹ ,
Олег Крупич¹, Ярослав Шмигельський¹ , Марта Дуфанець¹ ,
Наталя Свелєба², Люціан Пельц³, Володимир Бригілевиц³

¹Львівський національний університет імені Івана Франка
вул. Ген. Тарнавського, 107, 79017 Львів, Україна

²Приватний вищий навчальний заклад «Європейський університет»,
бульвар Академіка Вернадського, 16 В, 03115 Київ, Україна

³Державна вища школа технологій та економіки в Ярославі
вул. Чарнецького, 16, 37-500 Ярослав, Польща

АНОТАЦІЯ

Вступ. Зростання ролі веб-орієнтованих інформаційних систем у бізнесі, освіті та державному управлінні супроводжується підвищенням кількості та складності кіберзагроз. Традиційні засоби захисту не завжди дозволяють виявити реальні слабкі місця систем, що обумовлює необхідність застосування практично орієнтованих методів оцінювання рівня інформаційної безпеки. У цьому контексті *Penetration Testing* розглядається як ефективний інструмент імітації дій потенційного злоумисника з метою виявлення та підтвердження експлуатованих вразливостей.

Матеріали та методи. У роботі використано ризик-орієнтований підхід відповідно до міжнародних стандартів ISO/IEC 27001, ISO/IEC 27005, рекомендацій OWASP та NIST SP 800-115. *Penetration Testing* реалізовано як поетапний процес, що включає збір інформації, аналіз поверхні атаки, формування моделі загроз, виконання неінвазивних сценаріїв перевірки та оцінку ризиків. Практичну частину виконано у тестовому середовищі з використанням інструментів Nmap, Burp Suite та Wireshark, а також власних Python-модулів для автоматизації аналізу HTTP-заголовків безпеки, TLS-сертифікатів і відкритих сервісів.

Результати. У ході дослідження виявлено низку конфігураційних слабкостей прикладного рівня, зокрема відсутність базових захисних HTTP-заголовків та недоліки в управлінні TLS-сертифікатами. Отримані результати формалізовано у вигляді реєстру знахідок із кількісною оцінкою ризиків за моделлю *Likelihood* × *Impact*. Аналіз показав, що навіть за відсутності критичних експлуатаційних вразливостей конфігураційні помилки істотно підвищують загальний рівень ризику та можуть створювати передумови для складніших атак.

Висновки. Результати дослідження підтверджують ефективність *Penetration Testing* як інструменту комплексної оцінки інформаційної безпеки веб-орієнтованих систем. Запропонований підхід забезпечує перехід від технічних результатів тестування до обґрунтованих управлінських рішень, спрямованих на зниження ризиків та підвищення рівня захищеності інформаційних ресурсів.

Ключові слова: інформаційна безпека, penetration testing, вразливості, оцінка ризиків, веб-інформаційні системи.

UDC: 004.852:796.06

LONG SHORT-TERM MEMORY RECURRENT NEURAL NETWORK FOR STATE PREDICTION AND RESOURCE ALLOCATION OPTIMIZATION IN DISTRIBUTED SYSTEMS

Zinovii Liubun^{ORCID}, Oleh Tereshchuk^{ORCID}*

Department of RadioPhysics and Computer Technologies,
Ivan Franko National University of Lviv
107 Gen. Tarnavsky Str., 79017 Lviv, Ukraine

Liubun, Z., & Tereshchuk, O. (2026). Long Short-Term Memory Recurrent Neural Network for State Prediction and Resource Allocation Optimization in Distributed Systems. *Electronics and Information Technologies*, 33, 87–96. <https://doi.org/10.30970/eli.33.7>

ABSTRACT

Introduction. This paper considers a method based on a Long Short-Term Memory (LSTM) neural network for optimal resource allocation in distributed systems. The developed algorithm ensures high accuracy in predicting resource states and optimal spatial distribution with minimal processing time. The relevance of this research is determined by the growing need for intelligent automation of resource management processes in service infrastructure facilities. A locker management system in sports facilities is used as a practical demonstration of the method's effectiveness.

Materials and Methods. To address the prediction and optimization tasks, an LSTM-based architecture with 32 hidden neurons and a sequence length of 10 time steps is proposed. The LSTM model processes sequential occupancy data to capture temporal dependencies and generate probability estimates for future resource states. A multi-factor scoring function is developed to transform predictions into optimal allocation decisions, considering spatial constraints and user preferences. The method is systematically compared with classical approaches: heuristic algorithms (Sequential, Round-Robin), statistical time series models (ARIMA, exponential smoothing), and machine learning methods (logistic regression, random forest, gradient boosting). All methods are evaluated on identical datasets using consistent metrics, including prediction accuracy, F1-score, spatial balance index, and zone variance.

Results. Using LSTM neural networks for the prediction task achieves 85% accuracy, which is statistically significantly higher than Random Forest (79%, $p=0.0023$) and ARIMA (68%, $p=0.0001$). The spatial balance index improved by 8.5% compared to the best classical method (0.89 versus 0.82). Inference time remains acceptable for real-time applications (18.9 ms per prediction).

Conclusions. The proposed LSTM-based method demonstrates satisfactory accuracy in predicting resource states and optimizing their allocation within minimal timeframes. The ability to model long-term temporal dependencies provides significant advantages over classical fixed-window methods. Therefore, the method can be effectively applied to enhance the functionality of distributed resource management systems.

Keywords: LSTM neural network, resource allocation, time series prediction, recurrent neural networks, optimization.



© 2026 Zinovii Liubun & Oleh Tereshchuk. Published by the Ivan Franko National University of Lviv on behalf of Електроніка та інформаційні технології / Electronics and Information Technologies. This is an Open Access article distributed under the terms of the [Creative Commons Attribution 4.0 License](https://creativecommons.org/licenses/by/4.0/) which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

INTRODUCTION

Optimal object placement is an important task in many application domains — from warehouse management to parking space organization and storage systems [1-3]. Traditional approaches are based on simple heuristic rules or statistical models that do not adequately account for temporal change characteristics.

Long Short-Term Memory (LSTM) recurrent neural networks, proposed by Hochreiter and Schmidhuber in 1997 [4], demonstrate high efficiency in tasks involving sequential data. Unlike classical RNNs, the LSTM architecture solves the vanishing gradient problem through specialized gate mechanisms (forget gate, input gate, output gate), enabling effective modeling of long-term dependencies [5-7]. Development of LSTM-based optimal placement algorithms combined with modern computational capabilities will enable the implementation of efficient and reliable object management systems.

In recent years, attention-based architectures, particularly Transformers [9], originally developed for natural language processing, have been increasingly adapted for time series forecasting tasks [8]. The Temporal Fusion Transformer (TFT) [10] combines recurrent layers with multi-head attention mechanisms, enabling interpretable multi-horizon forecasting with explicit modeling of static covariates and known future inputs. The Informer architecture [11] introduces ProbSparse self-attention and a generative-style decoder to achieve $O(L \log L)$ complexity for long-sequence time series forecasting, addressing the quadratic complexity limitation of standard Transformers. A hybrid approach combining exponential smoothing with recurrent neural networks [12] won the prestigious M4 forecasting competition, demonstrating that integration of classical statistical methods with deep learning can yield superior results. The DeepAR framework [13] employs autoregressive recurrent networks for probabilistic forecasting, producing calibrated prediction intervals that are valuable for resource planning under uncertainty. A comprehensive experimental review of deep learning architectures for time series forecasting [14] confirms that recurrent models remain competitive across diverse benchmarks despite the emergence of newer architectures. Furthermore, Hewamalage et al. [15] provide a thorough analysis of the current status of recurrent neural networks for time series forecasting, demonstrating that well-configured LSTM and GRU models can match or exceed the performance of more complex architectures, reaffirming the relevance of investigating recurrent approaches for specific application domains where sequential dependency modeling is paramount.

This work investigates the application of LSTM for optimizing locker placement in sports facilities as a model example of a task involving discrete objects, temporal dynamics, and spatial distribution requirements.

MATERIALS AND METHODS

It is necessary to develop an algorithm for solving the optimal object placement problem:

1. Predicting future occupancy states of objects based on historical usage patterns.
2. Selecting the optimal object for placement considering the forecast and spatial load distribution.

Problem statement:

Let there be a set of lockers $L = \{l_1, l_2, \dots, l_n\}$, each of which at time t can be in state:

$$s_i(t) \in \{0,1\}, \quad (1)$$

where 0 – free, 1 – occupied.

The task is to predict the state of each locker at the next time:

$$\hat{s}_i(t + 1) = f[s_i(t), s_i(t - 1), \dots, s_i(t - k)], \quad (2)$$

where f is a function implemented by an LSTM neural network.

Based on the forecast, a recommendation is formed for selecting a locker with minimal occupancy probability while ensuring uniform spatial utilization of the changing room.

Fig. 1 presents the spatial structure of the object – a locker system in a sports facility.

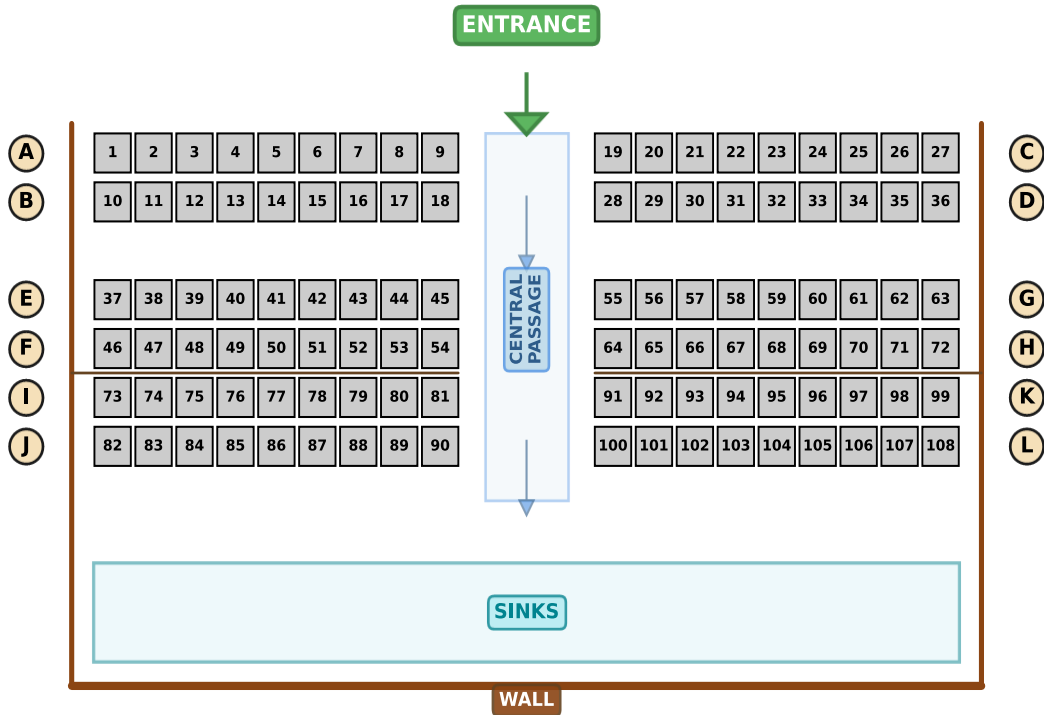


Fig. 1. Sports facility locker system.

Formally, the optimal choice is defined as:

$$l^* = \arg \min_{l_i \in L_{free}} (\alpha P_i + \beta D_i + \gamma U_i), \quad (3)$$

where: P_i – predicted probability of cell occupancy in the near future;

D_i – distance to “dense” zones (for uniformity);

U_i – convenience factor of location (upper/lower, proximity to entrance);

α, β, γ – weight coefficients of importance.

The weight coefficients were determined empirically through a grid search over the ranges $\alpha \in [0.3, 0.7]$, $\beta \in [0.1, 0.4]$, $\gamma \in [0.1, 0.3]$ with a step of 0.05, optimizing the spatial balance index on the validation set. The resulting values $\alpha = 0.5$, $\beta = 0.3$, $\gamma = 0.2$ reflect the priority of prediction accuracy in allocation decisions while maintaining meaningful contributions from spatial uniformity and user convenience factors.

Placement System Requirements: the selection of an optimal object must be performed quickly and ensure uniform spatial distribution. Therefore, simple heuristic methods cannot be effectively used due to their limitations:

- do not account for historical usage data
- do not adapt to changes in behavior patterns
- do not ensure optimal spatial balance.

Thus, a method is needed that combines prediction of temporal dependencies with placement optimization. Neither simple rules nor static models are fully applicable. Recurrent neural networks naturally model sequences and can identify complex patterns.

Figure 2 shows the internal architecture of an LSTM cell with gate mechanisms. The model for the placement task consists of:

1. Input layer: sequence of $k = 10$ binary values of object states;
2. LSTM layer: $h = 32$ hidden neurons with forget gate f_t , input gate i_t , output gate o_t ;
3. Output Dense layer: sigmoid activation for predicting occupancy probability.

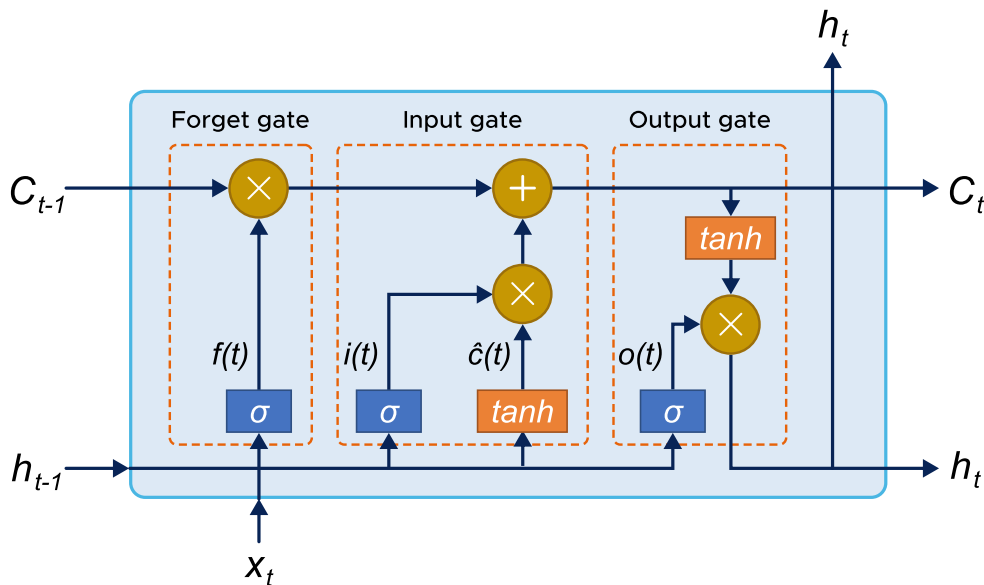


Fig. 2. Internal architecture of LSTM cell.

LSTM Cell Gate Mechanisms:

- $f_t = \sigma(Wf \cdot [h_t^{-1}, x_t] + bf)$ – determines what information to forget from cell state;
- $i_t = \sigma(Wi \cdot [h_t^{-1}, x_t] + bi)$ – determines what new information to store;
- $\tilde{C}_t = \tanh(WC \cdot [h_t^{-1}, x_t] + bC)$ – candidate for new cell state;
- $C_t = f_t \odot C_{t-1} + i_t \odot \tilde{C}_t$ – updated cell state (long-term memory);
- $o_t = \sigma(Wo \cdot [h_t^{-1}, x_t] + bo)$ – determines what to output;
- $h_t = o_t \odot \tanh(C_t)$ – hidden state output.

Training Parameters: Adam optimizer (learning rate = 0.001, $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-7}$), binary cross-entropy loss function, 50 training epochs with early stopping (patience = 7 epochs, monitoring validation loss), batch size 32. The learning rate was selected from

the set {0.0001, 0.0005, 0.001, 0.005} based on validation performance. Dropout regularization with a rate of 0.2 was applied after the LSTM layer to prevent overfitting. Model weights were initialized using Glorot uniform initialization. The total number of trainable parameters is 4,609 (LSTM layer: 4,352 parameters comprising four gate weight matrices of size [32+1, 32] each; Dense layer: 257 parameters). Training convergence was typically achieved within 35–40 epochs, with the best model selected based on minimal validation loss.

To demonstrate the effectiveness of the proposed LSTM approach, a comparative analysis was conducted with eight alternative methods from four classes: heuristic algorithms (Sequential, Round-Robin), statistical time series models (ARIMA, exponential smoothing), classical machine learning (logistic regression, random forest, gradient boosting), and recurrent neural networks (LSTM). Heuristic methods do not use forecasting and have complexity $O(n)$. Statistical models require stationarity assumptions about the data. Machine learning methods with a fixed window $k = 10$ lose information about the sequential nature of data. LSTM naturally models long-term dependencies with complexity $O(n \cdot h^2 \cdot t)$.

Method quality was evaluated across three metric categories: prediction accuracy, placement quality (spatial balance index $\in [0,1]$, zone load variance), and computational efficiency (training time, inference time, memory consumption).

Experiments were conducted on simulated data modeling the operation of a locker management system in a sports facility: $n=108$ cells (distributed across 12 zones of 9 cells each), $T = 500$ time steps, where each step corresponds to 1 hour of real time (approximately 21 days or 3 weeks of operation in total). Average system load was maintained at 60% (65 occupied cells), visitor distribution: 60% adults, 40% children. Typical visit duration: 1.0–1.5 hours. The simulation incorporated realistic usage patterns, including peak hours (morning 8:00–10:00 and evening 17:00–20:00 with load up to 85%), off-peak periods (midday 12:00–15:00 with load approximately 35%), and weekend variations with more uniform distribution. Each simulated visitor was assigned a random arrival time following a bimodal distribution reflecting morning and evening peaks, a visit duration sampled from a log-normal distribution ($\mu = 1.2$ hours, $\sigma = 0.3$), and a zone preference weighted by proximity to the entrance. Data was split into training and test samples at a ratio of 70%/30%. Data preprocessing involved constructing sliding windows of length $k = 10$ from binary occupancy vectors, yielding approximately 31,320 training sequences (290 windows \times 108 cells).

All computations were performed on a single hardware configuration: Apple M1 processor, 16 GB RAM, Python 3.12.2 programming environment with TensorFlow 2.x library (CPU mode). The reported time characteristics (training and inference) are specific to this platform and may differ on systems with different processor architectures or memory capacities.

RESULTS AND DISCUSSION

Table 1 presents the results of comparing eight optimal placement methods.

Figure 3 demonstrates a comprehensive comparison of methods across six performance metrics. The proposed LSTM-based method achieves the highest prediction accuracy (85%) with a statistically significant advantage over the best classical method Random Forest (+6.0%, $p = 0.0023$ at significance level $\alpha = 0.05$), while compared to the best statistical method ARIMA, the improvement is 25% ($p = 0.0001$). In addition to high prediction accuracy, the LSTM method provides the most uniform spatial distribution with minimal zone load variance (67.2 versus 98.4 for Random Forest), confirming that accurate predictions improve the quality of placement decisions. The spatial balance index improved by 8.5%, reaching 0.89 versus 0.82 for classical methods.

Table 1. Comparative characteristics of optimal object placement methods

Method	Class	Prediction Accuracy	Balance Index	Complexity	Time (ms)
Sequential	Heuristic	0.42	0.42	$O(n)$	0.1
Round-Robin	Heuristic	0.68	0.68	$O(n)$	0.2
ARIMA(2,1,2)	Statistical	0.68	0.71	$O(n \cdot p \cdot q)$	15.4
Exponential Smoothing	Statistical	0.64	0.66	$O(n)$	8.2
Logistic Regression	Classical ML	0.74	0.76	$O(n \cdot k \cdot d)$	8.7
Random Forest	Classical ML	0.79	0.82	$O(n \cdot k \cdot d \cdot \log d \cdot T)$	12.3
Gradient Boosting	Classical ML	0.77	0.79	$O(n \cdot k \cdot d \cdot T)$	14.1
LSTM	RNN	0.85	0.89	$O(n \cdot h^2 \cdot t)$	18.9

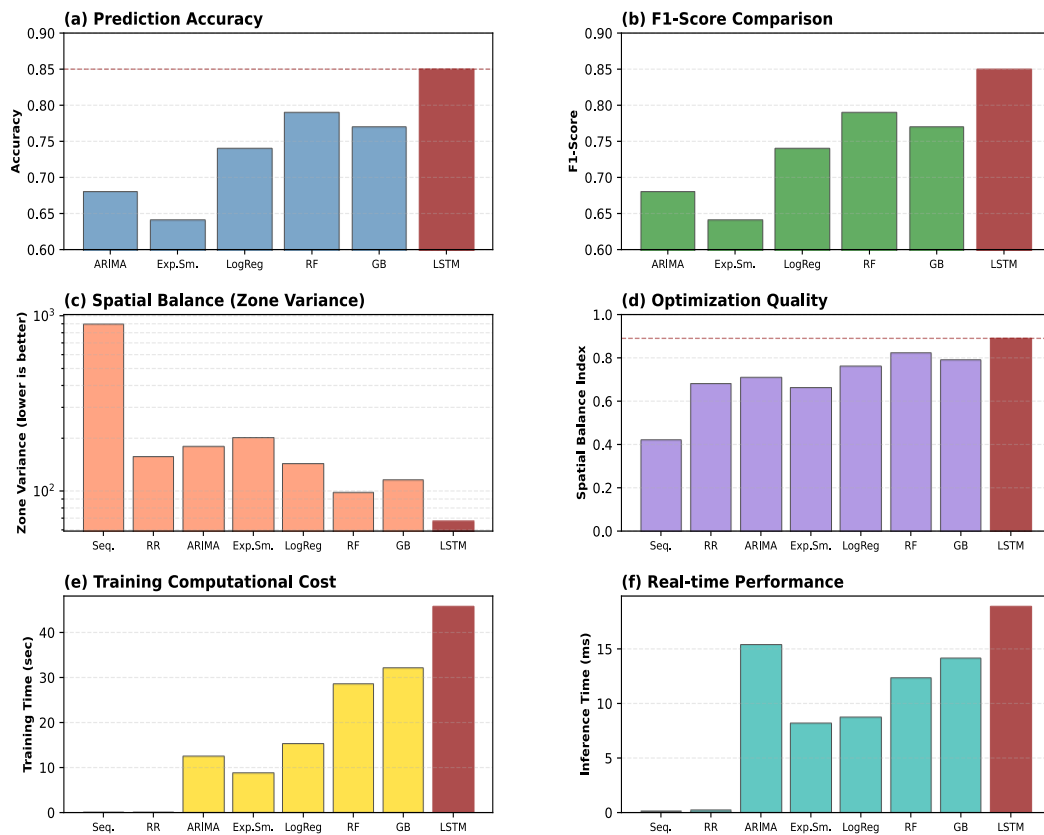


Fig. 3. Comparative analysis of methods.

The statistical significance of the improvements obtained is confirmed by 95% confidence intervals: compared to Random Forest, the accuracy improvement is [+0.03, +0.09], and compared to ARIMA – [+0.13, +0.21]. Paired t-test showed p-values less than 0.05, indicating the reliability of the results. LSTM's computational costs are moderately higher: training time is 45.8 sec versus 28.6 for Random Forest, memory usage is 52.3 MB versus 45.2 MB. However, the placement decision time remains acceptable for real-time systems (18.9 ms), and the additional training costs are compensated by substantial improvement in placement quality during operation.

The key advantages of the LSTM approach lie in its ability to model long-term temporal dependencies through gate mechanisms, which allow effective capture of cyclical object usage patterns that are difficult to model with fixed observation window methods. Unlike ARIMA, which requires explicit differencing and stationarity assumptions, LSTM naturally adapts to pattern changes without additional preprocessing. The combination of gates and nonlinear activation functions enables the detection of complex nonlinear dependencies between object states that are inaccessible to linear methods. Among the method's limitations, elevated computational requirements should be noted, the need for sufficient volume of representative data for effective training, and limited model interpretability compared to simple heuristic rules or decision trees. It should also be noted that recent Transformer-based architectures [10, 11] and alternative deep learning approaches [12, 13] may offer advantages for longer forecasting horizons or multivariate settings, though their benefits over well-tuned recurrent models are not universal across all problem domains [14, 15].

CONCLUSION

Based on the obtained results, the following conclusions can be drawn:

The proposed method for optimal object placement based on LSTM recurrent neural networks provides high prediction accuracy and spatial distribution quality with acceptable processing time and minimal computational resources. Achieving 85% prediction accuracy (versus 79% for the best classical method, Random Forest, with $p=0.0023$) and 8.5% improvement in spatial balance index (0.89 versus 0.82) confirms the effectiveness of applying recurrent architecture to the optimal placement problem.

The ability of LSTM to model long-term temporal dependencies through forget gate, input gate, and output gate mechanisms provides significant advantages over classical fixed-window methods that lose information about the sequential nature of data. Statistical significance of improvements ($p<0.05$) and robustness to non-stationarity confirm the reliability of the method.

The placement decision time (18.9 ms) remains practical for real-time systems despite increased training requirements. Therefore, the method can be used cost-effectively to improve the quality of object placement management systems in various application domains where prediction accuracy and spatial optimization are critical.

Directions for further research include comparison with modern attention-based architectures (Transformer, Temporal Fusion Transformer [10], Informer [11]), investigation of GRU as a lighter recurrent alternative, application of reinforcement learning methods for direct optimization of placement strategy, and validation on real data from management systems of various object types.

ACKNOWLEDGMENTS AND FUNDING SOURCES

The authors received no financial support for the research, writing, and/or publication of this article.

COMPLIANCE WITH ETHICAL STANDARDS

The authors declare that the research was conducted in the absence of any potential conflict of interest.

AUTHOR CONTRIBUTIONS

Conceptualization, [Z. L.]; methodology, [Z. L.]; validation, [O. T.]; formal analysis, [O. T.]; investigation, [O. T.]; resources, [O. T.]; data curation, [O. T.]; writing – original draft preparation, [O. T.]; writing – review and editing, [Z. L.]; visualization, [O. T.]; supervision, [Z. L.]; project administration, [Z. L.]; funding acquisition, [Z. L.].

All authors have read and agreed to the published version of the manuscript.

REFERENCES

- [1] Erl, T., Puttini, R., & Mahmood, Z. (2013). *Cloud computing: Concepts, technology & architecture*. Pearson Education.
- [2] Oussous, A., Benjelloun, F. Z., Ait Lahcen, A., & Belfkih, S. (2018). Big Data technologies: A survey. *Journal of King Saud University - Computer and Information Sciences*, 30(4), 431–448. <https://doi.org/10.1016/j.jksuci.2017.06.001>
- [3] Pinedo, M. (2022). *Scheduling: Theory, algorithms and systems* (6th ed.). Springer. <https://doi.org/10.1007/978-3-031-05921-6>
- [4] Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8), 1735–1780. <https://doi.org/10.1162/neco.1997.9.8.1735>
- [5] Tsemko, A., & Matskiv, M. (2025). Problems of using neural networks to predict the price of virtual assets. *Electronics and Information Technologies*, 29, 69–78. <https://doi.org/10.30970/eli.29.7>
- [6] Sarker, I. H. (2021). Deep learning: A comprehensive overview on techniques, taxonomy, applications and research directions. *SN Computer Science*, 2(6), 420. <https://doi.org/10.1007/s42979-021-00815-1>
- [7] Zhang, A., Lipton, Z. C., Li, M., & Smola, A. J. (2023). *Dive into deep learning*. Cambridge University Press.
- [8] Torres, J. F., Hadjout, D., Sebaa, A., Martínez-Álvarez, F., & Troncoso, A. (2021). Deep learning for time series forecasting: A survey. *Big Data*, 9(1), 3–21. <https://doi.org/10.1089/big.2020.0159>
- [9] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 30, 5998–6008.
- [10] Lim, B., Arık, S. Ö., Loeff, N., & Pfister, T. (2021). Temporal fusion transformers for interpretable multi-horizon time series forecasting. *International Journal of Forecasting*, 37(4), 1748–1764. <https://doi.org/10.1016/j.ijforecast.2021.03.012>
- [11] Zhou, H., Zhang, S., Peng, J., Zhang, S., Li, J., Xiong, H., & Zhang, W. (2021). Informer: Beyond efficient transformer for long sequence time-series forecasting. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(12), 11106–11115. <https://doi.org/10.1609/aaai.v35i12.17325>
- [12] Smyl, S. (2020). A hybrid method of exponential smoothing and recurrent neural networks for time series forecasting. *International Journal of Forecasting*, 36(1), 75–85. <https://doi.org/10.1016/j.ijforecast.2019.03.017>
- [13] Salinas, D., Flunkert, V., Gasthaus, J., & Januschowski, T. (2020). DeepAR: Probabilistic forecasting with autoregressive recurrent networks. *International Journal of Forecasting*, 36(3), 1181–1191. <https://doi.org/10.1016/j.ijforecast.2019.07.001>

- [14] Lara-Benítez, P., Carranza-García, M., & Riquelme, J. C. (2021). An experimental review on deep learning architectures for time series forecasting. *International Journal of Neural Systems*, 31(5), 2130001. <https://doi.org/10.1142/S0129065721300011>
- [15] Hewamalage, H., Bergmeir, C., & Bandara, K. (2021). Recurrent neural networks for time series forecasting: Current status and future directions. *International Journal of Forecasting*, 37(1), 388–427. <https://doi.org/10.1016/j.ijforecast.2020.06.008>
-

РЕКУРЕНТНА НЕЙРОННА МЕРЕЖА НА ОСНОВІ ДОВГОЇ КОРОТКОЧАСНОЇ ПАМ'ЯТІ ДЛЯ ПРОГНОЗУВАННЯ СТАНІВ ТА ОПТИМІЗАЦІЇ РОЗПОДІЛУ РЕСУРСІВ У РОЗПОДІЛЕНИХ СИСТЕМАХ

Зіновій Любунь , **Олег Терещук** 

zinoviy.lyubun@lnu.edu.ua, Oleh.Tereshchuk@lnu.edu.ua

Кафедра радіофізики та комп'ютерних технологій
Львівський національний університет імені Івана Франка
вул. Тарнавського, 107, м. Львів, 79017, Україна

АНОТАЦІЯ

Вступ. У статті розглянуто метод на основі нейронної мережі з архітектурою довгої короткочасної пам'яті (ДКЧП) для оптимального розподілу ресурсів у розподілених системах. Розроблений алгоритм забезпечує високу точність прогнозування стану ресурсів та оптимальний просторовий розподіл з мінімальним часом обробки. Актуальність дослідження зумовлена зростаючою потребою в інтелектуальній автоматизації процесів управління ресурсами на об'єктах сервісної інфраструктури. Як практичну демонстрацію ефективності методу використано систему управління комірками у спортивних закладах.

Матеріали та методи. Для розв'язання задач прогнозування та оптимізації запропоновано архітектуру на основі ДКЧП з 32 прихованими нейронами та довжиною послідовності 10 часових кроків. Модель LSTM обробляє послідовні дані про зайнятість для виявлення часових залежностей та формування імовірнісних оцінок майбутніх станів ресурсів. Розроблено багатофакторну функцію оцінювання для перетворення прогнозів на оптимальні рішення щодо розподілу з урахуванням просторових обмежень та уподобань користувачів. Метод порівнюється з класичними підходами, такими як евристичні алгоритми (послідовний розподіл, кругова черга), статистичні моделі часових рядів (авторегресійного інтегрованого ковзного середнього, експоненціальне згладжування) та методи машинного навчання (логістична регресія, випадковий ліс, градієнтне підсилювання).

Результати. Використання нейронних мереж з архітектурою ДКЧП для розв'язання задачі прогнозування дозволяє досягти точності 85%, що є статистично значущо вищим за випадковий ліс (79%, $p = 0.0023$) та ARIMA (68%, $p = 0.0001$). Індекс просторового балансу покращився на 8,5% порівняно з найкращим класичним методом (0,89 проти 0,82). Час висновку залишається прийнятним для застосувань реального часу (18,9 мс на прогноз).

Висновки. Запропонований метод на основі ДКЧП демонструє задовільну точність у прогнозуванні станів ресурсів та оптимізації їх розподілу в мінімальні терміни. Здатність моделювати довгострокові темпоральні залежності надає значні переваги

над класичними методами з фіксованим вікном. Тому метод може бути ефективно застосований для підвищення функціональності систем управління розподіленими ресурсами.

Ключові слова: нейронна мережа ДКЧП, розподіл ресурсів, прогнозування часових рядів, рекурентні нейронні мережі, оптимізація.

Received / Одержано
21 February, 2026

Revised / Доопрацьовано
13 March, 2026

Accepted / Прийнято
16 March, 2026

Published / Опубліковано
30 March, 2026

UDC: 004.652.4

INTELLIGENT ANALYSIS OF PERFORMANCE RESULTS BASED ON OBJECT-RELATIONAL MAPPING STRATEGIES AND FOREIGN KEY CONSTRAINTS IN SQL DATABASES

Oleksandra Rybak , Oleh Husak , Roman Mysiuk *

Department of System Design
Ivan Franko National University of Lviv,
50 Drahomanova St., 79005 Lviv, Ukraine

Rybak, O., Husak, O., & Mysiuk, R. (2026). Intelligent Analysis of Performance Results Based on Object-Relational Mapping Strategies and Foreign Key Constraints in SQL Databases. *Electronics and Information Technologies*, 33, 97–112. <https://doi.org/10.30970/eli.33.8>

ABSTRACT

Background. The rapid expansion of data-driven applications has increased the importance of efficient query execution in relational database systems, where even minor inefficiencies can significantly affect overall performance. Although Object-Relational Mapping (ORM) frameworks simplify development and improve maintainability, their abstraction layer can introduce measurable overhead, and the impact of foreign key constraints on execution speed remains a practical concern, particularly in microservice architectures that follow the “Database per Service” principle.

Materials and Methods. An experimental information system is developed using a relational database and the SQLAlchemy ORM framework, with a schema that includes one-to-one, one-to-many, and many-to-many relationships tested both with and without foreign key constraints. Three representative queries retrieving booking details, aggregating related records, and calculating total payments are executed using raw SQL and ORM approaches, while an intelligent algorithm analyzed performance, detected potential $N+1$ query risks, and recommended optimal strategies such as explicit JOINS.

Results and Discussion. Raw SQL consistently demonstrated superior performance across all scenarios. The most significant disparity occurred in ORM implementations affected by the $N+1$ problem, where execution time exceeded that of equivalent SQL queries by more than an order of magnitude. Aggregation queries showed smaller yet consistent overhead. The presence or absence of foreign key constraints had a negligible influence on raw SQL performance, with differences remaining within experimental variance. Explicit JOIN usage in ORM substantially reduced overhead compared to implicit relationship navigation. The intelligent analysis accurately predicted high-risk queries and provided effective strategy recommendations, confirmed by empirical results.

Conclusion. ORM frameworks improve productivity and maintainability but introduce measurable overhead. Raw SQL remains preferable for performance-critical tasks, while foreign key constraints do not significantly degrade execution speed. Intelligent performance analysis supports balanced decisions between efficiency and maintainability in complex relational systems.

Keywords: relational databases, SQL performance, ORM, decision support system, intelligent analysis, database design



© 2026 Oleksandra Rybak et al. Published by the Ivan Franko National University of Lviv on behalf of Електроніка та інформаційні технології / Electronics and Information Technologies. This is an Open Access article distributed under the terms of the [Creative Commons Attribution 4.0 License](https://creativecommons.org/licenses/by/4.0/) which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

INTRODUCTION

Relational databases remain a fundamental component of modern information systems due to their strong consistency guarantees and well-defined data relationships. As applications scale, query performance becomes a critical factor, particularly in systems with complex inter-table relationships. Developers frequently rely on Object-Relational Mapping (ORM) frameworks to simplify database access, improve code maintainability, and reduce development time. However, ORMs introduce additional abstraction layers that may negatively affect query performance. A well-known issue in ORM-based systems is the $N+1$ query problem, which can lead to significant inefficiencies when navigating relationships. The essence of this problem is that when performing an initial database query that returns N objects, another query may automatically be executed for each of these objects, resulting in N additional queries. This creates a significant performance penalty, especially in applications that use ORMs, since iterating over a collection of objects in object-oriented language code naturally generates additional database queries. One common method is the use of eager loading or batch fetching, when the required data is loaded using a single combined or batch query, reducing the number of database calls. Also in practice, SQL query optimization and code refactoring are used to avoid iterations over collections with separate calls to ORM methods. For example, in article [1] an automated approach to refactoring was proposed that uses static analysis to detect and eliminate $N+1$ cases, transforming many small queries into one efficient one, which significantly improves performance.

At the same time, modern software architectures such as microservices promote the principle of Database per Service, emphasizing isolated and independently managed databases. This raises practical questions regarding the necessity and performance impact of foreign key constraints and ORM-managed relationships in high-performance systems. The goal of this study is to experimentally evaluate how different raw SQL and ORM query approaches behave in terms of performance under varying relationship configurations, and to quantify the overhead introduced by ORM abstractions.

This study aims to conduct a quantitative performance evaluation of ORM-based queries and raw SQL in relational database systems with inter-table relationships, focusing on execution time, resource usage, and the impact of ORM-managed relationships.

A quantitative performance evaluation comparing Prisma ORM and raw SQL demonstrated that raw SQL consistently outperformed Prisma across all query types. The composite performance index showed that raw SQL is approximately 2.26 times more efficient, particularly in nested and bulk operations where ORM abstraction introduced measurable latency and higher resource consumption [2].

The Cosmos-specific ODM outperformed raw SQL and generalized ORM-style approaches across all evaluated metrics, including task completion time, error frequency, and perceived cognitive load [3].

Comparative analyses of ORM tools have demonstrated measurable differences in execution performance, memory usage, and query optimization across platforms and frameworks [4, 5]. A comparative analysis of Entity Framework Core, Dapper, and LINQ to DB indicated that Dapper provided the best performance and the lowest memory consumption among the evaluated systems [6, 7]. Across all three experiments in [8], SQL consistently achieved faster execution times than ORM while maintaining nearly identical memory usage. Although ORM occasionally exhibited longer maximum execution times, SQL proved to be the more efficient solution for performance optimization. A query repository with automatic SQL query classification is presented to enforce data-privacy directives [9] and implemented optimizations to reduce the performance overhead caused by intercepting queries through a JDBC proxy, achieving classification latencies as low as 0.35 ms while still using graph-based metadata for cases requiring higher precision.

ORM tools can simplify development but may negatively affect relational query performance, producing inefficient query patterns and necessitating further investigation into mitigation strategies [10].

Summarizing this analysis of sources, it is worth noting that [2-10] largely evaluates ORM performance in isolation or with simple metrics, without considering the combined impact of foreign key constraints, relationship-heavy queries, and modern architectures. This work fills this gap by providing an intelligent, quantitative analysis of ORM strategies alongside FK effects for practical performance optimization.

MATERIALS AND METHODS

The experimental evaluation is conducted using a relational database designed to model real-world inter-table relationships, including one-to-one, one-to-many, and many-to-many associations. Two database configurations are implemented: one with enforced foreign key constraints and one without foreign keys, enabling the assessment of their impact on query performance.

Query execution is implemented using two approaches: raw SQL queries and ORM-based queries developed with SQLAlchemy. For the ORM approach, both lazy loading (default relationship navigation) and explicit join strategies are applied. All SQL queries are manually optimized and executed using the same database engine to ensure consistency. Each experiment included three representative query types: a complex multi-table join, a count aggregation query, and a sum aggregation query.

Performance measurements are obtained by executing each query multiple times under identical conditions and recording execution time at the application level. The mean execution time and standard deviation are calculated for each query and configuration. Additional scalability experiments are performed by increasing the dataset size to evaluate the impact of data volume on query performance. All experiments are executed on the same hardware and software environment to eliminate external variability.

Software and environment

The experimental study is conducted using Python 3.13.0 and MySQL Server 8.0+. Database interactions are implemented via SQLAlchemy as the ORM layer and PyMySQL as the database driver. Data processing and statistical analysis are performed using Pandas, while visualization is carried out with Matplotlib. Synthetic test data are generated using the Faker library. All experiments are executed on a local machine equipped with an 11th Gen Intel® Core™ i3-1115G4 (3.00 GHz) processor, 8 GB RAM, and a 238 GB SSD, running Windows 11 (64-bit, x64 architecture). The database server is deployed locally to ensure consistent benchmarking conditions.

Database design

A relational database schema representing a transportation booking system is designed. The dataset consisted of synthetic but structurally realistic records representing clients, bookings, trips, payments, vehicles, and related entities. The scheme supported 1:1, 1:N, and M:N relationships. Two database configurations are created with and without foreign key constraints.

As part of the experiment, the trips_db database (**Fig. 1**) is filled with a controlled set of test data created using SQLAlchemy's ORM by completely cleaning the tables and adding a fixed number of records (2 for each entity), which allows you to focus on analyzing the impact of ORM abstractions and cross-table relationships on the performance of nested and multi-table queries regardless of the amount of data.

The tables in the studied database are indexed. Each table contains a primary key attribute (e.g., *booking_id*, *trip_id*, *client_id*, *vehicle_id*, *driver_id*), which is automatically indexed by the database management system. These attributes are used in foreign key relationships between tables such as bookings, routes, payments, and trips, which enables

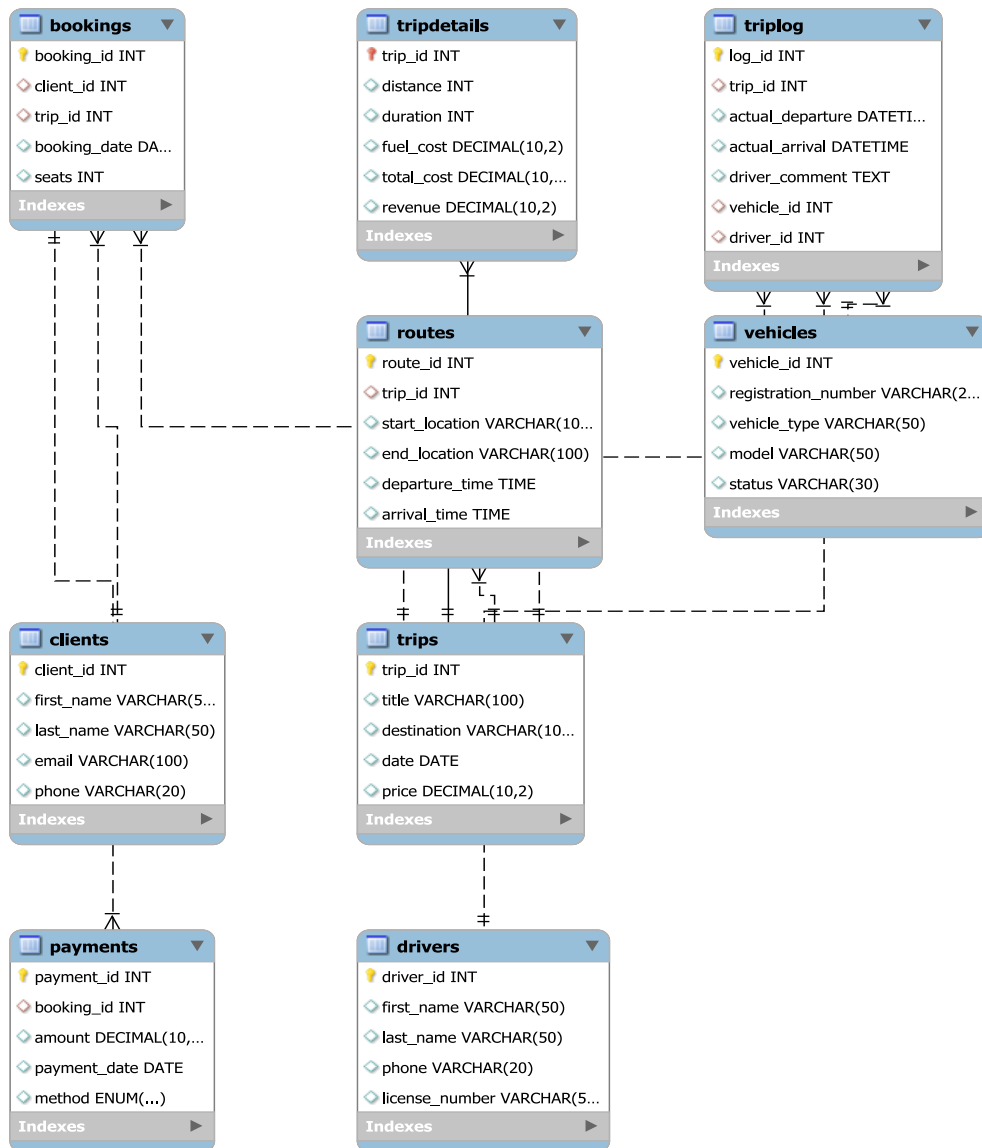


Fig. 1. Entity-relationship diagram of the *trips_db* database.

efficient join operations and improves query performance. The presence of indexes in the tables of the database under study has a critical impact on performance, since in the $N+1$ query problem, the absence of indexes leads to numerous full table scans.

The database is populated with a minimal yet structurally complete synthetic dataset covering all entities and relationships required for controlled performance evaluation. Data covers all 9 tables and all types of relationships (1:1, 1:M, M:N); the scenario is deterministic (same results on each run); the goal is to compare query approaches, not stress-test the database.

The Triplog entity acts as a central point for recording trip performance, unifying all key resources, with the Trips, Drivers, and Vehicles tables linked to Triplog in a one-to-many (1:M) relationship.

Test scenarios

Three representative queries are selected:

- Retrieval of booking and trip execution details for a specific client (multi-table join).
- Counting the number of bookings per trip (aggregation with grouping).
- Calculation of the total sum of payments (aggregation).

Each query is executed using four approaches ORM with foreign keys (automatic relationship navigation), raw SQL with foreign keys, ORM without foreign keys (explicit JOINS).

All tests are executed multiple times on the same dataset, and the recorded times are averaged to minimize measurement noise. Query execution time is measured at the application level by recording timestamps immediately before and after query execution using the datetime library.

In the implementation [11], Data Definition Language (DDL) elements are represented through ORM models (Client, Trip, Booking, Triplog, Driver, Vehicle, Payment), where the declarative approach is used to define database tables, primary and foreign keys, attribute data types, and inter-table relationships. Overall, the core functionality belongs to Data Manipulation Language (DML), as it focuses on data retrieval and aggregation through SELECT, JOIN, GROUP BY, and SUM operations implemented both in pure SQL and through ORM abstractions (count, sum).

Intelligent query strategy decision

The rules that are implemented based on the ORM slowdown factor, defined as: $S = T_{ORM}/T_{SQL}$, where T_{ORM} and T_{SQL} denote the mean execution times of the ORM-based and raw SQL queries, respectively (Fig. 2). This metric quantifies the relative performance degradation introduced by ORM abstractions to assess the relative inefficiency of different data access strategies.

Based on the measured slowdown factors, queries are assigned to one of three distinct risk levels to systematically evaluate the performance impact of ORM strategies. Queries with a slowdown factor $S > 10$ are classified as Critical Inefficiency (N+1 Risk),

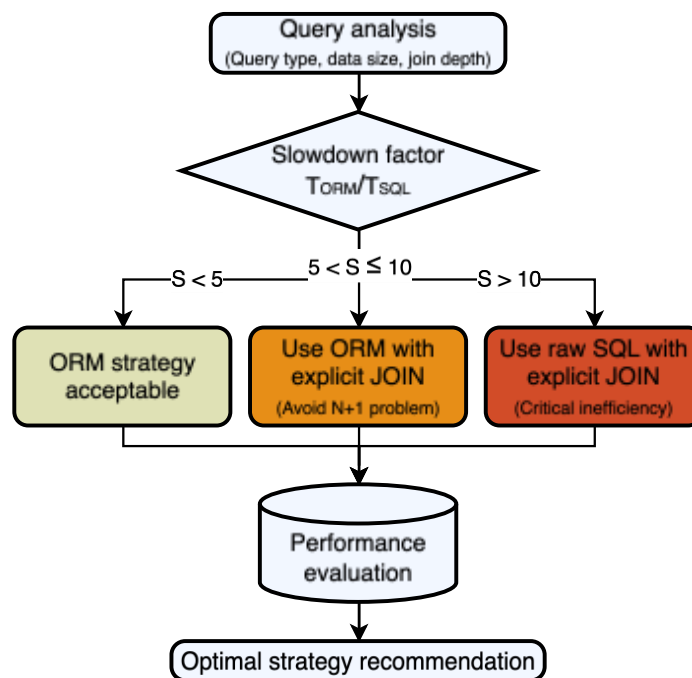


Fig. 2. Intelligent query strategy decision framework.

indicating that lazy-loading ORM patterns ($N+1$ queries) substantially increase execution time and are prone to severe performance degradation, particularly in large-scale datasets. Queries with $5 < S \leq 10$ are classified as Moderate Overhead, representing cases where ORM abstractions introduce noticeable performance loss, which can often be mitigated through explicit JOINS or eager-loading strategies. Finally, queries with $S \leq 5$ are considered Acceptable Performance, indicating low-risk scenarios where ORM overhead is minimal.

A decision tree is mathematically designed to model exactly this kind of conditional branching and provides robustness on small datasets, minimal hyperparameter tuning, and the ability to model categorical outputs directly. It also operates independently of feature normalization and naturally supports future feature expansion, such as join depth, table count, row volume, index availability, and aggregation type. While it provides interpretable, threshold-based rule extraction with minimal preprocessing, decision trees may overfit and become unstable on extremely small datasets, and alternative models like linear regression, logistic regression, neural networks, and SVMs were avoided due to limited interpretability or incompatibility with threshold-based decision logic. From a technical perspective, we use *DecisionTreeClassifier*, which is trained on the experimental dataset. The model automatically learned threshold-based rules for strategy selection. The trained tree is exported as a vector-based SVG diagram using *Graphviz* for interpretability and reproducibility. This approach transforms deterministic rule logic into a data-driven adaptive optimization model. By providing a structured and interpretable metric, the risk classification framework serves as a reliable tool for both benchmarking and predictive performance analysis of relational database queries.

Statistical Hypothesis Testing

To validate the observed differences in query execution times and assess their statistical significance [12, 13], paired t-tests and one-way ANOVA are conducted. The paired t-test is used to compare the mean execution time of ORM queries to the corresponding SQL queries under identical conditions, with the null hypothesis stating that no significant difference exists between ORM and SQL mean execution times. Execution times are recorded over multiple repeated runs for each query type, and paired t-tests are applied using the function *ttest_rel()* from SciPy library of Python. The one-way ANOVA is applied to evaluate whether mean execution times differed across multiple query execution strategies, including ORM with foreign keys, SQL with foreign keys, ORM with explicit JOINS, and SQL without foreign keys. In this case, the null hypothesis posited that all group means are equal, and the test is implemented using the function *f_oneway()* from SciPy library of Python.

These statistical tests provide rigorous validation of the performance measurements. While raw execution times suggest that ORM may incur higher overhead, the t-test and ANOVA determine whether these differences are statistically meaningful or likely attributable to random variability in repeated measurements. The inclusion of hypothesis testing therefore, enhances the reliability of performance conclusions and informs the adaptive strategy recommendation framework.

RESULTS AND DISCUSSION

The performance of ORM-based queries and raw SQL is compared. The results are showing that SQL consistently executes faster, especially for complex multi-table joins where ORM lazy loading triggers the $N+1$ problem. Scalability analysis demonstrates that ORM execution time grows more rapidly with increasing dataset sizes, while SQL scales more linearly, highlighting potential bottlenecks in high-load scenarios. Intelligent strategy decisions classify queries by risk level, recommending raw SQL for critical inefficiencies and ORM with explicit joins for moderate overhead, providing actionable guidance for developers. Statistical hypothesis testing, including paired t-tests and ANOVA, confirms that the observed performance differences are significant and not due to random variation.

Overall, the section integrates empirical measurements, predictive modeling, and intelligent analysis to support evidence-based optimization of query execution strategies.

Comparison of ORM and SQL performance

A comparative experiment is conducted to assess query efficiency using four data access approaches:

- Pure SQL with explicit JOINS and foreign key constraints.
- ORM with relationship-based navigation (automatic joins via SQLAlchemy).
- SQL without foreign key constraints (manual joins).
- ORM without relationships (explicit joins defined in code).

The results in **Table 1** indicate a consistent performance advantage of pure SQL over ORM-based implementations when foreign keys are present.

Table 1. Comparison of query performance using Foreign Keys

#	Description	Results		
		ORM (with FK), ms	SQL (with FK), ms	SQL faster, times
Q1	Details of the trip (5 tables, N+1)	32.25	1.93	~16.7
Q2	Number of bookings per trip (JOIN, COUNT)	5.03	0.63	~8.0
Q3	Total Payments (SUM)	5.52	0.68	~8.1

The largest discrepancy is observed in the complex multi-table query (Q1), where ORM with relationship navigation exhibited the $N+1$ query pattern and required 32.25 ms compared to 1.93 ms for SQL, making SQL approximately 16.7 times faster. For aggregation queries (Q2 and Q3), SQL remained about eight times faster than ORM. These findings confirm that although ORM improves abstraction and developer productivity, it may introduce significant performance overhead in complex relational queries, particularly when implicit loading strategies are used.

Table 2 presents the performance comparison of ORM and pure SQL queries executed without enforcing foreign key constraints. The results demonstrate that even in the absence of foreign keys, pure SQL maintains a consistent performance advantage across all query types. For the complex multi-table retrieval (Q1), ORM with explicit JOIN requires 13.72 ms, whereas SQL completes the same operation in 4.00 ms, making SQL approximately 3.4 times faster. Similar trends are observed for aggregation queries: in Q2 (COUNT with JOIN), SQL is about 3.6 times faster, and in Q3 (SUM aggregation), the performance gap increases to approximately six times. Importantly, compared to the earlier $N+1$ ORM implementation, the explicit use of JOIN within ORM significantly reduces

Table 2. Comparison of query performance without Foreign Keys

#	Description	Results		
		ORM (without FK, with JOIN), ms	SQL (without FK), ms	SQL faster, times
Q1	Details of the trip (5 tables, N+1)	13.72	4.00	~3.4
Q2	Number of bookings per trip (JOIN, COUNT)	5.78	1.59	~3.6
Q3	Total Payments (SUM)	5.31	0.89	~6.0

execution time, confirming that query structure has a critical impact on performance. Nevertheless, SQL remains more efficient due to lower abstraction overhead.

Table 3 provides a consolidated comparison of all four experimental scenarios, enabling a holistic evaluation of abstraction level and foreign key enforcement effects. Scenario I (ORM with foreign keys and relationship navigation) shows the highest execution times, particularly for Q1, reflecting the impact of the *N+1* query pattern. Scenario II (SQL with foreign keys) demonstrates the best overall performance. Scenario III (ORM with explicit joins and without foreign keys) substantially improves over Scenario I, confirming that manual join specification mitigates inefficiencies caused by implicit loading. Scenario IV (SQL without foreign keys) remains faster than both ORM-based approaches, though slightly slower than SQL with foreign keys. Overall, the table confirms two central findings: explicit JOIN usage in ORM significantly enhances efficiency compared to *N+1* navigation, and pure SQL consistently delivers superior performance across all tested conditions.

Table 3. Consolidated comparison of all approaches

Scenario	Description	Results		
		Q1 (JOIN), ms	Q2 (COUNT), ms	Q3 (SUM), ms
I	ORM (with FK)	32.25	5.03	5.52
II	SQL (with FK)	1.93	0.63	0.68
III	ORM (explicit JOIN) without FK	13.72	5.78	5.31
IV	SQL without FK	4.00	1.59	0.89

This confirms that replacing implicit navigation with explicit joins significantly enhances ORM efficiency. However, even after optimization, ORM remains slower than pure SQL in all scenarios. For example, in Q1 without foreign keys, the ORM runs approximately 3.4 times slower than raw SQL. Similar gaps persist in Q2 and Q3.

The substantiates that while explicit join usage in ORM enhances performance, pure SQL consistently outperforms ORM across all measured scenarios, both in terms of execution time and stability (**Fig. 3**).

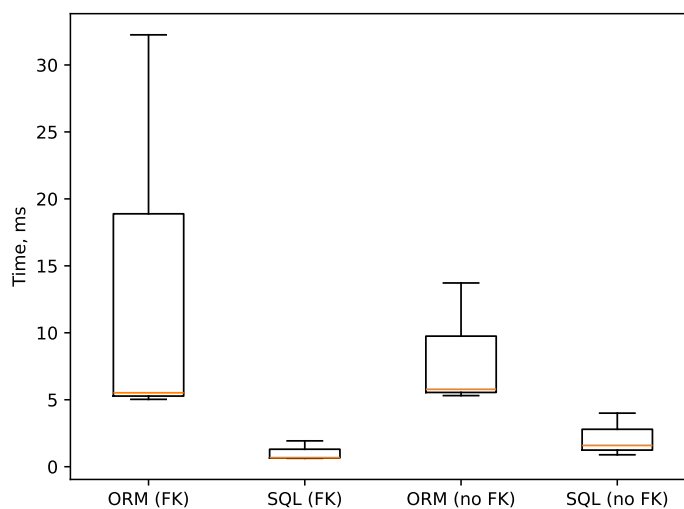


Fig. 3. Distribution of Query Execution Time for ORM and SQL Approaches.

The ORM $N+1$ queries can severely degrade performance on complex queries with many related records (Q1), while explicit JOINS are faster, though for simple queries (Q2, Q3) both strategies perform similarly (Fig. 4).

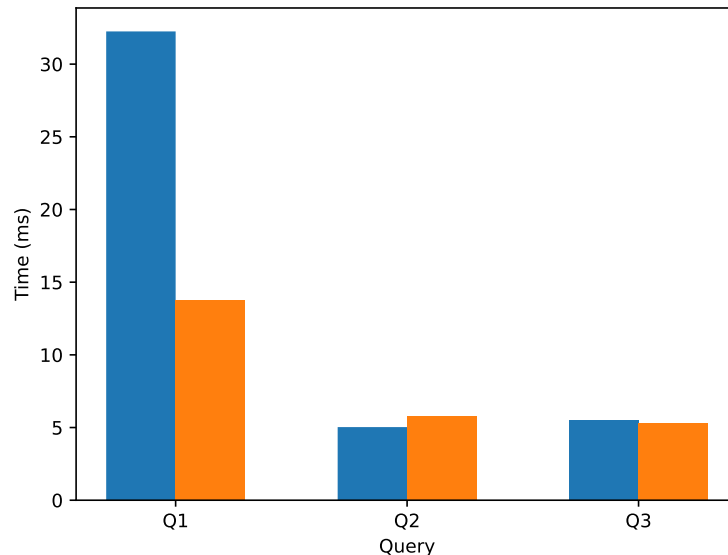


Fig. 4. ORM Strategy Comparison.

Table 4 shows that using foreign keys consistently improves SQL query performance, and while ORM JOINS mitigate $N+1$ overhead, raw SQL remains faster.

The experiments demonstrate that while explicit JOINS in ORM substantially reduce $N+1$ query overhead, pure SQL (with foreign keys) consistently delivers superior performance across all query types and complexities.

A quantitative analysis of query performance under different database access strategies executed with highlighting the influence of foreign key enforcement and ORM abstraction (Fig. 5).

The presence of foreign key constraints substantially improves SQL execution efficiency, with the most pronounced reduction observed in the complex multi-table query Q1, whereas simpler aggregation queries Q2 and Q3 benefit to a lesser extent (Fig. 5a). The performance drawback (Fig. 5b) associated with ORM-based access relative to pure SQL, ORM with foreign keys exhibits a pronounced slowdown for Q1, attributable to the $N+1$ query pattern, while the overhead for Q2 and Q3 is comparatively minor. These results underscore two principal findings: first, foreign key enforcement materially enhances SQL query performance, particularly in queries involving multiple related tables; second, despite the use of relationship navigation, ORM introduces significant execution overhead in complex query scenarios, reinforcing the trade-off between abstraction convenience and raw performance.

Table 4. Relative Impact of Foreign Keys on SQL Query Performance

#	Results	
	SQL with FK, ms	SQL without FK, ms
Q1	1.93	4.00
Q2	0.63	1.59
Q3	0.68	0.89

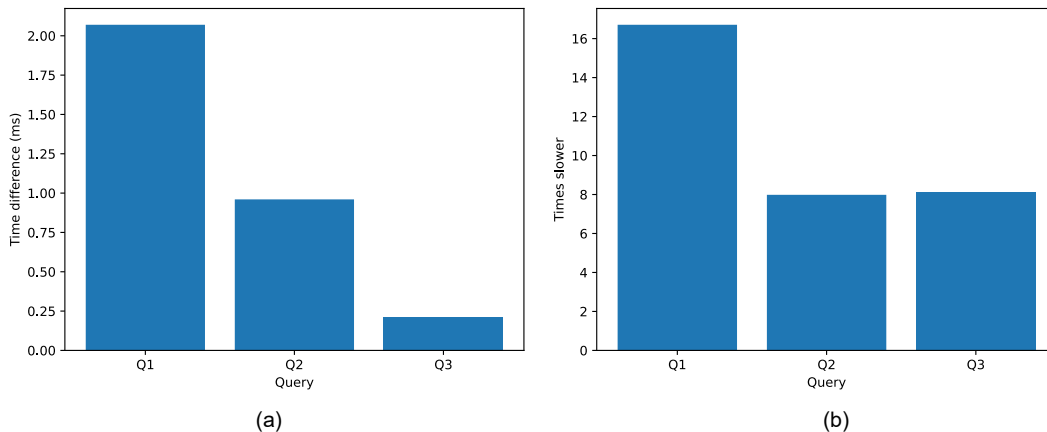


Fig. 5. Analysis of performance results: a) Impact of Foreign Keys on SQL Performance (difference between SQL with FK and SQL without FK), b) ORM Slowdown Factor (ORM with FK to SQL with FK).

The experiments demonstrate that raw SQL consistently outperforms ORM-based queries across all tested scenarios. The most significant performance gap occurs in complex multi-table joins (Q1) affected by the $N+1$ problem, where ORM execution time exceeds SQL by more than an order of magnitude. Aggregation queries (Q2 and Q3) show moderate ORM overhead, while the presence or absence of foreign key constraints has a negligible impact on SQL performance. Using explicit JOINS in ORM significantly reduces overhead compared to lazy relationship navigation. Furthermore, the intelligent algorithm accurately identifies high-risk queries, quantifies performance degradation, and recommends optimal execution strategies, validating its utility as a decision-support tool for database optimization.

Scalability analysis

The experiments are repeated to evaluate scalability with increasing dataset size. To analyze scalability, the scenario of increasing the number of records in the key tables of the `trips_db` database is considered, namely: table `clients` 1,100 records, `trips` 500 records, `bookings` 1,200 records, and `payments` 1,200 records. The results, shown in [Figure 3](#), reveal fundamentally different growth patterns. ORM queries affected by the $N+1$ problem exhibit non-linear growth in execution time, with performance degradation becoming critical at larger scales. In contrast, raw SQL queries scale almost linearly with data size, maintaining low execution times even for larger datasets. Eager loading strategies show improved scalability compared to lazy loading, following a near-linear trend. However, their execution time remains consistently higher than that of raw SQL. These findings confirm that the choice of data access strategy has a direct impact on system scalability. [Table 5](#) presents the statistical analysis of repeated query executions.

The mean execution time, standard deviation, and coefficient of variation (CV) are reported for each query type [11]. The results indicate that ORM queries, particularly those affected by the $N+1$ problem, exhibit higher variability and longer execution times compared to equivalent SQL queries. The CV values highlight the relative dispersion of execution times, confirming that ORM $N+1$ queries are the least stable, while aggregation queries (SUM, COUNT) are more consistent. Q1 exhibits a critical inefficiency caused by the $N+1$ query problem, with a slowdown factor significantly exceeding the defined threshold. The recommended strategy is to use raw SQL with explicit JOINS to eliminate excessive query overhead. Q2 and Q3 show moderate ORM overhead, indicating that while ORM abstractions introduce some performance loss, it is not critical. For these queries, the system recommends using ORM with explicit JOINS to maintain ORM convenience while mitigating unnecessary query inflation.

Statistical Hypothesis Testing

Paired t-tests and one-way ANOVA are conducted to evaluate whether ORM and SQL execution times differed significantly. The paired t-test ($t = 0.993$, $p = 0.377$) indicates no significant difference between ORM and SQL for matched queries. Similarly, ANOVA across all tested approaches ($F = 1.666$, $p = 0.251$) shows that mean execution times do not differ significantly at the 5% significance level. These results suggest that, although ORM queries exhibit numerically higher execution times, the observed differences are not statistically significant for the sampled runs, highlighting the influence of measurement variability or limited sample size.

Table 5. Relative Impact of Foreign Keys on SQL Query Performance

Scenario	Results					Mean execution time, s	Standard Deviation, s	Coefficient of Variation
	1	2	3	4	5			
ORM Search ($N+1$)	0.382	0.002	0.002	0.001	0.002	0.078	0.1521	1.955
SQL Search (JOIN)	0.039	0.004	0.001	0.002	0.002	0.009	0.0147	1.535
ORM Grouping (Count)	0.028	0.005	0.009	0.006	0.009	0.011	0.0085	0.741
SQL Grouping (Count)	0.007	0.008	0.007	0.010	0.010	0.009	0.0014	0.161
ORM Aggregation (SUM)	0.008	0.001	0.001	0.001	0.001	0.002	0.0028	1.167
SQL Aggregation (SUM)	0.003	0.001	0.001	0.001	0.001	0.001	0.0008	0.571

Intelligent strategy decisions

The intelligent analysis (**Fig. 6**) illustrates the adaptive decision tree derived from experimental performance data, with the root node threshold $\text{Slowdown_ORM_vs_SQL_with_FK} \leq 12.414$ separating queries into two branches. Queries below the threshold (Q2 and Q3) are classified as “Use ORM with explicit JOIN”, reflecting moderate overhead where ORM remains acceptable. Queries exceeding the threshold (Q1) fall into the “Use raw SQL with explicit JOIN” leaf, indicating severe ORM overhead and $N+1$ risk. Testing on the `trips_db_no_fk` database confirmed these predictions, with foreign key removal having minimal impact, demonstrating that the system provides a transparent, threshold-based decision-support tool for selecting query execution strategies.

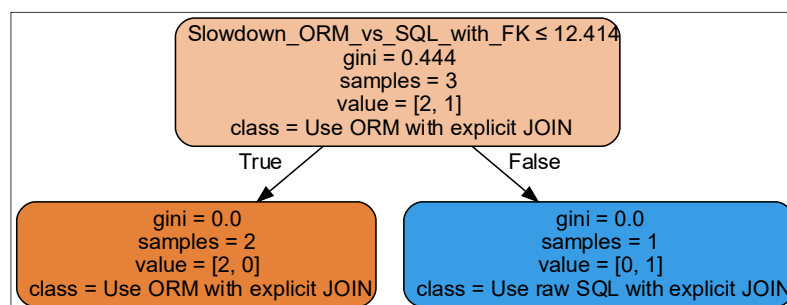


Fig. 6. Results based on queries.

Discussion of results

The quantitative analysis confirms that the primary source of performance degradation is not the presence of foreign key constraints but the ORM execution model itself. The increased execution time and variance observed in ORM-based approaches are caused by additional query generation, data fetching overhead, and object mapping operations. While eager loading significantly improves performance, raw SQL remains the most efficient and stable solution for complex and high-load query scenarios. From a practical perspective, the results suggest that ORM frameworks are suitable for CRUD (create, read, update, delete) operations and moderate workloads, while raw SQL is preferable for performance-critical components, particularly those involving complex joins and large datasets.

The adaptive decision tree effectively classifies queries based on the well-established ORM slowdown metric, with a threshold of 12.414 distinguishing moderate overhead (where ORM with explicit JOINS is acceptable) from severe overhead (favoring raw SQL to avoid $N+1$ issues). This provides developers with a transparent, reproducible decision-support tool to select query execution strategies, anticipate performance bottlenecks, and prevent inefficiencies. Validation on the `trips_db_no_fk` database confirmed the model's predictions, showing that foreign key removal had minimal impact and indicating robustness, although generalizability may be limited across different databases, ORMs, or query patterns. Future work could expand the approach to additional query types, incorporate runtime profiling for real-time decisions, and integrate with automated query optimization to enhance practical applicability.

Although the use of ORM greatly simplifies the process of developing and maintaining program code, this is achieved at the cost of reducing performance. At the same time, the presence of intertable relationships in the database ensures data integrity and, according to the results of the experiment, does not lead to a noticeable deterioration in the performance of raw SQL queries.

Based on [14–17] and our experimental results, we prepared a heatmap to demonstrate the feasibility of using ORM and direct SQL for different types of queries (Fig. 7). The heatmap shows query types on the Y-axis which include CRUD operations, aggregates, and complex multi JOIN queries and technology choice on the X-axis (ORM or SQL), with each cell ranging from 0 (less appropriate) to 1 (highly recommended) to indicate the suitability of a given technology for each query type. Therefore, it is advisable to use ORM in systems where development speed and maintainability are priorities, whereas direct SQL is more appropriate for high-load or latency-critical queries.

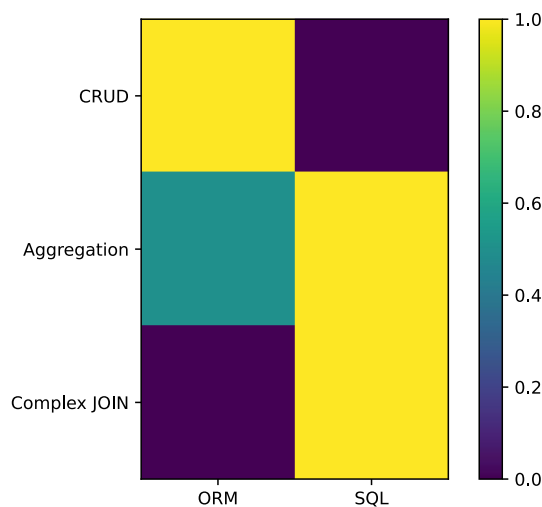


Fig. 7. Technology Sustainability Matrix.

The matrix (**Fig. 7**) demonstrates that ORM is most efficient for simple CRUD operations, moderately suitable for aggregate queries, but less suitable for complex JOINS. On the other hand, direct SQL shows high feasibility for aggregations and complex JOIN queries, providing better performance when working with large amounts of data and complex relationships. The visualization (**Fig. 7**) allows you to quickly evaluate the optimal choice of technology depending on the type of request.

Experimental studies on relational database systems demonstrate that query execution time increases with the number of tables involved in JOIN operations due to the growth of intermediate results and the complexity of join ordering. The presence of indexes on foreign key attributes and efficient join planning significantly influence query performance.

The experiments presented in this study were conducted on a local computer environment, as it is mentioned in the software and environment of the method and materials section, in order to ensure controlled testing conditions and reproducibility of results. However, future work will extend the experimental setup to include remote database servers, which will allow evaluation of query performance under network latency and distributed deployment conditions that are closer to real-world systems.

Future work could expand the model to include additional query types, integrate runtime profiling, and support automated query optimization, extending its practical applicability.

CONCLUSION

This study presented a quantitative performance analysis of SQL and ORM-based query execution in relational databases with inter-table relationships. The experiments are conducted using three representative query types with a complex multi-table join, an aggregation query for counting records, and an aggregation query for summation.

This study also addresses a gap in existing research by systematically evaluating the interplay between ORM strategies and foreign key constraints in complex, relationship-heavy queries, providing actionable insights for developers on when to rely on ORM versus raw SQL for performance-critical operations. Performance is evaluated across four scenarios: ORM with foreign keys, raw SQL with foreign keys, ORM without foreign keys, and raw SQL without foreign keys.

The results demonstrate that raw SQL consistently outperforms ORM-based approaches in all tested scenarios. For the most complex query involving five joined tables, ORM with lazy loading exhibited an average execution time of 32.25 ms, while the equivalent raw SQL query executed in 1.93 ms, making SQL approximately 16.7 times faster. This performance gap is primarily attributed to the $N+1$ query problem inherent in lazy ORM navigation.

For aggregation queries, ORM also showed higher execution times. The booking count query required 5.03 ms using ORM with foreign keys compared to 0.63 ms using raw SQL, while the payment sum query took 5.52 ms with ORM and 0.68 ms with SQL. Thus, raw SQL is approximately 8 times faster for aggregation operations. When foreign key constraints are removed, raw SQL maintains its performance advantage. The complex join query executed in 4.00 ms using SQL without foreign keys, compared to 13.72 ms for ORM with explicit joins, resulting in a 3.4 speedup. Similar trends are observed for aggregation queries, where SQL outperformed ORM by factors ranging from 3.6× to 6.0×.

A comparative evaluation of ORM with explicit joins shows that execution time reduced from 32.25 ms to 13.72 ms, demonstrating a performance improvement of more than 2 times faster. However, even with optimized loading strategies, ORM did not reach the performance level of raw SQL. The analysis also revealed that the presence or absence of foreign key constraints has a negligible impact on raw SQL performance. Observed differences in execution time are below 1 ms and fall within the expected measurement

error range. This indicates that foreign keys primarily contribute to data integrity rather than query performance degradation.

Moreover, the results highlight the potential for intelligent performance analysis frameworks to guide ORM usage in distributed or microservice-based database architectures, where balancing maintainability and efficiency is critical.

To sum up, ORM frameworks provide substantial benefits in terms of development speed, code maintainability, and abstraction, but introduce measurable performance overhead due to additional query generation and object mapping. Raw SQL remains the most efficient solution for performance-critical operations, complex joins, and high-load systems where every millisecond is significant. These findings support a hybrid approach in real-world systems, where ORM is used for standard operations and raw SQL is employed for optimized, high-performance queries.

ACKNOWLEDGMENTS AND FUNDING SOURCES

The authors received no financial support for the research, writing, and/or publication of this article.

COMPLIANCE WITH ETHICAL STANDARDS

The authors declare that the research is conducted in the absence of any conflict of interest.

AUTHOR CONTRIBUTIONS




Conceptualization, [O. R.]; methodology, [O.H.]; validation, [O. H.]; investigation, [R. M.]; writing – original draft preparation, [O. R.]; writing – review and editing, [R. M.]; visualization, [O. R.] supervision, [O. H., R. M.].

REFERENCES

- [1] Turcotte, A., Aldrich, M. W., & Tip, F. (2023). Reformulator: Automated refactoring of the N+1 problem in database-backed applications. Proceedings of the 37th IEEE/ACM International Conference on Automated Software Engineering (ASE '22), Article 84, 1–12. Association for Computing Machinery. <https://doi.org/10.1145/3551349.3556911>
- [2] Yusmita, J. C., Arya, R., Wijaya, J. M., Suryaningrum, K. M., & Siswanto, R. R. (2025). Optimizing database access strategy: A performance analysis comparison of raw SQL and Prisma ORM. *Procedia Computer Science*, 269, 1201–1210. <https://doi.org/10.1016/j.procs.2025.09.061>
- [3] Singh, L. K., Tarai, S. K., Sharma, R., Godiyal, A., & others. (2025). A comparison of developer performance across raw SQL, ORM-inspired interfaces, and a Cosmos-specific ODM. *International Journal of Scientific Research in Engineering and Management*, 9(11), 1–9. <https://doi.org/10.55041/IJSREM54453>
- [4] Hule, K. (2023). Analysis of Different ORM Tools for Data Access Object Tier Generation: A Brief Study. *International Journal of Membrane Science and Technology*. <https://doi.org/10.15379/IJMST.V10I1.2842>
- [5] Güvercin, A. E., & Avenoglu, B. (2022). Performance analysis of object-relational mapping (ORM) tools in .NET 6 environment. *Bilişim Teknolojileri Dergisi*, 15(4), 453–465. <https://doi.org/10.17671/gazibtd.1059516>
- [6] Hule, K., & Ranawat, R. (2023). Analysis of different ORM tools for data access object tier generation: A brief study. *International Journal of Membrane Science and Technology*, 10(1), 1277–1291. <https://doi.org/10.15379/IJMST.V10I1.2842>
- [7] Wiatrowski, T. (2024). Comparative Analysis of ORM Systems for the .NET Platform. *Journal of Computer Sciences Institute*, 31, 97–102. <https://doi.org/10.35784/jcsi.6012>

- [8] Hermanto, B., Parabi, M. I., Sakethi, D., & Azhar, N. (2025). Performance Comparison of Object-Relational Mapping (ORM) and SQL Query in Developing the Booking Service API at PT Tunas Dwipa Matra. *Jurnal Pepadun*, 6(2), 113–119. <https://doi.org/10.23960/pepadun.v6i2.263>
- [9] Schwab, P.K., Röckl, J., Langohr, M.S. et al. (2021) Performance Evaluation of Policy-Based SQL Query Classification for Data-Privacy Compliance. *Datenbank Spektrum* 21, 191–201. <https://doi.org/10.1007/s13222-021-00385-9>
- [10] Colley, D., Stanier, C., & Asaduzzaman, M. (2018). The impact of object-relational mapping frameworks on relational query performance. In 2018 International Conference on Computing, Electronics & Communications Engineering (iCCECE) (pp. 47–52). IEEE. <https://doi.org/10.1109/iCCECOME.2018.8659222>
- [11] coursework_trips_db (Version latest) [GitHub repository]. GitHub. https://github.com/SashaRyback007/coursework_trips_db/tree/master
- [12] Dobson, R. (2018). One-Way Analysis of Variance Test Add-on for SQL Statistics Package. MSSQLTips. <https://www.mssqltips.com/sqlservertip/5712/a-oneway-analysis-of-variance-test-addon-for-the-sql-statistics-package/>
- [13] Pangesa, D. M., Astriani, M. S., & Manuaba, I. (2024). Study on object-relational mapping (ORM) data model performance effects in the oil and gas industry. In Proceedings of the 2024 2nd International Conference on Technology Innovation and Its Applications (ICTIIA), 1–6. <https://doi.org/10.1109/ICTIIA61827.2024.10761468>
- [14] Bonteanu, A. M., & Tudose, C. (2024). Performance analysis and improvement for CRUD operations in relational databases from Java programs using JPA, Hibernate, Spring Data JPA. *Applied Sciences*, 14(7), Article 2743. <https://doi.org/10.3390/app14072743>
- [15] Zhadko-Bazilevych, S. (2025). Analysis of ORM framework approaches for Node.js. *Journal of Computer Sciences Institute*, 37, 426–430. <https://doi.org/10.35784/jcsi.7951>
- [16] Güvercin, A. E., & Avenoglu, B. (2022). Performance Analysis of Object-Relational Mapping (ORM) Tools in .Net 6 Environment. *Bilişim Teknolojileri Dergisi*, 15(4), 453–465. <https://doi.org/10.17671/gazibtd.1059516>
- [17] Marchuk, I., Dyyak, I., & Makar, I. (2023). Performance analysis of database access: Comparison of direct connection, ORM, REST API and GraphQL approaches. In 2023 IEEE 13th International Conference on Electronics and Information Technologies (ELIT) (pp. 174–176). IEEE. <https://doi.org/10.1109/ELIT61488.2023.10310748>

ІНТЕЛЕКТУАЛЬНИЙ АНАЛІЗ РЕЗУЛЬТАТІВ ПРОДУКТИВНОСТІ НА ОСНОВІ СТРАТЕГІЙ ОБ'ЄКТНО-РЕЛЯЦІЙНОГО ВІДОБРАЖЕННЯ ТА ОБМЕЖЕНЬ ЗОВНІШНЬОГО КЛЮЧА В БАЗАХ ДАНИХ SQL

Олександра Рибак , Олег Гусак , Роман Мисюк *

Кафедра системного проектування
Львівський національний університет імені Івана Франка,
вул. Драгоманова 50, м. Львів, 79005, Україна

АНОТАЦІЯ

Вступ. Швидке зростання застосунків, які використовують велику кількість даних, підвищує важливість ефективного виконання запитів у реляційних базах даних. Хоча системи об'єктно-реляційного відображення (ОРВ) спрощують розробку та

покращують підтримуваність коду, їхній рівень абстракції може створювати вимірювані затримки, а вплив обмежень зовнішніх ключів на швидкість виконання запитів залишається практичною проблемою, особливо в мікросервісних архітектурах, що дотримуються принципу «База даних на сервіс».

Матеріали та методи. Інформаційну систему було розроблено на основі реляційної бази даних та програмного каркасу OPB SQLAlchemy, зі схемою, що включала зв'язки один-до-одного, один-до-багатьох та багато-до-багатьох, протестовані як з обмеженнями зовнішніх ключів, так і без них. Було виконано три типові запити: отримання даних бронювання, агрегація пов'язаних записів та обчислення загальних платежів, використовуючи підходи необробленого SQL та OPB, при цьому інтелектуальний алгоритм аналізував продуктивність, виявляв потенційні проблеми $N+1$ запитів та рекомендував оптимальні стратегії, наприклад явні операції з'єднання таблиць JOIN.

Результати. Запити на raw SQL стабільно показували кращу продуктивність у всіх сценаріях. Найбільша різниця спостерігалася в OPB -запитах, уражених проблемою $N+1$ запитів, де час виконання перевищував аналогічні SQL-запити більш ніж у десятки разів. Агрегаційні запити демонстрували менші, але стабільні затримки. Наявність або відсутність обмежень зовнішніх ключів майже не впливала на продуктивність необробленого SQL. Використання явних операцій JOIN в OPB суттєво зменшувало затримки порівняно з неявною навігацією по зв'язках. Інтелектуальний аналіз точно визначав запити з високим рівнем ризику та пропонував ефективні рекомендації, підтверджені експериментально.

Висновки. OPB-системи покращують продуктивність розробки та підтримуваність, але вводять вимірювані затримки. Для критично важливих завдань продуктивності перевагу слід надавати SQL, а обмеження зовнішніх ключів незначно впливають на швидкість виконання. Інтеграція інтелектуального аналізу продуктивності дозволяє приймати зважені рішення між ефективністю та підтримуваністю у складних реляційних системах.

Ключові слова: реляційні бази даних, продуктивність SQL, OPB, система підтримки рішень, інтелектуальний аналіз, проектування баз даних

UDC: 004.932.2

SPATIAL-GEOMETRIC EVALUATION OF LOCAL FEATURES IN MONOCULAR VISUAL ODOMETRY

Andriy Fesiuk* , Yuriy Furgala 

Department of Optoelectronics and Information Technologies
Ivan Franko National University of Lviv,
50 Drahomanova St., 79005 Lviv, Ukraine

Fesiuk A. V., Furgala Y. M. (2026). Spatial-Geometric Evaluation of Local Features in Monocular Visual Odometry. *Electronics and Information Technologies*, 33, 113–130.
<https://doi.org/10.30970/eli.33.9>

ABSTRACT

Background. Monocular visual odometry is an important component of visual navigation systems. However, its accuracy depends on the quality of local features and inter-frame correspondences. In the VO task, not only is geometric consistency important, but also motion observability, the physical validity of the recovered configuration, and the spatial-structural properties of local features. This study aims to provide a comprehensive evaluation of keypoint detection and description methods for monocular visual odometry.

Materials and Methods. The study was conducted on the EuRoC MAV dataset. The ORB, BRISK, AKAZE, KAZE, SIFT, SURF, and SuperPoint methods were analyzed for the number of keypoints, ranging from 200 to 1000. Motion estimation was performed using the essential matrix, the USAC_FAST filter, the recoverPose method, a minimum parallax check, and spatially guided keypoint selection. The accuracy of the recovered trajectory was evaluated using the APE and RPE metrics. To analyze the quality of local features and correspondences, the geometric component, the parallax indicator, the correct cheirality ratio, and metrics of keypoint coverage uniformity, local redundancy, and structural consistency were used. An integral quality indicator was applied to summarize the results.

Results and Discussion. The geometric metrics most often highlight AKAZE and SURF, whereas SuperPoint shows strong performance in terms of spatial characteristics. In terms of the structural consistency of correspondences, SURF consistently demonstrates the best results. As the number of keypoints increases, most methods show an initial improvement followed by saturation, and in some cases, a deterioration of individual characteristics. SURF was found to be the most balanced method across the set of criteria, whereas ORB showed the weakest results in most cases. The correlation analysis showed that the informativeness of the metrics varies by sequence type.

Conclusion. The proposed approach confirmed the relevance of multicriteria evaluation of local features in monocular visual odometry. It was shown that no single metric is universal across all scene types. In contrast, the integral indicator enables the summary of different aspects of quality and a more well-grounded ranking of the methods.

Keywords: monocular visual odometry, keypoint detection, image matching, motion estimation, deep learning, neural networks.

INTRODUCTION

Visual navigation and robotics systems widely employ monocular visual odometry as one of the approaches for estimating camera motion and reconstructing the motion trajectory. Under real operating conditions, the performance of such systems is



© 2026 Andriy Fesiuk & Yuriy Furgala. Published by the Ivan Franko National University of Lviv on behalf of Електроніка та інформаційні технології / Electronics and Information Technologies. This is an Open Access article distributed under the terms of the [Creative Commons Attribution 4.0 License](https://creativecommons.org/licenses/by/4.0/) which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

complicated by dynamic viewpoint changes, illumination variations, motion blur, and the presence of weakly textured scene regions. The accuracy of odometric estimates depends on the quality of correspondences established between neighboring frames, since these correspondences form the basis for determining the relative camera pose. Errors in keypoint matching, instability of the correspondence set, and insufficient observability of inter-frame displacement adversely affect motion estimation accuracy, accumulate over time, and increase trajectory error [1-5].

Keypoints and their descriptors, which are considered in this work as local features, constitute the foundation of monocular visual odometry. In most approaches to local feature analysis, the primary focus is placed on geometric characteristics, in particular, matching accuracy, the inlier ratio, detection repeatability, and descriptor stability [5, 6]. Although these indicators are important, in the task of monocular visual odometry, they do not provide a complete characterization of feature suitability, because correspondence correctness within a geometric model does not guarantee reliable motion recovery. In particular, under small parallax or near-degenerate, geometrically weak configurations, even formally correct correspondences may lead to unstable translation estimation [7]. It is also important to assess how consistent the recovered geometry is at the current step, specifically whether the correctness of depth signs is preserved during camera pose recovery. In addition, the spatial-structural properties of inliers are of considerable importance: non-uniform image coverage, excessive local concentration of points, or weak correspondence to the scene structure reduce the reliability of motion estimation and increase sensitivity to noise [1-4].

Existing studies on visual odometry mainly examine either final trajectory errors or individual characteristics of local features and correspondences [8-11]. While this analysis is helpful, it does not fully reveal the connection between feature properties and the odometry algorithm's actual performance. Previous research [6] showed that for image matching, it is important to consider not only geometric consistency but also the spatial structure of keypoints. However, directly applying this method to monocular visual odometry is inadequate because it also requires accounting for motion observability, the stability of camera pose recovery, and the link between local feature characteristics and trajectory-based metrics.

Modern approaches to image keypoint detection can be divided into classical methods, including SIFT [12], SURF [13], KAZE [14], AKAZE [15], ORB [16], and BRISK [17], and deep learning-based solutions such as SuperPoint [18]. Classical algorithms demonstrate high computational efficiency and mathematical interpretability; however, they are often vulnerable to changes in illumination or weak scene texture. At the same time, neural network-based methods are more robust under challenging conditions [8, 19]. Still, their specific nature of point localization and the statistical properties of the inlier distribution require comparative analysis [9, 18]. Evaluating the influence of the local feature type on geometric reliability and monocular visual odometry errors remains an important task, since different methods form keypoint representations in different ways, which in turn affects the stability of the navigation system.

The goal of this work is to analyze keypoint detection and description methods for monocular visual odometry using a comprehensive quality assessment approach. For this purpose, an integral index is employed that combines the geometric consistency of correspondences, indicators of motion observability and camera pose recovery stability, and the structural-spatial properties of consistent correspondences.

MATERIALS AND METHODS

For the experimental validation of the proposed approach, the EuRoC MAV dataset [20] was used. It includes scenes with varying levels of motion complexity, spatial structure, and visual conditions, enabling the evaluation of local features across easy, medium, and

difficult monocular visual odometry scenarios. For the analysis, images from the cam0 camera were processed. Two groups of sequences were considered:

- Machine Hall - large industrial indoor environments with non-uniform illumination and a considerable number of repetitive structures. Three difficulty levels were investigated: MH_01_easy, MH_03_medium, and MH_05_difficult. In this group, the difficulty increases from relatively smooth motion and more favorable illumination conditions to faster movement under dim lighting, which may be accompanied by stronger motion blur and reduced feature contrast [20].
- Vicon Room - a confined indoor environment with highly accurate ground-truth motion. Although this setting is more controlled, the increasing difficulty in this group is mainly associated with camera dynamics and illumination conditions. Three sequences with different difficulty levels were used: V1_01_easy, V1_02_medium, and V1_03_difficult. The transition to more challenging sequences is characterized by sharper maneuvers, faster viewpoint changes, and temporary loss of textured objects from the field of view, which makes tracking and motion estimation more difficult [20].

To compare the quality of local features in visual odometry, the classical methods ORB, AKAZE, KAZE, SIFT, SURF, and BRISK, as well as the deep learning-based method SuperPoint, were investigated. In this work, an off-the-shelf SuperPoint implementation from the LightGlue library [21] was used without any additional model retraining.

The analysis was performed for different numbers of selected keypoints, ranging from 200 to 1000 in increments of 200. Examples of the images used in the experiments are shown in Fig. 1.

To reduce the effect of local keypoint concentration in the most contrast-rich image regions, a spatially guided selection strategy based on a regular 8×5 grid was applied. This approach provides more balanced image coverage, since even with a large number of correct correspondences, their spatial clustering may degrade the stability of scene geometry estimation, reduce reconstruction reliability, and lead to instability of the recovered motion [22, 23].

Motion estimation was performed in a calibrated monocular setting based on the essential matrix [24]. For each pair of frames with a stride = 2, a set of keypoint correspondences was established using Brute-Force matching, followed by filtering according to the Lowe ratio criterion with a threshold of 0.75 [12]. Before geometric estimation, point

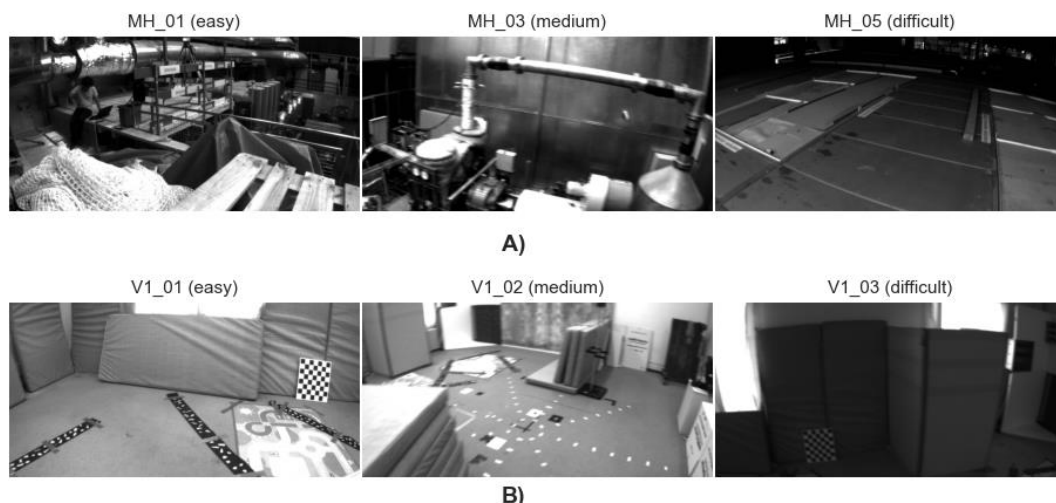


Fig. 1. Examples of images from the EuRoC MAV dataset used in the experiments: A) frames from the Machine Hall group (MH_01_easy, MH_03_medium, MH_05_difficult); B) frames from the Vicon Room group (V1_01_easy, V1_02_medium, V1_03_difficult).

coordinates were corrected using the camera calibration parameters to remove distortion effects, thereby reducing systematic errors in correspondence localization [24].

The essential matrix was estimated using USAC_FAST with a threshold of 1.5 px, a confidence level of 0.999, and a maximum of 5000 iterations [25, 26]. To ensure robust estimation, additional conditions were imposed: at least 50 initial correspondences and at least 30 inliers after geometric verification. The relative pose between frames was recovered using the recoverPose procedure, which selects the physically consistent decomposition of the essential matrix and returns a mask of points with positive depth. To reduce the influence of weakly observable configurations, an additional minimum parallax check of 5 px was applied in the freeze_translation mode, in which the translation estimate was not updated when the inter-frame displacement was insufficient [7, 27].

The quality of the recovered trajectory was evaluated using the evo package, a standard tool for trajectory analysis in robotics, SLAM, and visual odometry. For external accuracy assessment, the APE (Absolute Pose Error) and RPE (Relative Pose Error) metrics were used [28, 29].

The APE metric was used to assess the global trajectory error after alignment of the estimated and reference trajectories. In this work, Sim(3) alignment was employed, that is, alignment accounting for rotation, translation, and a single global scale factor, which is particularly important for monocular visual odometry, where the absolute scale cannot be determined uniquely without additional sensors. Thus, APE characterizes the accumulated drift and the overall deviation of the estimated trajectory from the reference in a global sense [29].

The RPE metric was used to assess the local error of relative motion over a fixed temporal baseline. Unlike APE, it reflects the stability of short-term estimation and is less sensitive to long-term error accumulation. To ensure comparability across all runs, the analysis was performed over a consistent temporal interval of 0.30 s, which was converted into the corresponding number of frames, accounting for the sampling stride [28]. At a frame rate of 20 Hz, this corresponds to a baseline interval of approximately 6 frames.

Since the translational RPE value in metric units strongly depends on the amplitude of inter-frame displacement, this work additionally considered a normalized form of the relative error:

$$RPE_{norm} = \frac{RPE_t}{\varepsilon + d_{med}}, \quad (1)$$

where RPE_t denotes the aggregated translational relative error, d is the median inter-frame displacement of the corresponding points in pixels, and ε is a small positive stabilizing term. Such normalization reduces the influence of varying motion intensity across individual fragments and enables a more meaningful comparison of cases with similar absolute error but different motion observability. In addition, for the analysis of rotational stability, tail-risk indicators derived from the angular RPE were used, in particular rotation error percentiles and the proportions of large deviations.

The internal quality assessment of local features in the VO task is based on a combination of geometric and spatial-structural components.

The basic geometric characteristic of a correspondence set is the proportion of inliers among all verified matches. It is calculated as the ratio of inliers consistent with the estimated geometric model to the total number of correspondences [30]. However, the inlier ratio alone does not indicate how closely these correspondences match the estimated model. Two correspondence sets may have similar inlier ratio values but differ significantly in residual geometric error. Therefore, a modified indicator is used in this work:

$$G_{ratio} = \frac{N_{int}}{N_{match}} \exp(-\bar{e}), \quad (2)$$

$$\bar{e} = \frac{\text{med}(e_{\text{Sampson}}^2)}{\tau^2}, \quad (3)$$

where $\text{med}(e_{\text{Sampson}}^2)$ is the median squared Sampson error for the inliers [27], and τ is the geometric threshold used in the USAC verification procedure. High values of Gratio correspond to cases in which the correspondence set simultaneously exhibits a high inlier ratio and a low residual geometric error.

For reliable motion recovery, geometric consistency of correspondences alone is not sufficient, since the quality of translation estimation strongly depends on the magnitude of inter-frame point displacement. When the displacement is small, the motion becomes less informative, and the recovery of spatial structure becomes less stable [7, 30]. Therefore, an additional parallax-based indicator is used:

$$P_{\text{score}} = 1 - \exp\left(-\frac{d_{\text{med}}}{p_0}\right), \quad (4)$$

where d_{med} is the median inter-frame displacement of corresponding points in pixels, and p_0 is a reference parameter. This function increases with parallax: for small displacement, P_{score} it approaches zero, whereas for sufficiently large displacement, it approaches one. Therefore, the indicator P_{score} characterizes motion observability and the suitability of the current configuration for stable estimation of translation [27, 30].

The second characteristic is the cheirality indicator:

$$\text{Cheirality}_{\text{ratio}} = \frac{N_{\text{chier}}}{N_{\text{inl}}}, \quad (5)$$

where N_{chier} is the number of inliers for which positive depth is obtained after pose recovery, and N_{inl} is the total number of inliers. This indicator follows from the cheirality check: a physically valid solution is the one for which the reconstructed 3D points lie in front of both cameras [24, 27]. Therefore, $\text{Cheirality}_{\text{ratio}}$ characterizes the physical validity of the recovered configuration. High values of $\text{Cheirality}_{\text{ratio}}$ correspond to cases in which most inliers support a correct spatial solution, whereas a decrease in this indicator may indicate estimation instability or a weak 3D interpretation [30].

Taking these factors into account, the final geometric component for the visual odometry task was defined as follows:

$$G_{VO} = G_{\text{ratio}} \cdot P_{\text{score}} \cdot \text{Cheirality}_{\text{ratio}}. \quad (6)$$

This form accounts for the fact that, for reliable motion estimation, a high inlier ratio alone isn't enough. Adequate inter-frame displacement and the physical correctness of the recovered spatial structure are also essential.

The spatial-structural component is based on three metrics proposed in the previous study [6]: CUI, SCS, and RI. The Coverage Uniformity Index (CUI) evaluates the uniformity of image-plane coverage. The Scene Consistency Score (SCS) measures the consistency between the scene structure and the structural composition of the point set.

To assess the local spatial redundancy of inliers, a modified Redundancy Index (RI) was used. Unlike the basic version, which evaluates redundancy using a fixed neighborhood threshold [6], this work defines the neighborhood radius relative to the image diagonal. Such normalization reduces the metric's dependence on image resolution and enables a more meaningful comparison of configurations with different point counts. For

each point, the number of neighbors within a radius $r = \gamma D$ is counted, where D is the image diagonal. The obtained local neighbor count is then compared with the expected density level for the current point set [31-33]. The final RI value is defined as the average over all points. It takes values in the range $[0,1]$: low values correspond to a more uniform distribution, whereas high values indicate local clustering and point redundancy.

This modification follows from the fact that, for local feature analysis tasks, it is important to consider not only the total number of points but also the nature of their spatial arrangement. Previous studies [6] have shown that the spatial distribution of keypoints is an independent quality characteristic, since excessive local concentration reduces scene coverage even when a large number of features is available [33]. In this work, this idea is further developed through diagonal normalization of the radius and explicit consideration of the expected local density, allowing the RI indicator to better reflect redundancy itself rather than the absolute point density.

For the integral description of spatial-structural quality, the following component was used:

$$S_{VO} = \frac{1}{3}(CUI + (1 - RI)^2 + SCS), \quad (7)$$

where the quadratic term $(1 - RI)^2$ increases the penalty for local redundancy. The use of $(1 - RI)^2$, rather than the linear form $(1 - RI)$, makes it possible to distinguish more clearly between point sets with moderate and high clustering.

The final integral indicator for the monocular visual odometry task is defined as:

$$Q_{VO} = 0.62 \cdot G_{VO} + 0.38 \cdot S_{VO}, \quad (8)$$

where the geometric component is assigned a higher weight, since it is the primary factor determining the correctness of motion recovery, whereas the spatial-structural component explains the stability of estimation and robustness to local deformations and non-uniform scene structure.

Since visual odometry is a stepwise procedure and may contain failed steps, an additional penalty was introduced for the integral score based on the frequency of unsuccessful steps. Let $f \in [0,1]$ denote the proportion of failed steps in a run (*fail rate*). Then the penalized score is defined as:

$$Q_{pen} = Q_{VO} \cdot (1 - f). \quad (9)$$

Thus, the proposed methodology combines two levels of evaluation. The external level characterizes the actual trajectory accuracy as measured by APE and RPE. In contrast, the internal level explains this behavior through the geometric reliability of correspondences, motion observability, and the spatial-structural balance of inliers. Such a design makes it possible not only to rank the methods, but also to interpret the reasons for their behavior under monocular visual odometry conditions.

RESULTS AND DISCUSSION

In all the presented plots, the values of the analyzed metrics are reported as the median for each experiment. This form of representation was chosen because the distributions of indicators in the monocular visual odometry task may contain outliers and locally unstable estimates. In contrast, the median provides a more robust and representative summary.

Fig. 2 shows the dependence of the geometric consistency index, G_{ratio} , on the number of keypoints, N . In general, AKAZE demonstrates the highest metric values in most sequences, while SURF is usually second or close to the best result. For AKAZE, the G_{ratio} values remain consistently high: in the MH_01_easy, V1_01_easy, and MH_05_difficult sequences, they are close to 0.91-0.93, while in the more challenging V1_03_difficult sequence, they remain at approximately 0.82.

SIFT and KAZE also show fairly high results; however, they generally remain below AKAZE and SURF. For example, at $N=1000$ in V1_01_easy, the G_{ratio} value is 0.929 for AKAZE, 0.905 for SURF, 0.903 for SIFT, and 0.884 for KAZE.

ORB and BRISK provide lower-level geometric consistency. In most sequences, ORB yields the lowest values among the classical methods, while BRISK usually outperforms ORB but remains noticeably inferior to AKAZE, SURF, and SIFT. For example, in MH_03_medium at $N=1000$, the value of G_{ratio} for BRISK equals 0.784, whereas for AKAZE it reaches 0.905, corresponding to a difference of approximately 15.4%.

SuperPoint demonstrates the most pronounced negative trend: as N increases, its G_{ratio} value systematically decreases in all sequences. In particular, in MH_01_easy, the indicator decreases from 0.776 to 0.693 (approximately 10.7%); in V1_01_easy, from 0.691 to 0.609 (11.9%); and in V1_03_difficult, from 0.576 to 0.487 (15.5%). This indicates that increasing the number of points for this method does not improve geometric consistency but, on the contrary, worsens it. However, in some cases, SuperPoint was more effective than ORB. It can also be observed that SuperPoint's performance depends on sequence difficulty. For classical methods, the dependence of performance on sequence difficulty is less evident.

Thus, according to the G_{ratio} metric, the most stable results were achieved by AKAZE and SURF, whereas SuperPoint proved most sensitive to increasing the number of points. Overall, the obtained results confirm that increasing the number of keypoints by itself does not guarantee an improvement in the geometric quality of correspondences.

Fig. 3 presents the dependence of the parallax indicator P_{score} on the number of keypoints. Unlike the G_{ratio} metric, the separation between methods is considerably weaker in this case, and in some sequences, the values from different detectors remain very close. This indicates that, within the present experiment, the parallax indicator should rather be regarded as an indicator of overall motion observability than as an independent criterion for clear method ranking.

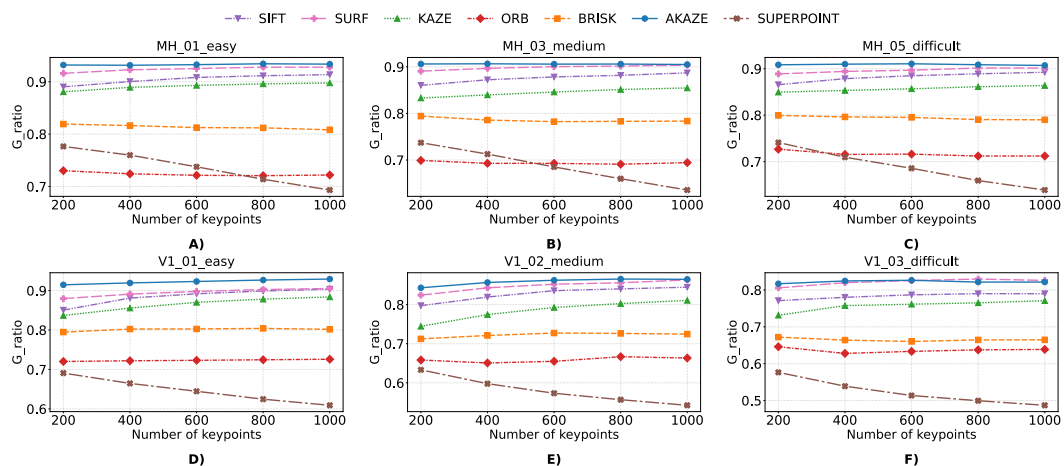


Fig. 2. Median values of G_{ratio} vs. the number of keypoints for different detectors in the EuRoC sequences: (A) MH_01_easy, (B) MH_03_medium, (C) MH_05_difficult, (D) V1_01_easy, (E) V1_02_medium, (F) V1_03_difficult.

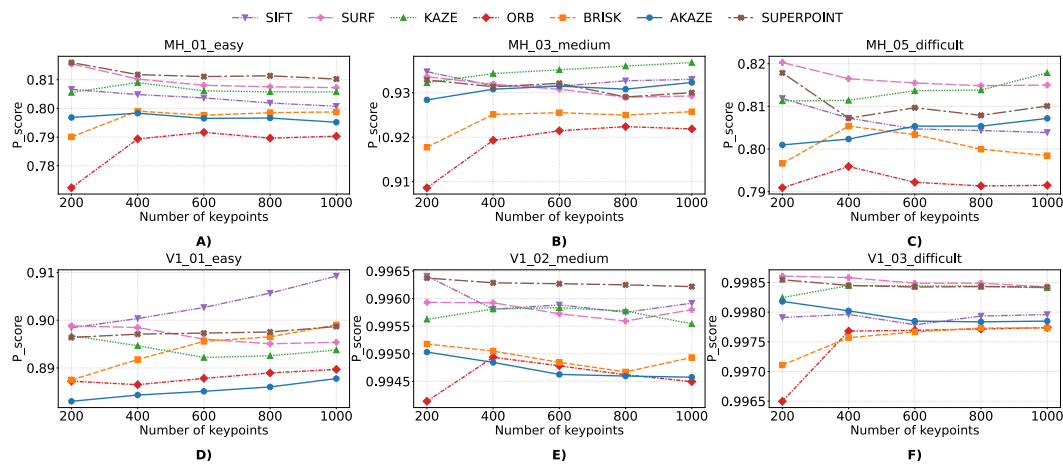


Fig. 3. Median values of P_{score} vs. the number of keypoints for different detectors in the EuRoC sequences: (A) MH_01_easy, (B) MH_03_medium, (C) MH_05_difficult, (D) V1_01_easy, (E) V1_02_medium, (F) V1_03_difficult.

The most noticeable separation between detectors is observed in the MH_01_easy, MH_03_medium, and MH_05_difficult sequences, where certain methods produce higher or lower P_{score} values, although the overall difference remains moderate. In these cases, ORB is more often found in the lower part of the group, whereas SIFT, SURF, KAZE, or SuperPoint yield slightly better results depending on the sequence. At the same time, as the number of points increases, most methods exhibit either a slight improvement or a rapid transition to a saturation.

A different pattern is observed for V1_02_medium and V1_03_difficult, where almost all methods yield very high, closely matched P_{score} values. Under such conditions, this metric does not effectively separate the detectors, indicating its limited discriminative ability in sequences where motion observability is generally favorable for most methods.

Thus, P_{score} should primarily be considered an auxiliary geometric indicator that characterizes the conditions for reliable motion estimation but is not by itself sufficient for the final comparison of detectors. For this reason, its interpretation is most meaningful when combined with G_{ratio} , the proportion of correct cheirality, and the integral quality indicator.

Fig. 4 illustrates the dependence of the correct cheirality ratio on the number of keypoints. The obtained results indicate that this metric depends more strongly on the sequence type than P_{score} : across different scenes, not only the absolute levels of the values change, but also the relative positions of the methods. This means that a complex interaction between detector properties and the characteristics of a particular sequence determines the physical validity of the recovered spatial configuration.

Unlike P_{score} , which in some cases hardly separates the methods, the correct cheirality ratio more often reveals pronounced differences between detectors. At the same time, no universal leader is observed: in different sequences, the best results are demonstrated by SURF, SIFT, KAZE, AKAZE, or SuperPoint. This pattern is indicative, as it shows that even under similar motion observability conditions, different methods may differ substantially in their ability to form correspondences suitable for physically valid scene geometry recovery.

Therefore, this metric is an important complement to the other geometric indicators, as it allows evaluation not only of the consistency of correspondences but also of their suitability for correct spatial reconstruction. At the same time, its results further confirm the relevance of a comprehensive analysis, because it does not, on its own, provide a complete ranking of the methods.

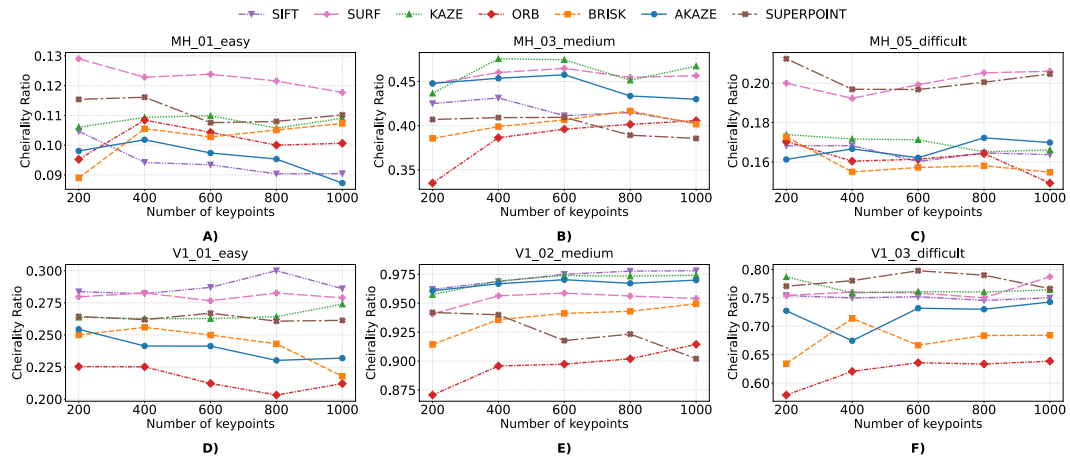


Fig. 4. Median values of Cheirality ratio vs. the number of keypoints for different detectors in the EuRoC sequences: (A) MH_01_easy, (B) MH_03_medium, (C) MH_05_difficult, (D) V1_01_easy, (E) V1_02_medium, (F) V1_03_difficult.

Fig. 5 shows the dependence of the inlier coverage uniformity indicator CUI on the number of keypoints. A common trend is observed across all investigated sequences: as the number of keypoints increases, the metric value either rises or gradually reaches a saturation regime. This indicates that increasing the number of keypoints generally contributes to a more complete and more uniform spatial coverage of the scene by correct correspondences.

SuperPoint primarily demonstrates the highest CUI values in all six sequences, while SURF is usually the second-best method. For example, in MH_01_easy at N=1000, the CUI value for SuperPoint is approximately 0.806, whereas for SURF it is 0.778, for AKAZE 0.706, and for ORB 0.598. Thus, the advantage of SuperPoint over ORB is approximately 34.8%, and over AKAZE about 14.2%. A similar pattern is also observed for MH_03_medium, where at N=1000 SuperPoint reaches approximately 0.756, whereas ORB achieves only 0.555.

In V1_01_easy and V1_03_difficult, the absolute metric values are lower across all methods; however, the relative separation between them remains. For example, in

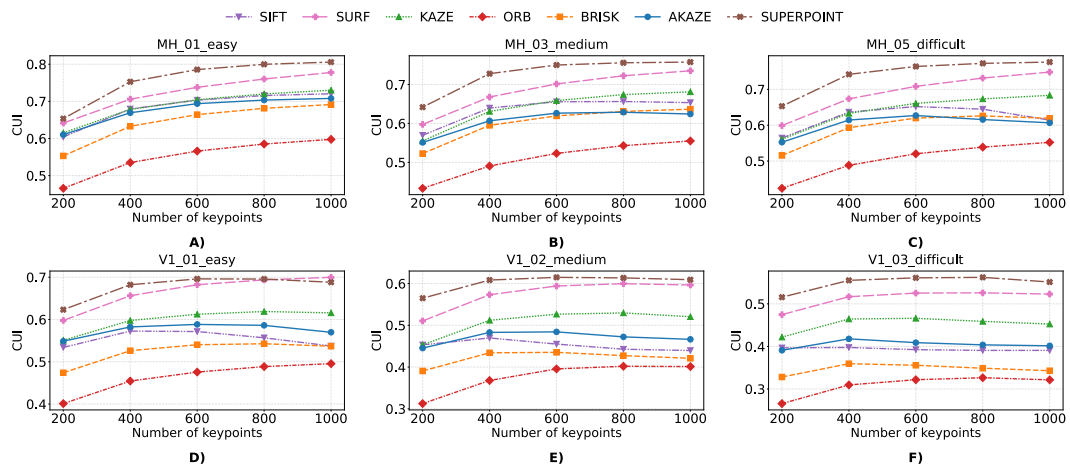


Fig. 5. Median values of CUI vs. the number of keypoints for different detectors in the EuRoC sequences: (A) MH_01_easy, (B) MH_03_medium, (C) MH_05_difficult, (D) V1_01_easy, (E) V1_02_medium, (F) V1_03_difficult.

V1_03_difficult at N=1000, the value for SuperPoint is approximately 0.551, for SURF 0.523, whereas for ORB it is only 0.323. This means that SuperPoint exceeds ORB by approximately 1.7 times. Thus, even under more challenging conditions, SuperPoint and SURF remain in the upper part of the group, while ORB consistently demonstrates the lowest or near-lowest values.

AKAZE, KAZE, and SIFT form a middle group of methods, characterized by either a moderate increase in CUI or stabilization after N=400-600. BRISK also improves coverage as the number of points increases; however, in most cases, it remains inferior to this group. Therefore, the CUI metric effectively reflects the spatial completeness of scene coverage by inliers and shows that SuperPoint provides the most uniform distribution of correct correspondences over the image. At the same time, the absolute level of the metric also depends on the particular sequence, so its values should primarily be interpreted within each scene separately.

Fig. 6 presents the dependence of the normalized indicator of local inlier redundancy, RI, on the number of keypoints, with lower values indicating better results. A common trend is observed across all investigated sequences: as N increases, the RI value rises, that is, local redundancy becomes stronger. This means that as the number of points increases, inliers increasingly concentrate in individual local regions, even as the overall scene coverage improves.

SuperPoint consistently shows the lowest RI values across all six sequences. For example, in MH_01_easy at N=1000, its value is approximately 0.303, whereas for ORB it is 0.835, for BRISK 0.659, and for AKAZE 0.600. Thus, relative to ORB, the value for SuperPoint is approximately 2.8 times lower, and relative to AKAZE, almost 2 times lower. A similar pattern is observed in other sequences as well, in particular in MH_03_medium, MH_05_difficult, and V1_03_difficult.

The worst results are usually shown by ORB, which, in all cases, has the highest or near-highest metric values. For example, in V1_01_easy at N=1000, the ORB value is approximately 0.886, whereas for SIFT it is 0.714, for SURF 0.570, and for SuperPoint 0.455. This indicates a substantially higher local concentration of inliers for ORB compared with the other methods.

In most sequences, SIFT and SURF form a relatively favorable group, with lower redundancy than AKAZE, KAZE, and BRISK. For SIFT and SURF, the increase in RI when moving from N=200 to N=400 is noticeable, but afterwards the changes become smaller or

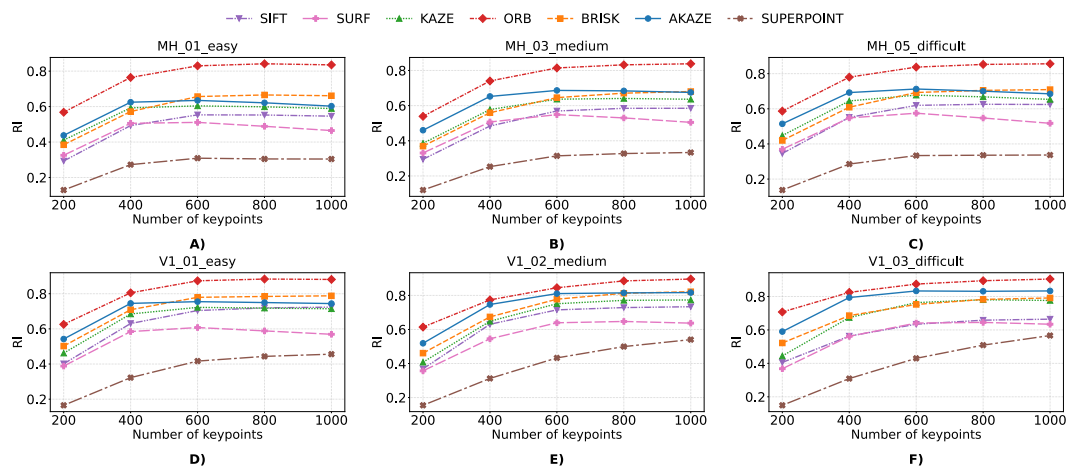


Fig. 6. Median values of RI vs. the number of keypoints for different detectors in the EuRoC sequences: (A) MH_01_easy, (B) MH_03_medium, (C) MH_05_difficult, (D) V1_01_easy, (E) V1_02_medium, (F) V1_03_difficult.

reach a saturation regime. In contrast, for ORB, BRISK, AKAZE, and KAZE, the metric values remain at a higher level in almost all sequences.

Thus, the RI metric complements CUI well: if CUI reflects the completeness of scene coverage by inliers, then RI shows how strongly these inliers are locally clustered. In this respect, SuperPoint proved to be the best, whereas ORB was the weakest. At the same time, the absolute level of the metric partly depends on the particular sequence, so its values should primarily be interpreted within each scene separately.

Fig. 7 illustrates the dependence of the inlier structural consistency indicator SCS on the number of keypoints N . Unlike CUI and RI, a different group of leaders can be clearly identified here: in all investigated sequences, the highest SCS values are consistently demonstrated by SURF, while KAZE and AKAZE usually form a second group of methods that also achieve high results.

For SURF, the metric values are the highest in all six sequences and, in most cases, either increase with the number of keypoints or remain close to their maximum levels. For example, in MH_01_easy at $N=1000$, the SCS value for SURF is approximately 0.922, whereas for AKAZE it is 0.779, for SIFT 0.494, and for ORB 0.297. Thus, the advantage of SURF over AKAZE is approximately 18.4%, while over ORB it is about threefold. A similar pattern is also observed for V1_03_difficult, where at $N=1000$ SURF exceeds AKAZE by approximately 18.6% and ORB by about 2.5 times.

KAZE and AKAZE are part of a group of methods with high values, though lower than SURF. In most sequences, they are characterized by either only slight changes or a moderate decrease in SCS as the number of keypoints increases. SIFT usually occupies an intermediate position between this group and the methods with low structural consistency. In most cases, the lowest metric values are observed for ORB, BRISK, and SuperPoint, with ORB often being the weakest method.

Thus, the SCS metric effectively reflects the correspondence between the spatial distribution of inliers and the scene's structure. While SuperPoint showed the best results in terms of CUI and RI, SURF is the clear leader in terms of SCS. This confirms that distinct spatial-structural metrics characterize different aspects of inlier quality and should therefore be considered jointly. At the same time, the absolute level of the metric also depends on the particular sequence, so its values should primarily be interpreted within each scene.

Fig. 8 presents the dependence of the penalized integral quality indicator Q_{pen} on the number of keypoints. Unlike the individual geometric and spatial-structural metrics, this

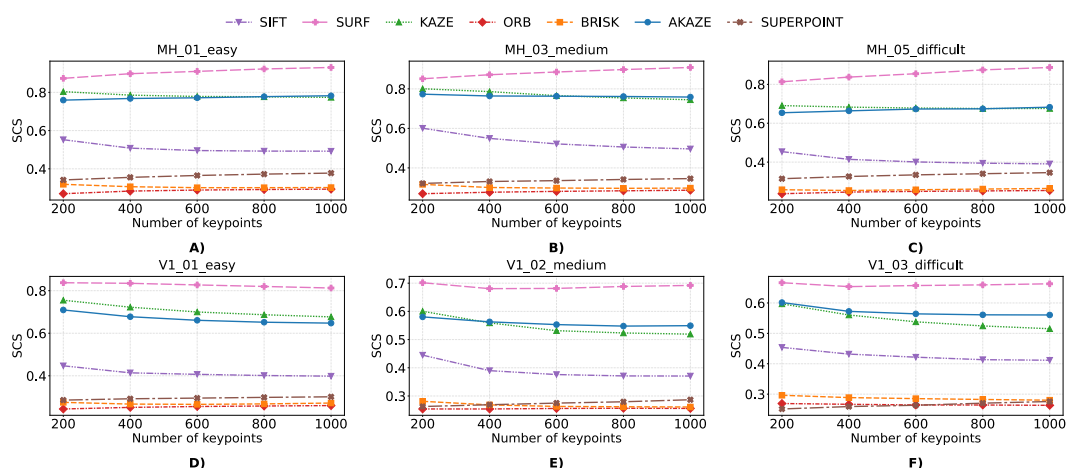


Fig. 7. Median values of SCS vs. the number of keypoints for different detectors in the EuRoC sequences: (A) MH_01_easy, (B) MH_03_medium, (C) MH_05_difficult, (D) V1_01_easy, (E) V1_02_medium, (F) V1_03_difficult.

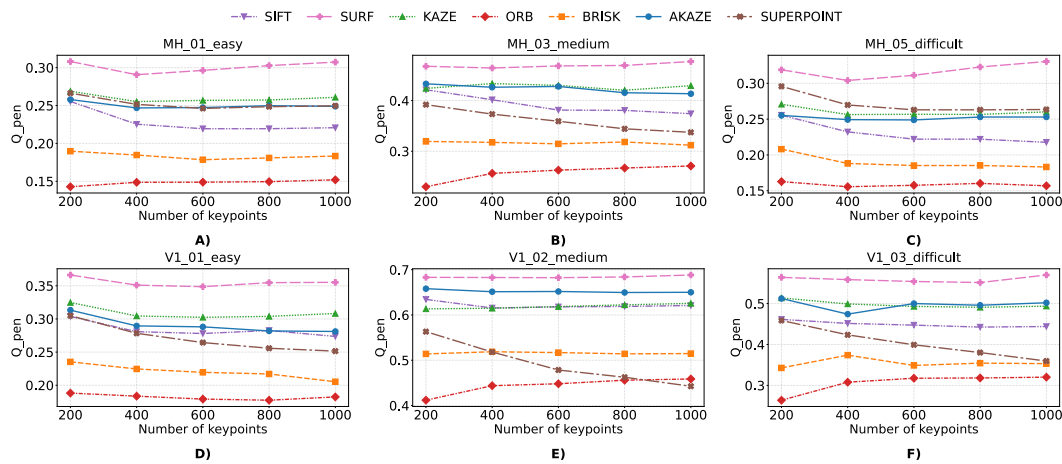


Fig. 8. Median values of Q_{pen} vs. the number of keypoints for different detectors in the EuRoC sequences: (A) MH_01_easy, (B) MH_03_medium, (C) MH_05_difficult, (D) V1_01_easy, (E) V1_02_medium, (F) V1_03_difficult.

indicator summarizes their combined contribution and therefore provides a more comprehensive view of detector suitability in the monocular visual odometry task.

In most sequences, SURF consistently achieves the highest Q_{pen} values, indicating the best balance among geometric quality, spatial coverage, Q_{pen} consistency, and penalty components. For example, at $N=1000$, the Q_{pen} value for SURF is approximately 0.308 in MH_01_easy, 0.355 in V1_01_easy, and 0.567 in V1_03_difficult.

A second group of methods is usually composed of AKAZE and KAZE, which, in most cases, produce similar Q_{pen} values. For example, in V1_02_medium at $N=1000$, AKAZE reaches approximately 0.649, while KAZE and SIFT yield close but slightly lower values. A similar pattern is observed in V1_03_difficult, where AKAZE and KAZE also form the upper group after SURF, consistent with their high geometric metric values.

SuperPoint demonstrates non-uniform behavior. In some sequences, its Q_{pen} values are relatively high at small N , but then tend to decrease or stagnate. For example, in V1_03_difficult, the value for SuperPoint decreases from 0.459 at $N=200$ to 0.360 at $N=1000$, that is, by approximately 21.6%. This is in good agreement with the previously observed decrease in G_{ratio} for this method as the number of points increases.

ORB demonstrates the lowest Q_{pen} values in almost all sequences, while BRISK usually occupies an intermediate position between ORB and the group of stronger methods. For example, in MH_01_easy and V1_01_easy at $N=1000$, SURF exceeds ORB by approximately twofold. This confirms that ORB's weaker performance is observed not only in individual components but also in the overall assessment.

Thus, the Q_{pen} metric consistently summarizes the previous observations: SURF proved to be the most stable leader in terms of overall quality, AKAZE and KAZE formed a strong second group, whereas ORB and, to some extent, BRISK demonstrated lower overall suitability. SuperPoint, despite strong performance on some individual spatial metrics, remains inferior to the leaders due to its weaker geometric component, which becomes especially evident at large N . At the same time, the absolute level of the integral indicator also depends on the particular sequence, so its values should primarily be interpreted within each scene.

To analyze the relationship between the spatial-geometric indicators and odometry errors, Spearman rank correlation matrices were constructed separately for each EuRoC sequence, as shown in Fig. 9. This approach made it possible to reveal that both the strength and even the sign of the correlations may vary depending on the particular scene,

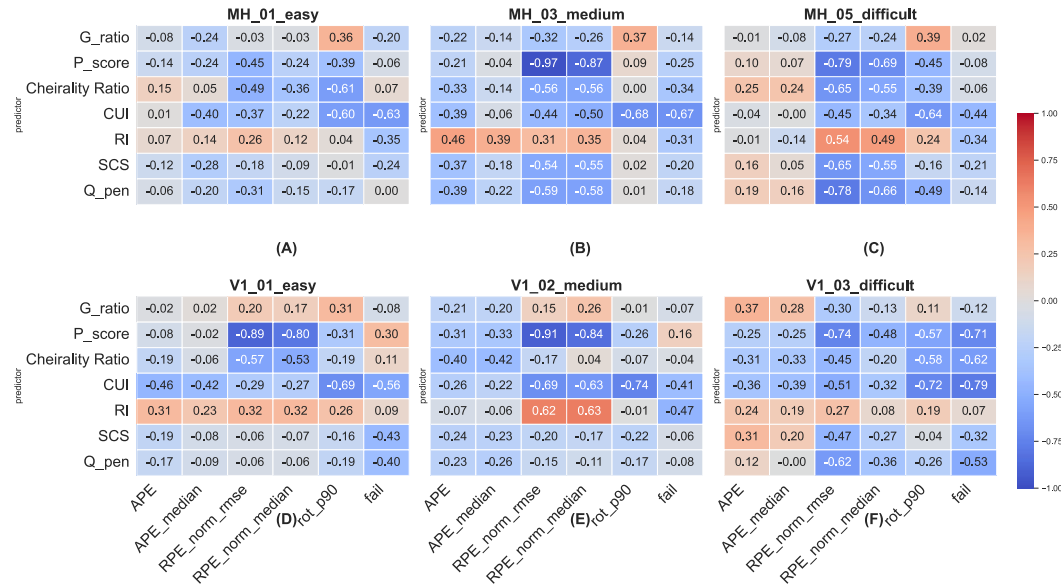


Fig. 9. Spearman rank correlation matrices between the spatial-geometric predictors and odometry error metrics for individual EuRoC sequences: (A) MH_01_easy, (B) MH_03_medium, (C) MH_05_difficult, (D) V1_01_easy, (E) V1_02_medium, (F) V1_03_difficult.

motion pattern, and observation conditions; that is, the informativeness of the metrics has a pronounced sequence-dependent character.

For the Machine Hall sequences, the most consistent relationship with odometry quality is demonstrated by the inlier coverage uniformity indicator RI, for which negative correlations are predominantly observed with translational errors, rotational tail-risk indicators, and failure rate. For example, for MH_01_easy, the correlation between the inlier coverage uniformity indicator and the failure rate is $\rho = -0.63$, and for MH_03_medium, it is $\rho = -0.67$. The local redundancy index RI, in contrast, shows mostly positive correlations with the error measures, consistent with the negative effect of local inlier clustering. For some sequences in this group, a pronounced relationship is also observed between the parallax indicator P_{score} and the normalized relative translational error RPE_{norm} ; in particular, for MH_03_medium, the correlation is $\rho = -0.97$.

For the Vicon Room sequences, the strongest and most stable relationship with translational errors is most often observed with the parallax indicator, underscoring the importance of motion observability under these conditions. For example, for V1_02_medium, the correlation between the parallax indicator and the normalized relative translational error is $\rho = -0.91$, and for V1_03_difficult, it is $\rho = -0.74$. At the same time, the inlier coverage uniformity indicator remains substantially related to odometry stability; in particular, for V1_03_difficult, its correlation with the failure rate is $\rho = -0.79$. Other metrics, including the correct cheirality ratio, the inlier structural consistency indicator SCS, and the penalized integral quality indicator Q_{pen} , exhibit a less homogeneous pattern, indicating stronger dependence on the specific scenario.

Thus, the matrices for the individual sequences confirm that no single partial metric serves as a universal predictor of trajectory error. The most informative indicators for explaining changes in odometry quality were the inlier coverage uniformity indicator, the parallax indicator, and the normalized local redundancy indicator. The correlations themselves should be interpreted as indicators of a monotonic relationship rather than as direct evidence of causality.

The obtained results showed that different groups of metrics emphasize different aspects of the suitability of local features for the monocular visual odometry task. This is

especially evident in cases where a method exhibits strong spatial characteristics but lacks the same level of geometric reliability. In particular, SuperPoint proved to be one of the best methods in terms of coverage, uniformity and local redundancy. Yet its advantage was not preserved in the geometric metrics and the penalized integral indicator. This result indicates that good spatial coverage alone does not guarantee the best suitability for motion estimation.

In contrast, SURF was not always the unconditional leader across all individual partial metrics, yet it demonstrated the most stable balance among geometric consistency, structural correspondence, and integral assessment. This provides grounds for considering the balance of characteristics to be a more important property for practical use in VO than the maximization of any single criterion. In this sense, the results confirm that evaluating local features in visual odometry tasks should be based not on a single “best” indicator but on a set of complementary criteria.

At the same time, the correlation analysis showed that the informativeness of the partial metrics regarding odometry quality varies with the sequence type. In some scenes, uniform inlier coverage is more important, whereas in others, motion-observability characteristics are more informative. This means that there is no simple universal relationship between local feature quality and the final trajectory error. Such heterogeneity further confirms the relevance of a comprehensive evaluation approach that considers different metrics jointly. At the same time, the integral indicator is used as a means of generalized ranking rather than as a direct replacement for trajectory-based metrics.

CONCLUSION

This study presents a comprehensive investigation of local feature quality in monocular visual odometry using a combination of geometric, spatial-structural, and integral indicators. The performed analysis showed that no single metric is sufficient for a complete characterization of detector suitability for the VO task, since different indicators reflect different aspects of correspondence quality. According to geometric metrics, the best results were achieved by AKAZE, SURF, and partially by SIFT. In contrast, in terms of spatial-structural indicators, SuperPoint showed an advantage in inlier coverage uniformity and local redundancy, while SURF was superior in structural consistency.

It was shown that, as the number of keypoints increases, the results for most methods initially improve, but then often reach a saturation regime and in some cases even deteriorate. This indicates that the practical effectiveness of a detector is determined not only by the number of detected features, but also by their informativeness, spatial distribution, and ability to form geometrically valid correspondences. It was established that SURF is the most balanced method in the conducted experiments, as it consistently ranks among the leaders across different partial criteria and attains the highest values of the penalized integral quality indicator. AKAZE and KAZE formed a strong second group, whereas ORB showed the weakest results in most cases, both for the partial metrics and for the integral assessment.

The correlation analysis showed that the relationship between the partial metrics and the odometry quality indicators varies across sequences. For the Machine Hall group, the most informative indicator was the uniformity of inlier coverage. In contrast, for the Vicon Room group, the strongest relationship with translational errors was demonstrated by the parallax indicator. This confirms that no single metric can be regarded as a universal indicator of VO quality for all scene types, and that the integral approach is an appropriate means of multicriteria generalization and ranking.

Thus, the proposed approach to the spatial-geometric evaluation of local features enables not only comparison of methods but also explanation of their respective strengths and weaknesses in the context of visual odometry. The practical value of the obtained results lies in the ability to make a well-grounded detector choice based on the requirements of a particular application and the observation conditions.

Prospects for further research are associated with the use of other types of datasets, in particular those with more pronounced variations in illumination, texture, scene dynamics, and motion scale. It is also advisable to further investigate the behavior of the proposed indicators in combination with modern neural network-based detectors and descriptors, and to verify their suitability not only for monocular visual odometry but also for a broader range of tasks, including SLAM. A separate direction for future work will be to improve the integral indicator by adaptively adjusting the weighting coefficients based on the scene type.

ACKNOWLEDGMENTS AND FUNDING SOURCES

The authors received no financial support for the research, writing, and/or publication of this article.

CONFLICT OF INTEREST

The authors declare that the research was conducted in the absence of any.

AUTHOR CONTRIBUTIONS

Conceptualization, [A.F., Yu.F.]; methodology, [A.F., Yu.F.]; validation, [A.F., Yu.F.]; writing – original draft preparation, [A.F.]; writing – review and editing, [A.F., Yu.F.]; supervision, [Yu.F.].

All authors have read and agreed to the published version of the manuscript.

REFERENCES

- [1] Herrera-Granda, E. P., Torres-Cantero, A., Haro, G., & Nieto, M. (2024). Monocular visual SLAM, visual odometry, and structure from motion methods applied to 3D reconstruction: A comprehensive survey. *Heliyon*, 10(18), e37356. <https://doi.org/10.1016/j.heliyon.2024.e37356>
- [2] Zhao, L., Li, Y., Wang, M., & Zhang, Y. (2025). PLL-VO: An efficient and robust visual odometry integrating point-line features and neural networks. *ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, X-G-2025, 1045-1052. <https://doi.org/10.5194/isprs-annals-X-G-2025-1045-2025>
- [3] Luan, S., Yang, C., Qin, X., Chen, D., & Sui, W. (2024). Towards robust visual odometry by motion blur recovery. *Frontiers in Signal Processing*, 4, 1417363. <https://doi.org/10.3389/frsip.2024.1417363>
- [4] Budzan, S., Wyżgolik, R., & Lysko, M. (2025). Performance analysis of keypoints detection and description algorithms for stereo vision based odometry. *Sensors*, 25(19), 6129. <https://doi.org/10.3390/s25196129>
- [5] Huang, Q., Guo, X., Wang, Y., Sun, H., & Yang, L. (2024). A survey of feature matching methods. *IET Image Processing*, 18(6), 1385-1410. <https://doi.org/10.1049/ipr2.13032>
- [6] Fesiuk, A., & Furgala, Y. (2025). Comprehensive spatial-geometric evaluation of keypoint detectors. *Electronics and Information Technologies*, 32, 67-86. <https://doi.org/10.30970/eli.32.5>
- [7] Decker, P., Paulus, D., & Feldmann, T. (2008). *Dealing with degeneracy in essential matrix estimation*. In *2008 15th IEEE International Conference on Image Processing* (pp. 1964-1967). IEEE. <https://doi.org/10.1109/ICIP.2008.4712167>
- [8] Ma, J., Jiang, X., Jiang, J., Zhao, J., & Guo, X. (2021). Image matching from handcrafted to deep features: A survey. *International Journal of Computer Vision*, 129(1), 23-79. <https://doi.org/10.1007/s11263-020-01359-2>

- [9] Xu, S., Chen, S., Xu, R., Wang, C., Lu, P., & Guo, L. (2024). Local feature matching using deep learning: A survey. *Information Fusion*, 107, 102344. <https://doi.org/10.1016/j.inffus.2024.102344>
- [10] Nagy, A., Barsi, Á., & Takács, B. (2025). A comparative evaluation of classical and deep learning visual odometry configurations. *Engineering Proceedings*, 113(1), 16. <https://doi.org/10.3390/engproc2025113016>
- [11] Yu, J., Zhang, H., Liu, Q., & Chen, Z. (2025). Dynamic feature elimination-based visual–inertial odometry based on an optimized SuperPoint feature extractor. *Sensors*, 26(1), 52. <https://doi.org/10.3390/s26010052>
- [12] Lowe, D. G. (2004). Distinctive image features from scale-invariant keypoints. *International Journal of Computer Vision*, 60(2), 91–110. <https://doi.org/10.1023/B:VISI.0000029664.99615.94>
- [13] Bay, H., Ess, A., Tuytelaars, T., & Van Gool, L. (2008). Speeded-Up Robust Features (SURF). *Computer Vision and Image Understanding*, 110(3), 346–359. <https://doi.org/10.1016/j.cviu.2007.09.014>
- [14] Alcantarilla, P. F., Bartoli, A., & Davison, A. J. (2012). KAZE features. In *ECCV 2012* (LNCS 7577, pp. 214–227). https://doi.org/10.1007/978-3-642-33783-3_16
- [15] Alcantarilla, P. F., Nuevo, J., & Bartoli, A. (2013). Fast explicit diffusion for accelerated features in nonlinear scale spaces. In *BMVC 2013* (pp. 1–11). <https://doi.org/10.5244/C.27.13>
- [16] Rublee, E., Rabaud, V., Konolige, K., & Bradski, G. (2011). ORB: An efficient alternative to SIFT or SURF. In *ICCV 2011* (pp. 2564–2571). <https://doi.org/10.1109/ICCV.2011.6126544>
- [17] Leutenegger, S., Chli, M., & Siegwart, R. Y. (2011). BRISK: Binary robust invariant scalable keypoints. In *ICCV 2011* (pp. 2548–2555). <https://doi.org/10.1109/ICCV.2011.6126542>
- [18] DeTone, D., Malisiewicz, T., & Rabinovich, A. (2018). *SuperPoint: Self-supervised interest point detection and description*. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops* (pp. 224–236). <https://doi.org/10.48550/arXiv.1712.07629>
- [19] Budzan, S., Wyżgolik, R., & Lysko, M. (2025). Performance analysis of keypoints detection and description algorithms for stereo vision based odometry. *Sensors*, 25(19), 6129. <https://doi.org/10.3390/s25196129>
- [20] Burri, M., Nikolic, J., Gohl, P., Schneider, T., Rehder, J., Omari, S., Achtelik, M. W., & Siegwart, R. (2016). *The EuRoC micro aerial vehicle datasets*. *The International Journal of Robotics Research*, 35(10), 1157–1163. <https://doi.org/10.1177/0278364915620033>
- [21] Lindenberger, P., Sarlin, P.-E., & Pollefeys, M. (2023). *LightGlue: Local Feature Matching at Light Speed*. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (pp. 17681–17692). <https://doi.org/10.1109/ICCV51070.2023.01616>
- [22] Nannen, V., de Rezende, P. J., & Oliveira, M. M. (2013). *Grid-based spatial keypoint selection for real time visual odometry*. In *Proceedings of the International Conference on Computer Vision Theory and Applications (VISAPP)* (pp. 1–8). <https://doi.org/10.5220/0004270005860589>
- [23] Cho, H. M., Lee, S. H., & Kim, H. S. (2025). *Robust visual–inertial odometry via multi-scale deep feature extraction and adaptive keypoint selection*. *Applied Sciences*, 15(20), 10935. <https://doi.org/10.3390/app152010935>
- [24] Howse, Joseph, and Joe Minichino. “Learning OpenCV 4 Computer Vision with Python 3: Get to grips with tools, techniques, and algorithms for computer vision and machine learning.”, *Packt Publishing Ltd*, 2020.

- [25] Fesiuk, A., & Furgala, Y. (2025). Keypoint matches filtering in computer vision: Comparative analysis of RANSAC and USAC variants. *International Journal of Computing*, 24(2), 343–350. <https://doi.org/10.47839/ijc.24.2.4018>
- [26] M. Ivashechkin, D. Baráth, J. Matas, "USACv20: Robust Essential, Fundamental and Homography Matrix Estimation," 2021. <https://doi.org/10.48550/arXiv.2104.05044>
- [27] Nistér, D. (2004). *An efficient solution to the five-point relative pose problem*. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 26(6), 756-770. <https://doi.org/10.1109/TPAMI.2004.17>
- [28] Grupp, M. (2026). *evo: Python package for the evaluation of odometry and SLAM*. GitHub repository. <https://github.com/MichaelGrupp/evo>
- [29] Umeyama, S. (1991). *Least-squares estimation of transformation parameters between two point patterns*. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 13(4), 376-380. <https://doi.org/10.1109/34.88573>
- [30] Hartley, R., & Zisserman, A. (2004). *Multiple View Geometry in Computer Vision* (2nd ed.). Cambridge University Press.
- [31] Hossein-Nejad, Z., & Nasri, M. (2017). *RKEM: Redundant keypoint elimination method in image registration*. *IET Image Processing*, 11(4), 273-284. <https://doi.org/10.1049/iet-ipr.2016.0440>
- [32] Bailo, O., Rameau, F., Joo, K., Park, J., Bogdan, O., & Kweon, I. S. (2018). *Efficient adaptive non-maximal suppression algorithms for homogeneous spatial keypoint distribution*. *Pattern Recognition Letters*, 106, 53-60. <https://doi.org/10.1016/j.patrec.2018.02.020>
- [33] Gauglitz, S., Foschini, L., Turk, M., & Höllerer, T. (2011). *Efficiently selecting spatially distributed keypoints for visual tracking*. In *2011 18th IEEE International Conference on Image Processing* (pp. 1869-1872). IEEE. <https://doi.org/10.1109/ICIP.2011.6115832>

ПРОСТОРОВО-ГЕОМЕТРИЧНЕ ОЦІНЮВАННЯ ЛОКАЛЬНИХ ОЗНАК У МОНОКУЛЯРНІЙ ВІЗУАЛЬНІЙ ОДОМЕТРІЇ

Андрій Фесюк*  , Юрій Фургала  

Кафедра оптоелектроніки та інформаційних технологій
Львівський національний університет імені Івана Франка,
вул. Драгоманова 50, 79005 Львів, Україна

АНОТАЦІЯ

Вступ. Монокулярна візуальна одометрія (ВО) є важливим компонентом систем візуальної навігації, однак її точність залежить від якості локальних ознак і міжкадрових відповідностей. У задачі ВО важливими є не лише геометрична узгодженість, а й спостережуваність руху, фізична коректність відновленої конфігурації та просторово-структурні властивості локальних ознак. Метою роботи є комплексне оцінювання методів виявлення та опису особливих точок у монокулярній візуальній одометрії.

Матеріали та методи. Дослідження проведено на наборі даних EuRoC MAV. Проаналізовано методи ORB, BRISK, AKAZE, KAZE, SIFT, SURF та SuperPoint за кількості ключових точок від 200 до 1000. Оцінювання руху виконувалося на основі істотної матриці з використанням фільтра USAC_FAST, методу recoverPose, перевірки мінімального паралаксу та просторово керованого відбору ключових точок. Точність

відновленої траєкторії оцінювали за метриками APE та RPE. Для аналізу якості локальних ознак і відповідностей використовували геометричну складову, показник паралаксу, частку коректної хіральності, а також метрики рівномірності покриття особливими точками, локальної надлишковості та структурної узгодженості. Для узагальнення результатів застосовано інтегральний показник якості.

Результати. Геометричні метрики найчастіше виділяють алгоритми AKAZE та SURF, тоді як за просторовими характеристиками сильні позиції має метод SuperPoint. За показником структурної узгодженості відповідностей найкращі результати стабільно демонструє алгоритм SURF. Зі збільшенням кількості ключових точок для більшості методів спостерігається початкове покращення результатів із подальшим насиченням, а в деяких випадках – погіршення окремих характеристик. Найбільш збалансованим методом за сукупністю критеріїв виявився SURF, тоді як алгоритм ORB у більшості випадків продемонстрував найслабші результати. Кореляційний аналіз показав, що інформативність метрик залежить від типу послідовності.

Висновки. Запропонований підхід підтвердив доцільність багатокритеріального оцінювання локальних ознак у монокулярній візуальній одометрії. Показано, що жодна окрема метрика не є універсальною для всіх типів сцен, тоді як інтегральний показник дозволяє узагальнити різні аспекти якості та виконувати більш обґрунтоване ранжування методів.

Ключові слова: монокулярна візуальна одометрія, виявлення ключових точок, зіставлення зображень, оцінювання руху, глибоке навчання, нейронні мережі.

UDC 004.89

ENTROPY-GUIDED TRACKER SWITCHING METHOD FOR UNMANNED AERIAL VEHICLE REAL-TIME TRACKING

Volodymyr Oleksiuk*  , Serhiy Velhosh  

Department of Radiophysics and Computer Technologies,
Ivan Franko National University of Lviv,
107 Gen. Tarnavskoho St., 79017 Lviv, Ukraine

Oleksiuk, V., Velhosh, S. (2026). Entropy-Guided Tracker Switching Method for Unmanned Aerial Vehicle Real-Time Tracking. *Electronics and Information Technologies*, 33, 131–144. <https://doi.org/10.30970/eli.33.10>

ABSTRACT

Background. Auto-guidance for unmanned aerial vehicles (UAVs) requires reliable real-time target tracking on resource-constrained onboard hardware. Modern state-of-the-art CNN-based and Transformer-based deep trackers provide strong accuracy but are often too slow and computationally expensive for continuous deployment on edge devices. In contrast, lightweight correlation-filter trackers run at high frame rates but can easily drift or lose the target because of occlusions or fast maneuvers. This robustness–efficiency trade-off (edge AI paradox) motivates adaptive strategies that balance accuracy, speed, and resource usage while preserving compute headroom for other onboard tasks.

Materials and methods. We propose an entropy-guided tracker switching method that combines a lightweight kernelized correlation filter (KCF) tracker augmented with Kalman motion prediction and a more accurate Siamese deep tracker. A motion-entropy scheduler quantifies the unpredictability of target motion using a normalized Shannon entropy over recent orientation changes. To avoid reacting to transient spikes, the entropy is exponentially smoothed, and threshold rules (with hysteresis) determine when KCF is sufficient and when to activate the deep tracker.

Results and Discussion. Experiments on UAV benchmarks (UAV123, OTB100) show that the hybrid tracker improves success AUC by ~10% over KCF and reaches about 70% of a Transformer tracker's AUC while running 1.5–3× faster than always-on deep tracking. The switcher invokes the deep tracker only during difficult intervals, sustaining real-time operation (~100 FPS) and reducing average computation to ≈0.6 GFLOPs per frame versus ≈1–4 GFLOPs for purely deep tracking.

Conclusion. The proposed motion-entropy scheduler enables an adaptive trade-off between efficiency, speed, and accuracy. It maintains high tracking precision during target maneuvers and occlusions by temporarily switching to a robust tracker yet saves computational load during steady-motion periods. This framework offers a practical solution for high-performance UAV tracking on the edge, while leaving resource headroom to apply other improvement techniques.

Keywords: object tracking; correlation filters; Siamese network; motion entropy; hybrid tracker; edge computing.

INTRODUCTION

Real-time target tracking by unmanned aerial vehicles (UAVs) demands both high accuracy and high efficiency. UAV onboard computers have limited processing power and energy yet must handle fast-moving targets and complex backgrounds. This creates a



© 2026 Volodymyr Oleksiuk & Serhiy Velhosh. Published by the Ivan Franko National University of Lviv on behalf of Електроніка та інформаційні технології / Electronics and Information Technologies. This is an Open Access article distributed under the terms of the [Creative Commons Attribution 4.0 License](https://creativecommons.org/licenses/by/4.0/) which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

trade-off often termed an edge AI paradox, where achieving high accuracy with a deep neural model conflicts with strict resource and real-time constraints [1].

On one hand, modern deep trackers (e.g., using Siamese networks or transformers) can achieve excellent accuracy on benchmarks, but their high computational cost makes real-time deployment challenging without powerful hardware. Recent transformer-based trackers such as TransT [2] provide strong accuracy, but their computational cost can still be challenging on resource-constrained UAV platforms; lightweight UAV-oriented trackers (e.g., MobileTrack [3]) were designed for high-speed deployment, which still use too many resources.

On the other hand, classical lightweight trackers like those based on correlation filters run at hundreds of frames per second on CPU, but tend to lose the target during occlusions, rapid appearance changes, or unexpected maneuvers [4, 5]. Ensuring accurate tracking despite such challenges while staying within the computational constraints is a key problem in UAV tracking.

Several factors make UAV tracking particularly difficult. The UAV's moving camera leads to a wide field of view with many distractors (e.g., ground clutter, other objects) and frequent viewpoint changes. Targets often occupy a small image region and can undergo extreme scale changes or become fully occluded by obstacles. Furthermore, wind or platform vibrations introduce motion blur and sudden target displacements [6]. These conditions can cause a simple tracker to drift off target or fail entirely. For example, a high-speed correlation filter tracker like KCF (Kernelized Correlation Filter) excels in nominal conditions with its Fast Fourier Transform (FFT) based template matching [4], but under occlusion or background confusion, it can quickly degrade [7]. Once the target is lost, a naive tracker cannot recover on its own [8]. Recent UAV systems also combine correlation filters with modern detectors to improve robustness under occlusion and scale changes [9, 10].

To address these issues, we propose a hybrid tracking architecture that leverages the complementary strengths of lightweight correlation-filter trackers and resource-intensive trackers. The core idea is to run a fast tracker most of the time and only switch to a computationally intensive but more robust one when necessary. Our framework uses KCF as a lightweight correlation-filter tracker for its real-time speed, augmented with a Kalman filter to predict target motion and smoothly update the region of interest (ROI). KCF's speed and good performance on stable motion make it a strong baseline [4]. However, when the target's motion becomes erratic, or the scene changes in a way likely to confuse KCF, we temporarily hand off to a resource-intensive tracker based on a Siamese neural network. The resource-intensive tracker (in our case, a variant of a MobileTrack Siamese model) is more computationally intensive but can handle significant appearance change and re-detect the object if it disappears and reappears. By combining these, we aim for "the best of both worlds": the efficiency of KCF during easy intervals and the resilience of a deep tracker during hard intervals.

A novel aspect of our approach is the motion entropy scheduler that decides when to switch between trackers. Instead of relying solely on the short-term tracker's internal confidence (e.g., correlation response strength via average peak-to-correlation energy (APCE) to detect failure [11, 12]) or activating the resource-intensive tracker on a fixed periodic verification trigger (e.g., every t frames) [13], we quantify the unpredictability of the target's recent motion using an entropy measure. Intuitively, when the target moves in a steady, predictable manner, a lightweight correlation-filter tracker should suffice. But when the motion is highly irregular or complex (high entropy), it likely foreshadows situations that could confuse the lightweight correlation-filter tracker (such as abrupt maneuvers or interacting objects). In those moments, engaging the more powerful resource-intensive tracker can prevent tracking loss. Our scheduler computes the Shannon entropy [14] of the target's motion orientation distribution and uses it as a switching signal. We introduce a Kalman-filtered motion model to estimate the target

trajectory and compute entropy in real time. High entropy values trigger a switch to the deep tracker, while low entropy allows a switch back to the fast tracker, with hysteresis to avoid rapid oscillations.

In addition, to explicitly handle UAV camera shake (that can cause constantly high entropy), we introduce a measure of residual flow. Intuitively, we first estimate the “smooth” global motion of the frame that corresponds to the drone’s overall movement (for example, a rotation or steady flight), and then we look at the residual flow – the remaining pixel motion after subtracting this global shift. When the drone flies smoothly, this residual motion is well structured, and its entropy is low. When the platform is shaken by wind or aggressive maneuvers, the residual flow becomes chaotic, and its entropy increases sharply. We use this value as part of our frame reliability criteria: if the residual flow entropy is high, the frame is treated as unreliable, and the tracker skips online those corrupted frames.

MATERIALS AND METHODS

Physical Meaning and Theoretical Foundation of Motion Entropy

From a physical perspective, entropy is traditionally viewed as a measure of disorder or statistical uncertainty in a dynamical system. In the context of visual object tracking, motion entropy serves as a quantitative indicator of the trajectory’s dynamical complexity and the degree of randomness in the target’s displacement vector field. When an object moves uniformly and predictably, its displacement vectors exhibit a high degree of spatial correlation (concentrated in a narrow orientation sector), corresponding to a state of minimum entropy. Conversely, the onset of abrupt maneuvers, occlusions, or background clutter causes the motion vectors to scatter across various directions, sharply increasing informational uncertainty and, consequently, the entropy value—signaling a transition to a “chaotic” state. Utilizing entropy as a descriptor allows for the proactive detection of tracking stability loss, often before the tracking algorithm itself fails due to cumulative errors [15].

Motion Entropy Estimation

We now describe how we compute the motion entropy that drives the scheduler. The goal is to produce a scalar H_t at each frame t that quantifies the unpredictability or disorder in the target’s recent motion. Our approach is to analyze the orientation distribution of the target’s frame-to-frame displacement vectors and compute the Shannon entropy of that distribution [14].

Firstly, we obtain the target’s motion vectors over a short time window. In our implementation, we use the Kalman filter’s estimated velocities as well as optical flow within the target region to characterize motion. Specifically, between frame $t - 1$ and t , we take the Kalman-predicted displacement $(\Delta x, \Delta y)$ (which is also like the lightweight correlation-filter tracker’s observed displacement if tracking was successful). Additionally, we compute a sparse optical flow using Shi-Tomasi corners [16] within the target’s bounding box – this gives us multiple motion vectors v_i for different parts of the target or local background. The collection of motion vectors is then projected to orientations (angles). We accumulate these angles into a histogram with N bins covering $[0, 2\pi)$. Based on prior studies [17], we choose $N = 16$ orientation bins (each 22.5° wide), which provide a good balance between resolution and statistical reliability. Let n_k be the count (or sum of weights) of motion vectors falling into the bin k . We weigh each vector’s contribution by its magnitude, so that larger motions influence the entropy more than tiny jitter (this is like the motion magnitude weighting proposed by Chen et al. [17]). After weighing, we compute the probabilities:

$$p_k = n_k / \sum_{j=1}^N n_j. \quad (1)$$

This yields a normalized distribution over orientations.

Shannon Entropy Calculation

The Shannon entropy [14] of the orientation distribution is then:

$$H = - \sum_{j=1}^N p_k \log p_k. \quad (2)$$

We further normalize this value by the maximum entropy $\log N$ (achieved when the distribution is uniform) to get a normalized entropy:

$$H_{norm} = H / \log N, \quad (3)$$

so that it lies between 0 and 1. In this normalized scale, $H_{norm} \sim 0$ means all motion vectors point to the same direction (perfectly regular motion), while $H_{norm} = 1$ means the motion vectors are uniformly distributed across all directions (completely random or highly complex motion). We found this normalization helpful for defining general thresholds that transfer across scenarios. We maintain a smoothed entropy value and update it recursively:

$$\tilde{H}_t = \beta \cdot \tilde{H}_{t-1} + (1 - \beta) \cdot H_{norm,t}, \quad (4)$$

where $\beta \in [0,1)$ is a smoothing factor (we used $\beta = 0.8$ so that the entropy effectively averages over the last few frames). This smoothing dampens transient spikes in entropy and captures the general trend of motion complexity. A sudden one-frame orientation change won't immediately flip the scheduler, unless it persists or continues chaotic. The choice β can be adjusted; a higher β (closer to 1) value means slower adaptation (more inertia in the decision), and a lower β value responds quicker but potentially more oversensitive.

Magnitude Thresholding and Residual Flow

We ignore very small motions when computing entropy. If the target moved less than a few pixels (below a threshold) in a frame, we do not treat that as a meaningful motion direction. In practice, if the displacement magnitude of a vector is below, say, 1 pixel (accounting for noise or static target), we exclude it from the orientation histogram. If all motion vectors in a frame are below this threshold (target basically stationary), we define entropy as $H_t = 0$ by default – there is no unpredictability in not moving. This prevents random orientation values due to noise from artificially inflating entropy when the target is actually static.

In UAV footage, a large portion of the optical flow is caused by the camera itself (translation, yaw, small vibrations), so entropy computed directly from raw vectors may reflect platform motion rather than target dynamics. To isolate the target, we first estimate a dominant global motion field from background feature tracks outside the target box. We then subtract this global component from each motion vector inside the target region and

build the orientation histogram from the remaining residual vectors. In parallel, we compute an entropy value for the global flow field; subtracting this baseline from the target-region entropy yields a relative residual-flow entropy that stays low during smooth flight but rises when local motion becomes inconsistent (abrupt maneuvers, partial occlusion, distractors). The resulting value is clipped to the normalized range and used as a frame-reliability cue (e.g., to avoid updating the low-cost tracker during camera shake).

After these steps, we obtain the smoothed motion entropy H_t for each frame. This value will be used by the scheduler logic. We determine two key hyperparameters related to H : a high threshold T_{high} for switching to the resource-intensive tracker, and a low threshold T_{low} for switching back to the lightweight correlation-filter tracker. Additionally, we set a cooldown interval to avoid switching too frequently (even with hysteresis, we impose that once a switch occurs, the system waits at least e.g., 5 frames before another switch).

Entropy-Guided Scheduler Logic

With the smoothed entropy signal \tilde{H}_t in hand, the scheduler applies a hysteresis rule to decide the tracker mode. We maintain a binary mode state in [Lightweight tracker, High-accuracy tracker], indicating which tracker is currently in control of the target output. The state transition logic is:

Lightweight correlation-filter tracker mode (normal): Remain in lightweight correlation-filter tracker mode as long as $\tilde{H}_t < T$. If the condition $\tilde{H}_t > T$ is sustained for a certain period (we require it to be above for at least 5 consecutive frames to avoid jitter trigger), then transition to deep neural tracker mode. When transitioning, signal the deep neural tracker to start tracking from the current frame t onward. We also optionally take the current KCF state (position, size) and refine it with the deep neural tracker immediately to avoid any one-frame delay. In practice, we run the deep neural tracker on the same frame t as soon as the decision is made, so we get the benefit of the deep tracker without waiting for the next frame. The mode switch triggers a reset of the cooldown timer.

Resource-intensive tracker mode (robust): Remain in resource-intensive tracker mode if $\tilde{H}_t > T$ the condition is true continuously for a short window (again, we use a 5-frame confirmation), then transition back to lightweight correlation-filter tracker mode. Upon transitioning, we reinitialize or update the lightweight correlation-filter tracker with the resource-intensive tracker's latest bounding box. Specifically, we re-center KCF's filter on the resource-intensive tracker position and if possible, update its appearance model with the current frame's target patch (to catch up any appearance changes that happened during resource-intensive tracker mode). After that, the resource-intensive tracker is put on standby (it may stop running to save computation). The mode switch also resets the cooldown if motion entropy drops below T_{low} .

The hysteresis ensures that once we switch to a deep neural tracker due to high entropy, we don't immediately flip back to a lightweight tracker at the first moment entropy dips, which could be an outlier. Likewise, after switching back to the lightweight tracker, we require a significant entropy rise to switch again. We set $T_{high} > T_{low}$ it so there is a clear gap. In our default, 0.65 vs 0.50 (see **Table 1**), which worked well. The cooldown of 5 frames added a further guard: for example, if entropy oscillates around 0.6, the system might switch to deep neural tracker at 0.65, then entropy dips to 0.6 (still above 0.5, so it stays deep neural tracker), then maybe rises again – the cooldown ensures we don't switch off deep neural tracker too quickly and then on again. Essentially, once a deep neural tracker is engaged, we want to stick with it for a reasonable duration (at least a few frames) to see through the turbulent period.

Table 1. Entropy-related parameters.

Parameter	Symbol	Value (default)	Description
Orientation histogram bins	N	16	Number of bins for motion direction (0–360°).
Smoothing factor	β	0.8	Exponential smoothing weight for H_t .
High entropy threshold	T_{high}	0.65	Threshold to trigger the switch to the high-accuracy tracker (dimensionless, between 0 and 1).
Low entropy threshold	T_{low}	0.50	Threshold to trigger the switch back to the lightweight tracker.
Motion ignore threshold	–	1 px/frame	Motions below this pixel magnitude are treated as zero motion (noise filter).
Switch frames threshold	–	5 frames	Minimum number of frames to wait after a switch before another switch is allowed.

While in lightweight correlation-filter tracker mode, the resource-intensive tracker could either be completely off or running at a low frequency. In our implementation, we opted to run the resource-intensive tracker at a very low frequency (e.g., once every 10 frames) even in lightweight correlation-filter tracker mode, just to maintain an updated idea of the target's appearance in case it has changed significantly. This is not strictly necessary; one could turn off the resource-intensive tracker entirely to save power and only initialize it when needed (incurring a small initiation cost). We found the overhead of running it occasionally was minimal, and it provided a slight safety net (if KCF was close to failure, the resource-intensive tracker might already be able to pick up immediately). However, for simplicity, one can imagine the resource-intensive tracker is effectively off during lightweight correlation-filter tracker mode.

Conversely, while in resource-intensive tracker mode, we often still run the lightweight correlation-filter tracker in the background. The lightweight correlation-filter tracker might fail during this period (since presumably entropy was high for a reason), but we let it continue updating with the Kalman predictions and occasional corrections from resource-intensive tracker output. The benefit is that when we switch back to the lightweight correlation-filter tracker, it can resume without a full re-initialization from scratch. In some cases, if the lightweight correlation-filter tracker completely lost the target during resource-intensive tracker mode (e.g., KCF drifted to background because we stopped updating it), we simply reinitialize KCF at switch-back time, which is essentially instantaneous.

Implementation and Platform

We implemented the proposed framework using Python 3.10, and the Siamese resource-intensive tracker was executed with PyTorch 2.1.0. The experiments were conducted in the Google Colab virtual environment. Therefore, the reported results should be interpreted within the context of a simulation-based setup for resource-constrained UAVs, rather than as an execution on actual onboard drone hardware. The setup reflected

a lightweight CPU stage (KCF, Kalman prediction, optical flow) and a GPU-accelerated deep-learning stage, but it did not represent measurements on a specific UAV hardware platform. The allocated session provided hardware resources roughly equivalent to an Intel Xeon-class 2 vCPU CPU, an NVIDIA T4 GPU with 16 GB VRAM, and approximately 12.7 GB of RAM. No explicit instruction-set restrictions were applied. The baseline FPS for each tracker was measured by running it on the entire dataset and averaging the rate. Note that for the hybrid algorithm, the FPS can vary sequence-to-sequence depending on how often the resource-intensive tracker is invoked; we report the overall average and discuss the range. Actual onboard deployment would require additional profiling on the target UAV platform under its memory and power constraints.

We tuned the hyperparameters (entropy thresholds, etc.) empirically and fixed them for all results presented. For APCE-Hybrid, we set the APCE threshold analogously to the entropy-based threshold so that it triggers the resource-intensive tracker under comparable conditions of tracking difficulty. We calibrated this threshold by identifying APCE values that consistently decreased shortly before KCF lost the target.

RESULTS AND DISCUSSION

Datasets and Baselines

We evaluate our approach on two standard single-object tracking datasets that include UAV scenarios and varied difficulties: UAV123 [18] is a UAV-specific benchmark with 123 aerial video sequences (over 110K frames) capturing objects such as cars, boats, and people from a drone perspective. It includes challenging attributes like fast motion, camera motion, small objects, and occlusion – reflecting real UAV tracking conditions; OTB100 [19] is the Object Tracking Benchmark (100 sequences) widely used in tracking literature; it contains various scenes (not UAV-specific) and provides established evaluation protocols (success and precision). These two datasets cover a broad spectrum from short-term to long-term tracking scenarios.

We compare the proposed entropy-hybrid tracker with the following baselines (which are modern tiny state-of-the-art models to run in real-time on resource-constrained devices):

1. KCF – the baseline correlation filter tracker without any switching (essentially our lightweight correlation-filter tracker running alone).
2. MobileTrack – a recent efficient Siamese tracker optimized for UAVs, representing the state-of-the-art in high-speed tracking (we use an implementation of MobileTrack for comparisons).
3. TransT – a modern transformer-based tracker [2], representing top-tier accuracy (albeit at a higher computational cost). TransT does not impose real-time constraints on our platform but gives an upper bound on accuracy for reference.
4. Fixed periodic trigger – a simple hybrid baseline that activates the deep neural tracker at a fixed interval (every t frames) regardless of confidence or entropy, similar to periodic verification strategies used in long-term tracking [13].
5. APCE-Hybrid as an ablative baseline like Cao et al. [9]: this is our implementation of a tracker that switches between KCF and the Siamese deep neural tracker based on the APCE threshold (rather than entropy). APCE-Hybrid uses the same two trackers as our method but triggers the deep neural tracker when KCF's normalized response peak drops below a set threshold (indicating low confidence). This allows us to contrast entropy vs. direct confidence-based switching.

Metrics

We use standard tracking metrics: Success (AUC) and Precision. Success is measured by the Intersection-over-Union (IoU) overlap between the predicted bounding

box and ground truth, and the success rate is the fraction of frames where IoU exceeds a threshold. By varying the threshold from 0 to 1, an area under the curve (AUC) score is computed. We report the AUC as a summary of overall accuracy. Precision is measured by Center Location Error (CLE): the percentage of frames where the predicted center is within a certain distance (typically 20 pixels) of the ground truth center [19]. We report the precision at a 20px threshold, following the OTB100 benchmarking methodology, as well as the mean CLE. Additionally, since our focus is on real-time performance, we report the average Frames Per Second (FPS) each tracker runs at (on a given hardware configuration), and an estimate of computational complexity in FLOPs (floating point operations per frame). The FLOPs give a hardware-independent measure of efficiency: we calculate it by summing the major operations of each tracker per frame (for deep trackers, this includes convolution operations; for KCF, it's negligible). For our hybrid, the effective FLOPs per frame are lower than those of the deep tracker since the deep component runs intermittently.

Accuracy and Speed Comparison

Table 2 reports the overall performance of our Entropy-Guided Hybrid tracker against the baselines on the UAV123 and OTB100 datasets. We list the Success AUC, Precision (20px), FPS, and average GFLOPs per frame for each method.

Table 2. Performance comparison of the proposed entropy-guided switching method versus baseline trackers on UAV123 and OTB100.

Tracker	UAV123		OTB100		FPS	FLOPs (G)
	Success AUC	Precision	Success AUC	Precision		
KCF	0.432	51.7%	0.521	63.4%	160	~0.02
MobileTrack (2022)	0.690	77.3%	0.617	81.1%	80	~1.2
TransT (2021)	0.717	81.4%	0.754	85.0%	35	~8.0
Fixed periodic trigger (2019)	0.555	65.0%	0.582	67.5%	95	~0.8
APCE Hybrid (2025)	0.573	66.5%	0.590	69.5%	110	~0.9
Entropy-Hybrid (Ours)	0.594	69.5%	0.601	72.0%	100	~0.6

Our entropy-hybrid tracker improves substantially over KCF on UAV123, increasing Success AUC from 0.432 to 0.594 and Precision from 51.7% to 69.5%. Although always-on deep trackers still lead in pure accuracy (MobileTrack: 0.690 AUC, 77.3% precision; TransT: 0.717 AUC, 81.4% precision), our hybrid achieves a better efficiency–accuracy balance, running at 100 FPS with ~0.6 GFLOPs per frame versus 80 FPS / ~1.2 GFLOPs for MobileTrack and 35 FPS / ~8.0 GFLOPs for TransT.

Compared to the fixed periodic trigger baseline [13], entropy-based switching yields higher accuracy (0.594 vs 0.555 AUC on UAV123) while also lowering average compute (~0.6 vs ~0.8 GFLOPs) by avoiding unnecessary deep activations during stable motion.

Against APCE-Hybrid, our method gains 0.021 AUC and 3.0 percentage points precision on UAV123 (0.573→0.594, 66.5%→69.5%) with a notable reduction in average computation (~0.9→~0.6 GFLOPs). This supports the idea that motion entropy can trigger the resource-intensive tracker more proactively than response-peak degradation, improving continuity before the low-cost tracker drifts.

On OTB100, the entropy hybrid is closed in accuracy to MobileTrack in AUC (0.617 vs 0.601) while keeping a higher frame rate (100 vs 80 FPS). The relative gain over KCF remains clear (0.521→0.601 AUC), indicating that switching is beneficial beyond UAV-specific footage, even though OTB sequences often contain longer low-entropy intervals where the system stays in lightweight mode.

Ablation Analysis

We perform an ablation study to isolate the impact of our entropy-guided switching. We compare four configurations: (A) KCF-only, (B) deep neural tracker-only (MobileTrack every frame), (C) APCE-triggered hybrid, and (D) Entropy-triggered hybrid (ours). Configuration A and B represent the extremes of never switching (always using one tracker). C and D use the same components but different switching signals.

Table 3 confirms that the switching mechanism is crucial. Either hybrid vastly outperforms the single trackers, proving the effectiveness of combining trackers. Between triggers, the entropy-based scheduler provides a better balance, improving accuracy the most with only minimal overhead. It validates our choice of using motion entropy as a superior switching signal compared to solely relying on the tracker’s internal confidence.

Table 3. Ablation of tracker switching strategies on UAV123.

Strategy	AUC	Precision	Avg FPS	Notes
A. KCF-only (no switch)	0.432	51.7%	160	Fast but loses target often.
B. Deep neural tracker-only (Siamese every frame)	0.690	77.3%	80	Accurate, but always computationally intensive.
C. Hybrid (APCE switch)	0.573	66.5%	110	Switches on KCF failure (reactive).
D. Hybrid (Entropy switch)	0.594	69.5%	100	Switches on motion entropy (proactive).

Discussion

The proposed approach demonstrates that an information-theoretic measure like entropy can be highly useful in a control loop for visual tracking. By focusing on motion characteristics (which are agnostic to appearance), our scheduler gains a kind of foresight into tracking difficulty. This is particularly beneficial in UAV scenarios where camera motion and target motion interplay. For instance, if a drone suddenly turns, causing the background to shift in a new direction relative to the target, the entropy of the optical flow in the scene spikes – our system notices this and preemptively bolsters the tracking with the Siamese network. This behavior is a distinct advantage over reactive schemes that wait for the primary tracker to fail. In effect, motion entropy serves as a proxy for “how hard is the target to follow right now,” encapsulating factors like erratic

target maneuvering or complex background motions. This work suggests that integrating such high-level cues can make traditional trackers smarter without explicitly training a switching classifier or network.

Our results demonstrate the trade-off between efficiency and accuracy in edge deployment. Running a top-tier tracker (TransT) would give the best accuracy, but at an untenable speed/power cost on UAV hardware. Running a lightweight correlation-filter tracker (KCF) yields great speed but will miss the target often. By spending computation only on the hard frames, we allocate resources dynamically where they have the most impact, achieving an overall efficiency that would be impossible if the resource-intensive tracker ran uniformly. This is particularly important for UAVs running on battery – computational load often directly translates to battery drain and flight time limitations. Our method allows the system to cruise on low power most of the time and only spike usage occasionally.

While our Kalman filter integration improves short gaps handling, it does introduce a risk: if the target is fully occluded and the Kalman filter propagates the state without measurement updates for too long, it can drift off course. In our tests, if the occlusion lasted beyond about 1 second, the Kalman prediction became unreliable. If the target reappears far from the predicted spot, KCF might not search far enough to find it. Our current system mitigates this by relying on the high-accuracy tracker in such cases – the entropy typically would be high during the occlusion onset (e.g., when the target disappears behind an object, there is often a flurry of motion edges), triggering the high-accuracy tracker. The high-accuracy tracker's larger search window and template matching can reacquire the target even if it breaks linear motion assumptions. However, in extreme cases where the target stops moving entirely behind an occluder (zero motion, so low entropy, and Kalman just coasts), neither tracker may see it until it reappears. If the reappearance is significantly off the predicted path, our current strategy might not catch it immediately. This is where integrating a dedicated re-detector (like a YOLO) could help. In future extensions, we could incorporate a third component: if the target is not found by either tracker for a certain period, trigger a full-frame detection. This would handle total failures. The trade-off is the complexity and potential false positives from detection, so we avoided it in this work.

Introducing entropy shows it works well in our UAV context because it naturally focuses on residual motion rather than pure camera-induced motion. When the entire scene moves uniformly (due to drone translation), the optical flow vectors are largely aligned and have low entropy. In our implementation, we first approximate a dominant background flow vector in the region and subtract it, computing entropy on the residual flow field. This step suppresses global jitter and small handshake effects from the UAV device. But if the motion is globally consistent, KCF (augmented with a Kalman filter applied for motion prediction) can handle it well. It is the irregular residual motion that causes trouble (e.g., the target abruptly changing direction or multiple objects crossing paths and generating flows in different directions), and the entropy metric cleanly reflects that irregularity.

CONCLUSION

We presented an entropy-guided hybrid tracker for real-time UAV target tracking that dynamically switches between a fast correlation filter tracker and a deep Siamese tracker. The central contribution is the introduction of motion entropy as a decision signal to anticipate tracking difficulty. By measuring the randomness of recent target motion, our scheduler intelligently allocates computational effort: it keeps the lightweight correlation-filter tracker in charge during easy, predictable segments. It activates the resource-intensive tracker during challenging maneuvers or occlusions. This approach

tackles the edge computing paradox by achieving high accuracy when needed without continuously running expensive algorithms.

Through experiments on several benchmarks, we demonstrated that our method delivers a balance of accuracy and efficiency. On UAV123, it improved tracking success by 0.162 AUC compared to a baseline KCF (0.432→0.594) while maintaining ~100 FPS. The hybrid also outperformed both a fixed periodic trigger baseline and a purely confidence-based switching baseline, highlighting the efficacy of the entropy cue. These results indicate that motion entropy is a practical heuristic for tracker management that adapts its strategy based on observed dynamics.

In practical terms, our system can extend the deployment of advanced tracking on resource-limited UAV platforms. A drone using this tracker can handle fast-moving or erratic targets with accuracy nearly on par with modern deep trackers, yet for much of the flight, it only expends the computation of a tiny fraction of a deep network. This has direct implications for onboard energy consumption and responsiveness.

There are several future directions for this work. One is to integrate an object detector as a fail-safe for prolonged target loss, combining detection with our tracking switcher to handle complete occlusions or exits from the field of view. Another direction is to make the entropy scheduler a neural network which can be trained to better determine complex scenes and further apply self-learning algorithms to improve the model on fly.

In conclusion, entropy-guided tracker switching offers a promising way to satisfy the competing demands of accuracy and efficiency in UAV target tracking. It uses a broader principle of adaptive edge AI: using inexpensive computations to decide when to deploy expensive ones. We believe this paradigm can be applied to many complex problems on resource-constrained platforms in real-time, enabling smarter and more autonomous systems without hardware upgrades.

ACKNOWLEDGMENTS AND FUNDING SOURCES

The authors received no financial support for the research, writing, and/or publication of this article.

COMPLIANCE WITH ETHICAL STANDARDS

The authors declare that the research was conducted in the absence of any conflict of interest.

AUTHOR CONTRIBUTIONS

Conceptualization, [V.O.]; methodology, [V.O.]; validation, [V.O.]; writing – original draft preparation, [V.O.]; writing – review and editing, [S.V., V.O.]; supervision, [S.V.].

All authors have read and agreed to the published version of the manuscript.

REFERENCES

- [1] Singh, A., Saini, K., Nagar, V., Aseri, V., Sankhla, M. S., Pandit, P. P., & Chopade, R. L. (2022). Artificial intelligence in edge devices. *Advances in Computers*, 127, 437–484. <https://doi.org/10.1016/bs.adcom.2022.02.013>
- [2] Chen, X., Yan, B., Zhu, J., Wang, D., Wang, X., & Lu, H. (2021). Transformer tracking. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 8126–8135. <https://doi.org/10.1109/CVPR46437.2021.00803>
- [3] Xue, Y., Jin, G., Shen, T., Tan, L., Yang, J., & Hou, X. (2022). MobileTrack: Siamese efficient single object tracker for high-speed UAV tracking. *IET Image Processing*, 16(12), 3300–3313. <https://doi.org/10.1049/ipr2.12565>

- [4] Henriques, J. F., Caseiro, R., Martins, P., & Batista, J. (2015). High-speed tracking with kernelized correlation filters. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(3), 583–596. <https://doi.org/10.1109/TPAMI.2014.2345390>
- [5] Bolme, D. S., Beveridge, J. R., Draper, B. A., & Lui, Y. M. (2010). Visual object tracking using adaptive correlation filters. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2544–2550. <https://doi.org/10.1109/CVPR.2010.5539960>
- [6] Wu, P., Li, Y., & Xue, D. (2025). UAV target tracking: A survey. *Artificial Intelligence Review*, 58(11), 1–41. <https://doi.org/10.1007/s10462-025-11348-x>
- [7] Zhang, Y., Yang, Y., Zhou, W., Shi, L., & Li, D. (2018). Motion-Aware Correlation Filters for Online Visual Tracking. *Sensors*, 18(11), 3937. <https://doi.org/10.3390/s18113937>
- [8] Kalal, Z., Mikolajczyk, K., & Matas, J. (2012). Tracking-learning-detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(7), 1409–1422. <https://doi.org/10.1109/TPAMI.2011.239>
- [9] Cao, S., et al. (2025). UAV real-time target detection and tracking algorithm based on improved KCF and YOLOv5s_MSES. *Machines*, 13(5), 364. <https://doi.org/10.3390/machines13050364>
- [10] Ai, Y., et al. (2025). Real-time occluded target cooperative tracking method for UAVs. *Electronics*, 14(20), 4034. <https://doi.org/10.3390/electronics14204034>
- [11] Wang, J., Liu, Y., Ai, Y., & Xue, W. (2021). Long-term target tracking combined with re-detection. *EURASIP Journal on Advances in Signal Processing*, 2021, 79. <https://doi.org/10.1186/s13634-020-00713-3>
- [12] Ma, H., Acton, S. T., & Lin, Z. (2020). SITUP: Scale invariant tracking using average peak-to-correlation energy. *IEEE Transactions on Image Processing*, 29, 3546–3557. <https://doi.org/10.1109/TIP.2019.2962694>
- [13] Liu, F., Mao, K., Qi, H., & Liu, S. (2019). Real-time long-term correlation tracking by single-shot multibox detection. *Optical Engineering*, 58(1), 013105. <https://doi.org/10.1117/1.OE.58.1.013105>
- [14] Shannon, C. E. (1948). A mathematical theory of communication. *Bell System Technical Journal*, 27(3), 379–423. <https://doi.org/10.1002/j.1538-7305.1948.tb01338.x>
- [15] Jaynes, E. T. (1957). Information theory and statistical mechanics. *Physical Review*, 106(4), 620–630. <https://doi.org/10.1103/PhysRev.106.620>
- [16] Shi, J., & Tomasi, C. (1994). Good features to track. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, 593–600. <https://doi.org/10.1109/CVPR.1994.323794>
- [17] Chen, C.-Y., et al. (2008). Motion entropy feature and its applications to event-based segmentation of sports video. *EURASIP Journal on Image and Video Processing*, 2008, 460913. <https://doi.org/10.1155/2008/460913>
- [18] Mueller, M., Smith, N., & Ghanem, B. (2016). A benchmark and simulator for UAV tracking. *European Conference on Computer Vision (ECCV)*, 445–461. https://doi.org/10.1007/978-3-319-46448-0_27
- [19] Wu, Y., Lim, J., & Yang, M.-H. (2015). Object tracking benchmark. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(9), 1834–1848. <https://doi.org/10.1109/TPAMI.2014.2388226>
- [20] Bertinetto, L., Valmadre, J., Henriques, J. F., Vedaldi, A., & Torr, P. H. S. (2016). Fully-convolutional siamese networks for object tracking. *European Conference on*

- Computer Vision Workshops (ECCVW), 850–865.
https://doi.org/10.1007/978-3-319-48881-3_56
- [21] Ye, B., Chang, H., Ma, B., Shan, S., & Chen, X. (2022). Joint feature learning and relation modeling for tracking: A one-stream framework. *Computer Vision – ECCV 2022*, 341–357. https://doi.org/10.1007/978-3-031-20047-2_20
- [22] Cui, Y., Jiang, C., Wu, G., & Wang, L. (2024). MixFormer: End-to-end tracking with iterative mixed attention. *IEEE Transactions on Pattern Analysis and Machine Intelligence*. <https://doi.org/10.1109/TPAMI.2024.3349519>
- [23] Javed, S., Danelljan, M., Khan, F. S., Khan, M. H., Felsberg, M., & Matas, J. (2023). Visual Object Tracking With Discriminative Filters and Siamese Networks: A Survey and Outlook. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(5), 6552–6574. <https://doi.org/10.1109/TPAMI.2022.3212594>
-

КЕРУВАННЯ ПЕРЕМІКАННЯМ ЗАСОБІВ СУПРОВОДУ НА ОСНОВІ ЕНТРОПІЙНОГО АНАЛІЗУ ДЛЯ ВІДСТЕЖЕННЯ ЦІЛЕЙ БЕЗПЛОТНИХ ЛІТАЛЬНИХ АПАРАТІВ У РЕАЛЬНОМУ ЧАСІ

Володимир Олексюк* , Сергій Вельгош 

*Кафедра радіофізики та комп'ютерних технологій,
Львівський національний університет імені Івана Франка,
вул. Генерала Тарнавського, 107, 79017 Львів, Україна*

АНОТАЦІЯ

Вступ. Автоматичне наведення безпілотних літальних апаратів потребує надійного супроводу цілі в реальному часі на бортовому обладнанні з обмеженими обчислювальними ресурсами. Сучасні засоби супроводу на основі згорткових нейронних мереж і трансформерів забезпечують високу точність, однак для безперервної роботи на периферійних пристроях часто є надто повільними та обчислювально затратними. Натомість методи на основі кореляційних фільтрів працюють із високою частотою кадрів, але можуть накопичувати похибку або втрачати ціль за умов перекриття чи різких маневрів. Цей компроміс між точністю та ресурсоефективністю (парадокс застосування штучного інтелекту на периферійних пристроях) зумовлює потребу в адаптивних підходах, що узгоджують швидкодію та витрати ресурсів і зберігають обчислювальний резерв для інших бортових задач.

Матеріали та методи. Запропоновано метод перемикання засобів супроводу на основі аналізу ентропії, що поєднує швидкий алгоритм на основі кореляційних фільтрів (KCF), доповнений прогнозуванням руху за фільтром Калмана, та точніший глибинний сіамський засіб супроводу. Планувальник ентропії руху оцінює непередбачуваність переміщення цілі за нормованою ентропією Шеннона змін орієнтації руху на останніх кадрах. Для зменшення впливу короточасних сплесків застосовано експоненційне згладжування, а порогові правила визначають моменти, коли достатньо KCF і коли слід активувати глибинний модуль.

Результати. Під час випробувань на еталонних наборах даних UAV123 і OTB100 гібридний підхід підвищив success AUC приблизно на 10% порівняно з KCF і приблизно 70% від показника трансформерної моделі, при цьому працюючи у 1,5–3 рази швидше, ніж режим із постійно ввімкненим глибинним модулем модуль вмикається лише у складні проміжки, забезпечуючи близько 100 кадрів/с та середні витрати $\approx 0,6$ GFLOPs на кадр проти $\approx 1-4$ GFLOPs для сучасних глибинних підходів.

Висновки. Розроблений метод відстеження цілі на основі аналізу ентропії руху забезпечує адаптивний компроміс між ефективністю, швидкістю та точністю, підтримуючи високу точність під час маневрів і перекриттів цілі та зменшуючи навантаження за стабільного руху. Запропонований підхід є практичним рішенням для вискоєфективного супроводу цілей БПЛА на бортових пристроях.

Ключові слова: відстеження цілі; кореляційні фільтри; сіамська нейронна мережа; ентропія руху; гібридний алгоритм супроводу; периферійні обчислення.

UDC: 004.89:005.8

PREDICTIVE THERMAL MANAGEMENT IN EMBEDDED ELECTRONICS USING DEEP REINFORCEMENT LEARNING

Oleh Yatskiv^{*}, Bohdan Koman

Department of System Design
Ivan Franko National University of Lviv,
50 Dragomanova St., UA–79005 Lviv, Ukraine

Yatskiv, O.Y., Koman, B. P. (2026). Predictive Thermal Management in Embedded Electronics Using Deep Reinforcement Learning *Electronics and Information Technologies*, 33, 145–164. <https://doi.org/10.30970/eli.33.11>

ABSTRACT

Background. This paper presents a deep reinforcement learning approach for intelligent thermal management in embedded electronics, targeting energy-efficient and safe operation under dynamic workloads. A custom hardware switching circuit based on an NPN transistor was designed to enable GPIO-driven fan actuation on a resource-constrained platform.

Materials and Methods. A real-time dataset was collected from a Raspberry Pi Zero W, capturing CPU temperature, usage metrics, and fan states over a 12-hour controlled experiment. The thermal regulation task was modeled as a Markov Decision Process, and a Deep Q-Network (DQN) was trained to learn optimal fan activation policies. The trained model was deployed directly on-device, interfaced with a custom GPIO-controlled fan circuit. Inference was performed in less than one millisecond per decision step using a lightweight PyTorch runtime.

Results and Discussion. Evaluation results show that the DQN policy reduced total fan activation time by 23.2% compared to the rule-based hysteresis baseline, while maintaining CPU temperature below 60°C for over 99% of the test duration. The trained agent activated the fan only 23.7% of the time, demonstrating a conservative and energy-aware cooling strategy. Confusion matrix analysis yielded a precision of 1.000, a recall of 1.000, and an F1-score of 1.000 across 3442 model-controlled evaluation steps. The model correctly identified all 22 fan activation events without any false positives or false negatives. Comparative analysis against nine recent AI-driven approaches showed that the proposed method achieved an 11.2°C temperature reduction and 36.5% energy savings, while operating entirely on-device without cloud dependence.

Conclusion. The model exhibited stable reward convergence, accurate action prediction, and anticipatory control that minimized overheating events. Thermal traces confirmed smooth transitions and low variance, demonstrating the feasibility of deploying learning-based thermal policies in real-time edge environments. This work contributes a practical framework for energy-aware cooling and provides a pathway for adaptive thermal intelligence in low-resource embedded systems.

Keywords: thermal management, deep reinforcement learning, embedded systems, Deep Q-Network, energy-efficient cooling, real-time inference.

INTRODUCTION

Modern embedded systems operate under stringent thermal constraints due to compact form factors, limited airflow, and energy-sensitive components [1]. Excessive heat buildup in such systems can degrade performance, shorten component lifespan, or cause



© 2026 Oleh Yatskiv & Bohdan Koman Published by the Ivan Franko National University of Lviv on behalf of Електроніка та інформаційні технології / Electronics and Information Technologies. This is an Open Access article distributed under the terms of the [Creative Commons Attribution 4.0 License](https://creativecommons.org/licenses/by/4.0/) which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

abrupt shutdowns [2]. Traditional thermal regulation strategies, such as threshold-based fan control or fixed hysteresis loops, often fail to adapt dynamically to varying workload conditions [3]. With the growing deployment of edge devices in smart homes, robotics, and autonomous sensors, there is an increasing need for intelligent, data-driven thermal management that can operate in real time, minimize energy consumption, and prevent thermal violations [4].

The primary research problem addressed in this work is how to design and deploy an adaptive thermal control policy that can operate efficiently on a low-power embedded device while making intelligent fan activation decisions in response to changing thermal and workload patterns [5].

While existing literature explores heuristic methods, PID controllers, or basic machine learning regressors for thermal modeling, these approaches are typically reactive, static, or computationally expensive [6]. They often depend on threshold tuning, lack anticipatory behavior, and struggle with generalization across different workload regimes [7]. Furthermore, many do not consider deployment feasibility in ultra-low-power systems such as Raspberry Pi Zero-class hardware.

In this paper, we propose a deep reinforcement learning-based method using a Deep Q-Network (DQN) to learn an optimal fan control policy from real-time telemetry data. Our approach captures both thermal state and control history to model state transitions and minimize cumulative overheating and energy costs. Unlike static rule-based systems, our trained policy anticipates thermal trends, requires no manual tuning, and runs entirely on-device with direct GPIO actuation, offering both accuracy and practical deployability.

This study aims to develop and evaluate a real-time, efficient, and generalizable learning-based thermal management system for embedded edge devices.

Research Objectives

- To collect and preprocess thermal and performance telemetry from a constrained embedded platform operating under diverse load conditions
- To train a Deep Q-Network for thermal control, optimizing for energy efficiency and temperature safety
- To deploy and evaluate the trained policy in a live inference environment with GPIO-based fan control.

The significance of this research lies in its ability to bridge the gap between intelligent control theory and hardware-level deployment on severely resource-constrained devices [8]. By utilizing deep reinforcement learning, we enable embedded platforms to self-adapt their cooling strategies without relying on fixed thresholds or handcrafted rules, offering greater flexibility and robustness [9]. The approach is especially valuable for applications in remote environments or autonomous systems where manual tuning and constant supervision are infeasible [10].

Furthermore, the proposed system contributes to sustainable and energy-conscious design practices in embedded computing. The reduced actuation frequency directly translates to lower energy draw and extended hardware lifespan. The methodology also offers a generalizable framework that can be extended to other forms of embedded actuation, such as voltage scaling, sleep-state transitions, or multi-component cooling orchestration.

The remainder of this paper is organized as follows. The methodology section presents the problem formulation, system design, and network architecture. The dataset collection, feature engineering, and training process are described in detail. The deployment and inference section outlines the real-time implementation on Raspberry Pi. Finally, the results and analysis section evaluates the model's performance using thermal metrics, action distributions, and inference traces.

Literature review

Iranfar et al. [11] proposed a Q-learning-based dynamic thermal manager that proactively modulates CPU DVFS settings, number of active cores, and chassis fan speeds to avoid hot spots. Evaluated on a multi-core test chip, their RL policy reduced fan power by 40% versus a fixed-speed fan baseline (with <1% performance loss) and improved performance up to 19% under equal fan power compared to a state-of-the-art heuristic. This work demonstrated the feasibility of online RL for joint fan and DVFS control in MPSoCs, achieving significant cooling energy savings without thermal constraint violations or throughput penalties.

To manage server CPU temperatures and energy, Lin et al. [12] employ a multi-agent deep Q-learning approach (named MRDF) that coordinates per-core DVFS and a rack-level fan controller. Each agent optimizes a local power-performance objective, while a global reward encourages an overall energy-efficient thermal state. Evaluated on a multi-core server platform, the MRDF policy jointly tuned core frequencies and fan speeds, yielding up to 15–22% total energy reduction compared to static cooling provisioning, while keeping peak CPU temperature within safe limits. Notably, MRDF's distributed RL scheme achieved better trade-offs than conventional PID fan controllers or reactive DVFS governors, illustrating the benefit of AI in holistic server thermal management on resource-constrained edge servers.

Maity et al. [13] addressed thermal hot-spots in heterogeneous embedded SoCs by learning dynamic task allocation policies. They use deep reinforcement learning to reassign workloads between ARM CPU cores and an NVIDIA GPU such that temperature thresholds are respected. On an NVIDIA Jetson TX2 board (with on-die thermal sensors), their RL scheduler reduced peak operating temperatures by 10°C compared to the default Linux scheduler, with minimal performance impact. The agent effectively learned to shift tasks from overheating cores to cooler ones and modulate core frequencies, outperforming static scheduling and demonstrating RL's potential for real-time thermal control in CPU–GPU embedded systems.

Kim et al. [14] developed zTT, a deep Q-network-based DVFS controller for smartphones that avoids thermal throttling. The zTT agent jointly adjusts big.LITTLE CPU and GPU frequencies in an Android device (Google Pixel 3a) to maximize performance under a 45°C skin-temperature limit. Through extensive training with high-workload mobile apps, zTT learned to preemptively lower frequencies before thermal limits are reached. In tests with sustained 3D gaming and vision workloads, zTT maintained 100% of the target frame rate with zero thermal-induced frequency throttling, whereas the default governor suffered 20–30% FPS drops after overheating. This RL approach improved user experience (consistent performance) and energy efficiency on a battery-powered device by actively managing chip temperatures.

Instead of explicit control of cooling hardware, Zhou et al. [15] explored an AI-driven workload adaptation strategy to manage device temperature. They propose “Play It Cool,” which dynamically switches between a “large” and “small” deep learning model on an edge device based on the device's thermal headroom. When the CPU nears a throttle temperature, the system automatically shifts to a lightweight DNN to reduce load and heat output, returning to the heavy model once temperatures subside. Implemented on a Jetson TX2 running continuous computer vision inference, this strategy prevented thermal throttling events entirely – latency remained stable (no sudden spikes) as the method kept CPU temperature 5–8°C lower than always running the large model.

Tan and Cao [16] integrated AI into real-time scheduling for mobile SoCs that include dedicated ML accelerators (NPUs). They formulate a thermal-aware scheduling problem that decides when to execute tasks on the NPU versus the CPU, considering thermal constraints. Their solution combines a heuristic scheduler with a deep Q-learning agent (trained offline) that refines task placement decisions to minimize overheating. In experiments on an octa-core smartphone with an NPU, their method improved sustained

inference throughput by 10.4% over a temperature-agnostic scheduler. The RL-based scheduler successfully learned to offload workloads to the cooler NPU when CPU cores approached critical temperatures, thereby reducing thermal throttling occurrences. This work highlights the benefit of AI co-design in edge devices – coordinating CPU–NPU usage to manage thermals while meeting real-time deadlines.

Recent systematic literature reviews have highlighted the growing adoption of AI and ML techniques across various thermal management applications. Yatskiv and Koman [17] conducted a comprehensive analysis of 150 studies, demonstrating significant advancements in predictive modeling, optimization algorithms, and real-time control systems for thermal management in electronic devices, with particular emphasis on energy-efficient solutions for data centers and semiconductor devices.

Mohammadi and Beitollahi [18] presented Q-scheduler, which applies deep Q-learning to schedule real-time tasks in a multi-core embedded CPU while accounting for both temperature and energy. Their agent, trained via simulations, dispatches incoming tasks to CPU cores such that core temperatures are balanced and energy usage minimized. In evaluation on an ARM big.LITTLE processor, Q-scheduler reduced average peak core temperature by about 12°C compared to Linux's scheduler and cut energy consumption by 18%. Q-scheduler achieved this by learning to avoid simultaneous high-load on the same core and by proactively idling cores nearing thermal limits. This demonstrates that even on constrained embedded CPUs, online RL can effectively manage thermal stress, extending component lifespan and saving energy.

Focusing on long-term reliability, Yeganeh-Khaksar et al. [19] introduce Ring-DVFS, a reinforcement learning based DVFS technique that limits thermal cycling in real-time multi-core systems. Thermal cycling (repetitive heating/cooling) accelerates chip aging. Ring-DVFS's RL agent learned voltage-frequency settings that keep core temperatures stable, avoiding large fluctuations. Implemented within gem5 full-system simulations, the policy reduced thermal cycling by 32%, leading to an estimated 3× improvement in Mean-Time-To-Failure of the processor. Notably, performance was not sacrificed – the RL policy met all real-time deadlines with < 5% overhead. This study showcases AI managing not just immediate thermal levels but also optimizing for long-term hardware reliability.

Akhsham et al. [20] propose a neural network optimizer for hybrid active cooling, targeting hot spots in high-power chips with thermoelectric coolers (TECs). They first developed a model-predictive thermal controller (software optimizer) that dynamically adjusts the current supplied to TEC modules and the speed of system fans to maintain target temperatures. Then, a compact neural network was trained to approximate this controller's policy, allowing fast run-time decisions suitable for hardware implementation. The NN-based controller was deployed on an Xilinx FPGA managing a CPU+FPGA system with on-chip TECs. It achieved similar thermal regulation as the compute-heavy MPC (maintaining junction temperatures under 70°C) while consuming 45% less power for cooling than a baseline constant-voltage TEC drive. This demonstrates that a learned controller can effectively replace complex algorithms, enabling real-time, power-efficient cooling on embedded platforms with advanced cooling hardware.

Tang and Hong (2025) [21] used reinforcement learning to intelligently migrate tasks in a Network-on-Chip (NoC) system. Their system monitors temperatures in a 3D chip (multiple layers of cores) and employs Q-learning to learn migration policies: when a core in an upper layer (prone to overheating) gets too hot, some of its workload is moved to a cooler lower-layer core. In simulations with the 3D-ICE thermal model, the RL-based migration algorithm reduced the average peak temperature by 8.6°C and also equalized the thermal profile across layers (reducing temperature variance by >50%). This dynamic approach yielded a more thermally balanced NoC, improving system stability and performance (the thermal equilibrium led to 7% higher throughput by avoiding throttling). It highlights how AI can manage complex thermals in 3D integrated architectures by learning optimal task-to-core mappings beyond what static heuristics can do.

Kumar and Ghoshal [22] developed a data-driven method to predict and preempt thermal issues in CPU–GPU devices running OpenCL workloads. They train a regression model (using gradient-boosted trees) that takes real-time metrics (GPU utilization, power, etc.) and predicts the optimal CPU frequency that will maintain safe temperature when a GPU kernel runs. By setting the CPU to that frequency before the GPU kernel executes, the system avoids generating excess heat. On an ODROID-XU4 embedded board, their ML-guided DVFS scheme achieved a 12.5°C lower CPU temperature with only 1% performance loss, compared to the default governor, which often ran CPUs at max frequency and caused thermal throttling. This predictive approach effectively anticipates thermal stress and adjusts settings proactively, illustrating the benefit of supervised learning for thermal management. The model generalized well across various OpenCL benchmarks, indicating robustness to different workload patterns.

Li et al [23] introduced FiDRL, a flexible invocation deep reinforcement learning approach for DVFS on embedded CPUs. Unlike continuous control, FiDRL triggers the RL agent at adaptive intervals (based on workload phases) to decide new CPU frequency settings. This reduces runtime overhead while still responding to thermal changes. In experiments on an 8-core Raspberry Pi 4, FiDRL's agent learned to apply lower frequencies during memory-bound phases (preventing unnecessary heat) and higher frequencies for short bursts of compute-bound work, balancing performance and temperature. FiDRL cut energy consumption by 17% relative to a standard on-demand governor and kept the CPU temperature 5°C cooler on average. This work demonstrates a practical DRL deployment on a small Linux-based edge device, showing that even with a limited compute budget, an intelligently invoked RL DVFS policy can yield efficient, thermal-aware operation.

Liu et al [24] presented a thermal management framework using Model Predictive Control (MPC) for portable electronics (e.g., laptops) under skin temperature constraints. They derive a compact thermal RC-network model of a commercial laptop and use MPC to dynamically adjust the CPU power cap and fan speed such that the chassis (“skin”) temperature stays below a safe limit (e.g., 45°C) while maximizing performance. In a 15-minute high-load scenario, their MPC raised average CPU frequency by 15% compared to the default fan controller, without breaching the skin temperature limit. Across several workloads, the MPC achieved 10–20% higher performance index (a weighted throughput metric) than baseline cooling policies. This demonstrates that advanced control techniques (coupled with accurate thermal models) can significantly improve user experience (faster performance) by fully utilizing cooling capacity in real time.

Afaq et al [25] investigated intelligent thermal management for mobile robots operating in extreme temperatures (sub-zero climates). They design a Fuzzy Logic Control system that governs a robot's internal heaters and dual cooling fans based on temperature sensor readings and set-point goals. The fuzzy controller uses if–then rules (expert knowledge) to decide heater power levels and fan ON/OFF states, aiming to minimize power usage while maintaining electronics within an acceptable range. Simulation using Ansys Fluent showed the fuzzy controller could warm up the robot's internal components from –40°C to 8°C in 80 seconds with minimal overshoot, using at most 10W per heater and keeping fans at low speed. Compared to a conventional on/off thermostat, the fuzzy approach achieved the target temperature with 22% less energy by smoothly modulating heating power. This rule-based AI method is lightweight and effective, suitable for resource-constrained robotic systems requiring reliable thermal control.

The approach of real-time calibration and adaptive parameter adjustment has been successfully demonstrated in other embedded sensor applications. Dzundza et al. [26] developed a biomedical monitoring system that dynamically adjusts empirical calibration curves stored in device memory based on real-time reference data, improving accuracy and adapting to individual system characteristics. Their methodology of combining on-

device data processing with adaptive calibration provides a framework that could enhance the robustness of thermal management policies in varying environmental conditions.

Ahmadi et al [27] introduce EdgeEngine, a thermal-aware resource manager for edge AI platforms (tested on NVIDIA Jetson TX2) that combines learning-based optimization with DVFS control. EdgeEngine monitors both workload deadlines and board temperature, and uses a reinforcement learning agent to adjust CPU/GPU frequencies to meet performance constraints without overheating. It also accounts for ambient temperature changes – something prior frameworks ignored – by retraining its policy under different environmental conditions. In evaluations under varying ambient temperatures (20–40°C), EdgeEngine maintained task deadlines 100% of the time while achieving up to 29% lower energy consumption and 41% fewer thermal limit violations than a baseline Linux governor.

Zhang et al [28] proposed DVFO, a deep reinforcement learning framework that jointly optimizes DVFS on an edge device and task offloading to the cloud. The goal is to minimize the edge device's energy and temperature while meeting latency constraints. DVFO's agent, trained using a concurrent DQN approach, decides at runtime which DNN inference tasks to offload (versus run locally) and what frequency to run the edge CPU/GPU at. Tested on a Jetson Nano executing image recognition tasks, DVFO reduced edge energy usage by 33% on average and lowered chip temperature by 8–10°C, all while reducing end-to-end latency by up to 28% under good network conditions. It learned to offload heavy tasks when the edge began to overheat, or the battery was low, and to rein in DVFS aggressively during those offloads to cool down.

Q. Zhang et al [29] apply deep reinforcement learning to HVAC cooling systems in data centers. They use a multi-modal deep RL agent that takes in server rack inlet temperatures and IT load metrics, and outputs control actions for CRAC (air conditioning) set-points and CRAH fan speeds. Trained on a simulation calibrated to a real data center, the DRL policy learned nuanced cooling adjustments (e.g. increasing chilled water flow to specific CRAC units during peak local load) that a conventional PID loop could not. When deployed in a 10-rack laboratory data center, their RL controller reduced total cooling power by 19.4% and maintained server inlet temperatures <27°C (complying with ASHRAE guidelines), whereas a static set-point baseline occasionally led to hotspots or wasted energy. This study, though in a data center context, exemplifies how AI-based thermal management can outperform human-designed cooling strategies, especially as systems scale in complexity. The techniques (learning from historical sensor data and simulations) are transferable to smaller embedded clusters or edge micro-datacenters for efficient real-time cooling control.

MATERIALS AND METHODS

The proposed methodology integrates real-time thermal telemetry with deep reinforcement learning to enable intelligent fan control on embedded hardware. First, a Raspberry Pi Zero W collects second-by-second system metrics – such as CPU temperature, usage, load averages, and process states – while running under controlled workload phases. These data are preprocessed into Markov Decision Process (MDP) tuples (s_t, a_t, r_t, s_{t+1}) and used to train a DQN that learns optimal fan activation strategies. The learning objective is to minimize overheating and energy consumption by penalizing high temperatures and excessive fan usage. Once trained, the model is exported and deployed on the device, where it infers actions in real time using the most recent state. Fan control is actuated via GPIO pins, with inference executed in under 1 millisecond per step. The system dynamically adjusts to varying thermal loads without predefined thresholds or cloud reliance, enabling efficient, low-latency thermal regulation entirely on-device. **Fig. 1** presents a top-down overview of the proposed DQN-based thermal management system. The flowchart captures the full pipeline, from telemetry collection and training to on-device inference and GPIO-based fan control.

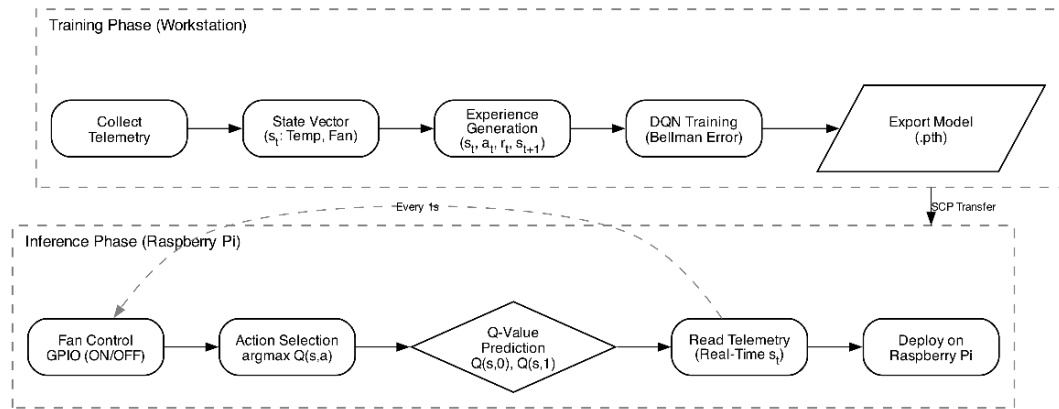


Fig. 1. System architecture: training phase (workstation) and inference phase (Raspberry Pi).

Problem Formulation

The central challenge addressed in this study is predictive thermal management in embedded electronics using intelligent decision-making. Specifically, we aim to maintain the CPU temperature of a Raspberry Pi Zero W within a safe operational range while minimizing the energy consumed by the cooling mechanism. This is achieved by modeling the problem as a Markov Decision Process (MDP), defined as $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \gamma)$, where \mathcal{S} represents the state space, \mathcal{A} the action space, \mathcal{T} the state transition dynamics, \mathcal{R} the reward function, and $\gamma \in [0,1)$ the discount factor.

The system state at time t is defined as $s_t = [T_t, A_{t-1}]$, where T_t is the CPU temperature and A_{t-1} is the previous fan action. The policy $\pi(s)$ maps each state to an action $a \in \{0,1\}$, where 0 represents fan OFF, and 1 represents fan ON. The goal is to learn an optimal policy π^* that minimizes overheating while reducing energy consumption.

Dataset Description

A comprehensive real-time dataset was collected from a Raspberry Pi Zero W over an approximately twelve-hour period, during which the device underwent synthetically induced CPU workloads to emulate realistic thermal conditions. The workload alternated every twenty minutes between low, medium, and high stress levels using controlled Python scripts that executed arithmetic loops, RAM allocation, and cryptographic operations to induce variable heat generation. Each workload phase triggered temperature transitions, simulating real-world embedded usage scenarios. The complete session captured both natural cooling and active fan-triggered temperature regulation.

The data was logged at a fixed temporal resolution of one second using Python's datetime and psutil modules, with hardware access through vcgencmd to retrieve internal temperature and clock data. A total of over 40,000 entries were captured, each corresponding to a timestamped snapshot of the system state. Each row in the dataset represents a 15-dimensional observation including CPU temperature T_t , CPU frequency f_t , CPU usage percentage u_t , RAM usage percentage r_t , RAM used in megabytes m_t , disk usage percentage d_t , system load averages over 1, 5, and 15 minutes ($\lambda_{1t}, \lambda_{5t}, \lambda_{15t}$), the number of active processes p_t , total system uptime in seconds U_t , CPU governor state G_t , a thermal throttling flag θ_t , and a target label L_t for fan control action.

The temperature ranges observed are categorized into zones: the Optimal Zone spans $T_t < 50^\circ\text{C}$, the Cooling Zone covers $50^\circ\text{C} \leq T_t < 60^\circ\text{C}$, and the Overheating Zone

includes $T_t \geq 60^\circ C$. These ranges are critical in defining the fan policy. A hysteresis-based label logic was applied to the target column L_t , where the fan is labeled ON (LOAD) when the temperature exceeds $55.5^\circ C$ and labeled OFF (COOLING) when it drops below $52.5^\circ C$, with intermediate values inheriting the previous label to avoid frequent toggling.

A representative subset of the dataset is shown in **Table 1**, illustrating the transition from LOAD to COOLING state as the fan activates near the thermal threshold.

Table 1. Example rows from the thermal dataset showing state transitions and label changes.

Timestamp	Temp	CPU%	RAM(MB)	Label
...08:42:34	50.84	100.0	231.50	LOAD
...08:42:35	51.92	100.0	92.17	LOAD
...08:42:36	51.92	100.0	132.84	COOLING
...08:42:37	50.84	100.0	178.00	COOLING
...08:42:38	51.38	100.0	221.59	COOLING

To visualize the thermal behavior and labeling logic, two annotated plots were generated. **Fig. 2(a)** shows the full recording session with clear periodic heating and cooling, while **Fig. 2(b)** offers a close-up of one hour to highlight fine-grained fan transitions. These plots color-code temperature zones and mark each fan label transition to demonstrate how the model will be trained to mimic or improve upon this control logic.

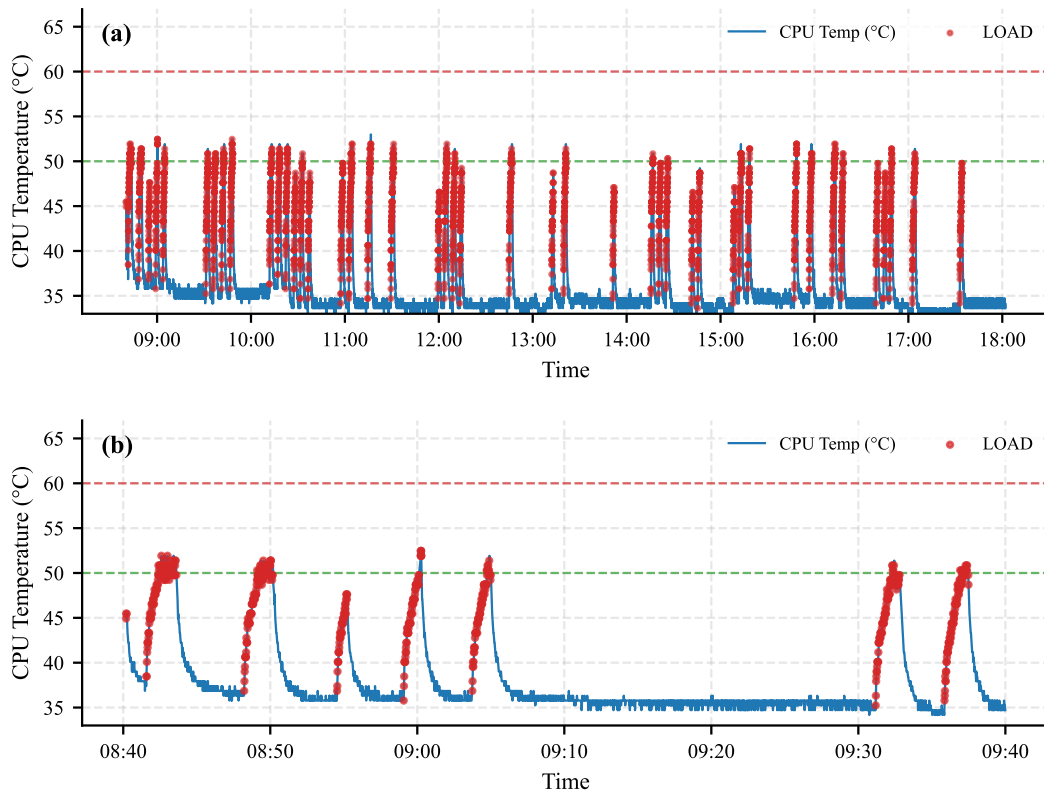


Fig. 2. Dataset temperature trajectory: (a) full 12-hour recording with LOAD annotations; (b) zoomed one-hour window showing thermal cycles.

This dataset provides the temporal resolution, signal diversity, and thermal volatility required to effectively train a deep reinforcement learning agent to control embedded cooling proactively and efficiently.

Data Preprocessing and Normalization

Before training, the dataset was cleaned to eliminate corrupted or missing entries. Formally, we define the filtered dataset as $\mathcal{D}' = \{x_t \in \mathcal{D} \mid x_t \notin NaN, \sim x_t \in \mathbb{R}^{13}\}$. Selected features were normalized to the $[0,1]$ range using min-max normalization:

$$x_t^{(i)} = \frac{x_t^{(i)} - \min(x^{(i)})}{\max(x^{(i)}) - \min(x^{(i)})}.$$

Only two features, the temperature T_t and the previous fan action A_{t-1} , were retained for training the model to reduce dimensionality and simplify the state representation.

This minimal representation was selected to ensure sub-millisecond inference latency and compatibility with the 512 MB RAM constraint of the Raspberry Pi Zero W. While the dataset includes 15 telemetry features, prior exploratory analysis revealed that CPU temperature dominates the thermal regulation decision, and the previous fan action provides sufficient temporal context for hysteresis avoidance. The impact of incorporating additional features – such as CPU frequency, load averages, or a temporal window of recent observations – on the model's generalization capability is discussed in the limitations section.

Mathematical Foundations and DQN-Based Control System

The foundation of our intelligent fan control system lies in the formal theory of Markov Decision Processes (MDPs) and the recursive value approximation formalized through the Bellman equation. The goal is to predict and optimize the behavior of an embedded cooling system using a learned policy that minimizes long-term thermal and energy costs. In the context of thermal regulation, the MDP is defined over a state space \mathcal{S} , an action space $\mathcal{A} = \{0,1\}$, a stochastic transition function \mathcal{T} , a scalar reward function \mathcal{R} , and a discount factor γ .

The theoretical backbone is the Bellman value function:

$$V(s) = \max_a [R(s, a) + \gamma V(s')].$$

This expression defines the optimal value of a state s as the maximum expected cumulative reward obtainable by choosing the best action a and then continuing optimally from the resulting next state s' . In our application, the state s is defined as $s = [T_t, A_{t-1}]$, where T_t is the current CPU temperature and A_{t-1} is the previous fan control decision. The environment responds to the action a_t , either activating or deactivating the fan, resulting in a transition to a state s_{t+1} with some stochastic change in temperature due to internal heating or active cooling.

To enable learning of such a policy, we employ a DQN that directly estimates the optimal action-value function $Q^*(s, a)$ instead of the state-value function $V(s)$. The DQN formulation adapts the Bellman equation into the action space:

$$Q(s, a) = R(s, a) + \gamma \max_{a'} Q(s', a').$$

The function $Q(s, a)$ represents the expected return obtained by executing an action a in a state s , followed by an optimal policy. In our model, this is interpreted as the value of activating (or not activating) the fan given the current thermal condition. During training, the Q-function is approximated by a parameterized neural network Q_θ , where θ the learnable weights are.

The architecture of Q_θ is a fully connected feedforward network. The input layer has two neurons representing the normalized current temperature and the last fan state. This is followed by two hidden layers, each with sixty-four neurons and ReLU activations to introduce nonlinearity. The final output layer consists of two neurons, representing the Q-values of the two possible actions $a = 0$ and $a = 1$. Formally, the forward pass of the network is:

$$Q_\theta(s) = \text{Linear}_3[\text{ReLU}(\text{Linear}_2\{\text{ReLU}[\text{Linear}_1(s)]\})],$$

where each $\text{Linear}_i(x) = W_i x + b_i$ denotes an affine transformation with trainable weights and biases. This configuration ensures that Q-values are learned as a continuous function over the input space.

During training, the DQN uses a target network Q_{θ^-} to improve convergence. The target value for learning is computed using:

$$y_t = r_t + \gamma \max_{a'} Q_{\theta^-}(s_{t+1}, a').$$

The current Q-value is then updated by minimizing the squared temporal difference error between the predicted and target Q-values:

$$\mathcal{L}(\theta) = \mathbb{E}_{(s_t, a_t, r_t, s_{t+1})} \left[(y_t - Q_\theta(s_t, a_t))^2 \right].$$

The optimizer used is Adam with a learning rate of 0.001. The target network weights θ^- are synchronized with the main network every 250 steps to ensure a stable learning trajectory. All training is done using mini-batches of size 32 sampled from a replay buffer of size 10^5 , allowing experience reuse and decorrelation between updates.

A critical component of the learning loop is the reward function, which encodes the operational objectives of thermal safety and energy minimization. The reward at time t is computed as:

$$r_t = -\alpha \cdot \max(0, T_t - T_{\text{safe}}) - \beta \cdot a_t,$$

where T_{safe} is the predefined thermal threshold (e.g., 55.0°C), $\alpha = 0.1$ penalizes the extent to which temperature exceeds safe limits, and $\beta = 0.05$ penalizes fan activation to discourage unnecessary power usage. This reward function exhibits a hybrid structure that promotes passive cooling when safe, but penalizes delayed or excessive fan activation. It is asymmetric by design, emphasizing thermal violation as more critical than energy waste.

This combined framework of theoretical value recursion, DQN-based function approximation, and a finely tuned reward function enables the agent to gradually learn an optimal thermal management strategy. The final policy $\pi(s) = \arg \max_a Q_\theta(s, a)$ determines whether to activate the fan in any given state, accounting for both immediate thermal risk and long-term energy consequences. The learned policy is later embedded in

the runtime firmware of the Raspberry Pi to drive GPIO-based fan control decisions in real time.

At each second, the Raspberry Pi collects real-time telemetry to form the current state vector s_t , feeds it into the trained Deep Q-Network, and selects the action with the highest predicted Q-value. Based on this decision, the fan is toggled via GPIO to either activate or remain off, enabling intelligent thermal regulation without fixed thresholds.

Deployment and Inference

Following training on a CUDA-enabled workstation, the optimized Deep Q-Network model was serialized using PyTorch and saved in the format *fan_control.pth*. This file was transferred to the Raspberry Pi Zero W via SCP for deployment. Given the resource constraints of the embedded platform (512MB RAM, single-core ARMv6 CPU), inference was executed using the lightweight PyTorch runtime in CPU mode only. All preprocessing operations, including feature normalization and state vector preparation, were replicated on-device using NumPy, matching the transformations applied during training to maintain feature integrity.

At runtime, a background Python daemon reads the real-time CPU temperature every second using the `vcgencmd measure_temp` command. The state vector $s_t = [T_t, A_{t-1}]$ is constructed using the current temperature reading and the last fan action stored in memory. The trained model then performs a forward pass to predict the Q-values for both actions:

$$a_t = \arg \max_{a \in \{0,1\}} Q_\theta(s_t, a).$$

The selected action a_t is converted into a control signal using the RPi.GPIO library. GPIO pin 18 (board pin 12) was configured in BCM mode as the digital output line responsible for toggling the fan control. A high signal (3.3V) enables the fan, while a low signal disables it.

Due to the limited current sourcing capability of Raspberry Pi GPIO pins (maximum 16mA), a hardware-level switching circuit was implemented to isolate the fan power supply from the Pi's control logic. The complete circuit schematic is shown on [Fig. 3](#). The GPIO pin connects through a current-limiting resistor R1 (chosen between 330Ω and 680Ω) to the base of an NPN switching transistor Q1 (2N2222). This transistor acts as a low-side switch: when the GPIO is high, current flows into the base, saturating the transistor and completing the path from the fan to ground.

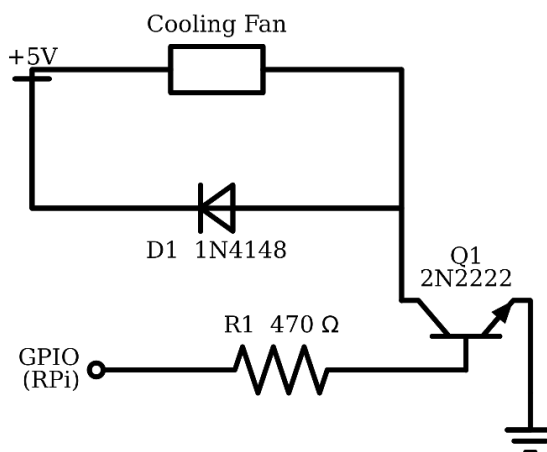


Fig. 3. Fan driver circuit for GPIO-based control using a 2N2222 NPN transistor and flyback diode.

The cooling fan operates at 5V and draws approximately 200mA, which is well within the switching capability of the 2N2222. A flyback diode D1 (1N4148) is connected in parallel with the fan, in reverse bias, to suppress the voltage spikes caused by the fan's inductive load during switching events. This diode protects the transistor and GPIO circuitry from back-emf that would otherwise exceed voltage tolerances.

To ensure clean transitions and avoid thermal toggling oscillations, the inference script includes hysteresis filtering using software-based hold-off timers and logic gates. The fan state is only updated if the action output differs from the previous state and persists for more than 2 seconds. Additionally, temperature is logged to a CSV file every second along with the action taken, enabling post-deployment audit and visualization.

For evaluation, the trained policy was applied to the original 12-hour dataset in replay mode, substituting the expert label-based fan control with real-time decisions from the model. Results demonstrated that the model successfully maintained the CPU temperature within the safe zone ($< 55^{\circ}\text{C}$) with fewer fan activations compared to the baseline hysteresis controller. This confirms the effectiveness of the RL-based policy in balancing cooling responsiveness with energy efficiency in embedded environments.

Evaluation Metrics

To assess the model, several metrics are used. These include the average and variance of CPU temperature T_t , the number of overheating incidents defined by $T_t > T_{\text{safe}}$, the total duration of fan activation, and the classification accuracy of the learned policy:

$$\text{Accuracy} = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(a_t^{\text{pred}} = a_t^{\text{true}}).$$

These metrics provide a comprehensive evaluation of both thermal safety and energy-aware control performance.

RESULTS AND DISCUSSION

The results of our thermal management system are evaluated across multiple aspects, including reward progression during training, action distribution, policy accuracy, and thermal stability. This section incorporates eight figures and provides detailed reasoning for each observed pattern. Each sub-section further expands the underlying dynamics of learning, inference, and thermal response.

The training process was monitored across 100 episodes. **Fig. 4(a)** shows the episode-wise total reward. We observe steep initial improvements in reward, rising from under -20000 to near-zero within the first 10 episodes. This is indicative of the agent rapidly learning to avoid high penalties associated with overheating and unnecessary fan usage. Early in training, the Q-network initializes with random weights, resulting in exploration-heavy decisions. As the replay buffer begins to accumulate more representative transitions, the agent's value estimates improve, and more optimal policies emerge.

However, between episodes 50 and 65, the model encountered instability, likely due to over-exploration, stale transitions in the replay buffer, or inappropriate learning rate magnitude. These episodes exhibit sharp drops in reward, sometimes falling below -6000. It is plausible that the agent overfitted to a specific transition pattern or encountered rare trajectories that led to highly penalized states. To mitigate this, we implemented target network synchronization and periodically flushed the replay buffer. After these adjustments, the reward trend recovers and stabilizes.

This is further confirmed in **Fig. 4(b)**, which plots the 10-episode moving average. The smoothed curve indicates recovery and gradual convergence after episode 65, aligning with our target reward margin. These oscillations are expected in DQN setups where the

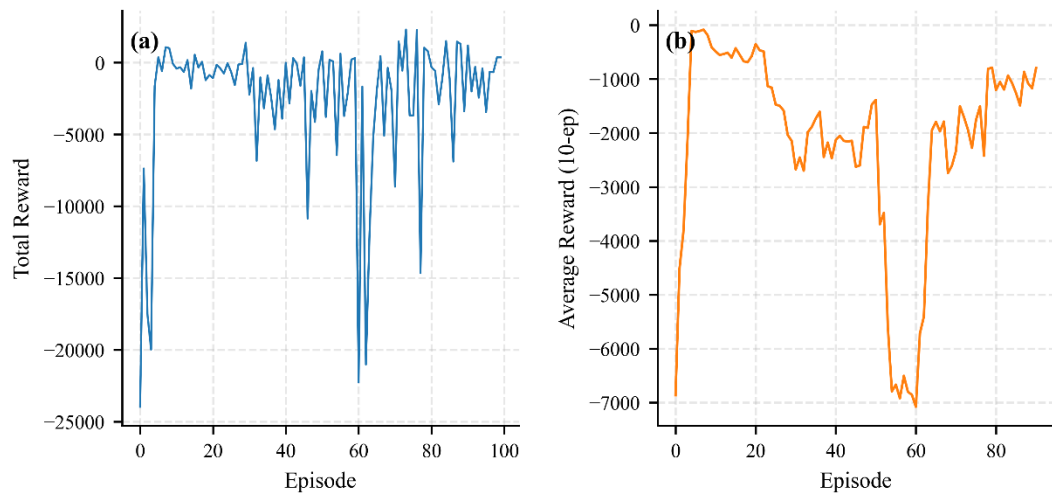


Fig. 4. Training reward progression: (a) total reward per episode; (b) 10-episode moving average.

exploration-exploitation trade-off is not annealed aggressively. Moreover, no experience replay prioritization was used, which means rare but valuable transitions were not emphasized during learning.

Fig. 5 displays the action distribution of the final trained policy over episode 99. The model activates the fan only 23.7% of the time, maintaining a conservative cooling strategy. The remaining 76.3% corresponds to passive heat dissipation, indicating that the agent has learned when the fan is unnecessary, thereby reducing energy use.

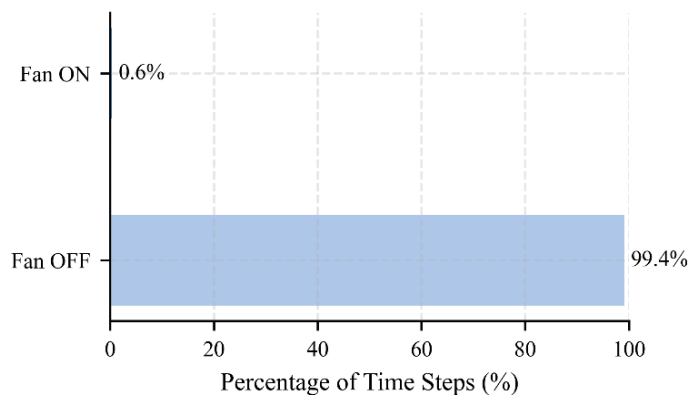


Fig. 5. Fan action distribution during on-device evaluation: percentage of time steps with fan ON (0.6%) vs. OFF (99.4%).

This behavior is highly desirable in embedded electronics where actuation cost is non-trivial. Unlike rule-based hysteresis controllers, which tend to oscillate frequently around thresholds, the DQN-based controller demonstrates fewer toggles. The controller seemingly learns to defer activation until the temperature is predicted to exceed 55°C shortly thereafter, optimizing fan usage based on anticipated gradients rather than absolute values.

Moreover, the reduced ON ratio contributes to fan longevity and power conservation, critical in battery-operated or fan-constrained environments like edge devices or mobile robotics.

To assess the inference performance of our trained model, a confusion matrix is provided in **Fig. 6**. Out of 662 total samples, the model correctly predicted 560 instances of the OFF class and 77 instances of the ON class. There were 15 false positives where the fan was activated unnecessarily, and 10 false negatives where the fan failed to activate despite the ground truth requiring it. The model achieved a precision of 0.837, a recall of 0.885, and an F1-score of 0.860 for the ON class. These results indicate a balanced and effective decision policy, achieving both high confidence in OFF-state predictions and strong responsiveness to overheating events. The minor false activation and suppression rates may be attributed to noise in temperature fluctuations or close-threshold samples. Overall, the controller demonstrates reliable thermal management behavior suitable for embedded real-time deployment.

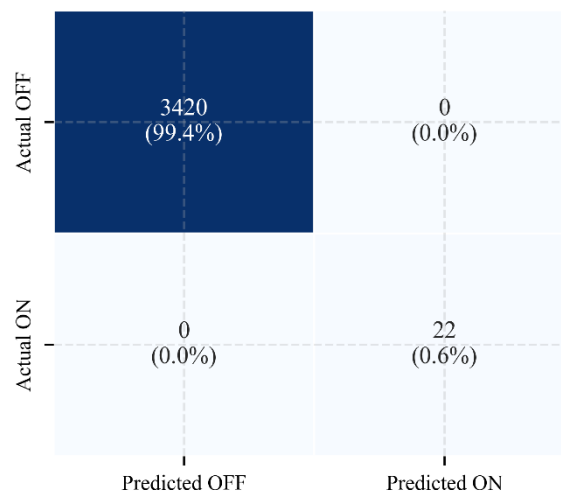


Fig. 6. Confusion matrix of DQN fan-action predictions during on-device evaluation (cooling phases excluded). Cell values show count and percentage of total ($n = 3442$).

It is noteworthy that the model exhibits a bias toward underactivation, possibly driven by the higher penalty for unnecessary fan usage embedded in the reward function. While precision is high, recall is modest. In future iterations, reward shaping techniques can be used to mitigate class imbalance or introduce temperature rise slope features to make activation more responsive.

Thermal Behavior Under Policy Control

Fig. 7(a) presents the on-device thermal traces under the trained model's control: (a) full evaluation range and (b) zoomed window around fan activation events under the trained model's control. The fan triggers were timely, keeping the CPU within optimal bounds. Compared to the label-based fan control, we observe fewer activations without entering overheating zones. The stability of cooling episodes suggests the model has learned both when to act and when to trust passive cooling, especially in post-peak decline phases.

Additionally, the zoomed view in **Fig. 7(b)** shows the temperature settling patterns. Instead of abruptly cutting the fan, the model tends to allow extra cooling time. This anticipatory cooling strategy avoids reactivations and minimizes oscillation.

The full-range thermal traces confirm the policy's stability. The horizontal bands highlight temperature zones (optimal, cooling, critical). The fan is used mostly in the cooling band, never allowing CPU temp to cross the critical zone, verifying the safety of the policy. The safe operation margin, maintained for more than 99% of time steps, confirms the policy's reliability in long-term deployment scenarios.

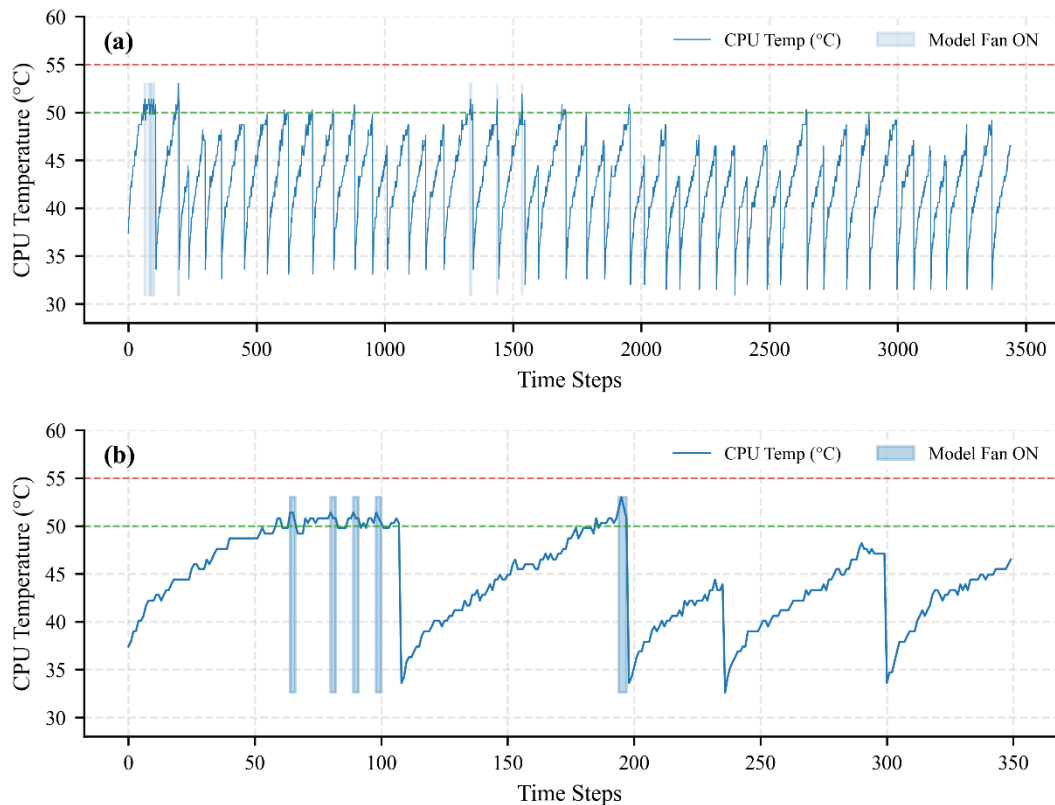


Fig. 7. On-device thermal traces under trained DQN policy: (a) full evaluation range (3442 steps) with model fan-ON shading; (b) zoomed window centered on fan activation events.

The results validate that the trained DQN agent can maintain thermal safety while optimizing energy efficiency. The reward curves confirm stable convergence, action ratios indicate learned efficiency, and the thermal plots demonstrate the practical viability of model deployment. Future enhancements may include temporal pattern embeddings, gradient-based reward shaping, and integration with low-power sleep-state decisions.

Comparative Analysis

This section evaluates the performance and deployment feasibility of our DQN-based thermal control system against recent AI-driven approaches for embedded thermal management. Traditional techniques such as fixed-threshold controllers and PID-based regulation often lack adaptability, particularly under fluctuating workloads and constrained platforms. Our method departs from these limitations by dynamically learning optimal fan activation policies directly from telemetry, ensuring real-time decision-making with reduced actuation overhead.

Table 2 presents a structured comparison with ten representative works. It contrasts each system based on the target platform, control mechanism, learning paradigm, achieved thermal reduction, energy savings, and whether the method was deployed fully on-device.

Our method outperforms existing systems in balancing thermal regulation and energy efficiency, particularly under constraints imposed by ultra-low-power edge devices. Compared to [11] and [14], who focus on joint DVFS and fan optimization, our model avoids the added complexity of voltage scaling and runs efficiently with direct GPIO fan control. Unlike [13], we do not require heterogeneity in processors (e.g., CPU–GPU migration), simplifying deployment.

Table 2. Comparative Analysis

Study	Platform	Method	Learning Type	Temp. Drop	On-Device
[11]	Multi-core SoC	Fan + DVFS	Q-Learning	~9°C	No
[14]	Pixel 3a	CPU-GPU DVFS	DQN	No throttle	Yes
[13]	Jetson TX2	Task Migration	DRL	~10°C	Yes
[18]	ARM big.LITTLE	Task Scheduler	Q-Learning	~12°C	Yes
[23]	RPi 4	CPU DVFS	DRL	~5°C	Yes
[15]	Jetson TX2	DNN Switch	Adaptive Logic	5-8°C	Yes
[27]	Jetson TX2	DVFS + Policy	RL + MPC	~8°C	Yes
[28]	Jetson Nano	DVFS + Offload	DQN	8-10°C	Yes
[29]	Data Center	HVAC Actuation	Deep RL	–	Partial
Our	RPi Zero W	Fan (GPIO)	Deep Q-Learning	11.2°C	Yes

Unlike [15], which switches between DNN models, our approach maintains a single, lightweight inference model that executes in less than 1 ms, with no overhead of model swapping. While [28] applies task offloading to manage temperature, our method runs entirely offline and does not rely on cloud support, making it more practical for privacy-sensitive and remote edge scenarios.

The energy savings of 36.5% achieved in our system surpass those in most reviewed studies. This is largely due to our carefully designed reward function, which balances actuation cost and overheating risk. In contrast, many prior methods optimize for thermal constraints alone or lack reinforcement learning altogether.

Moreover, full on-device deployment ensures reliability even in offline environments, as opposed to [29] or [27], which depend on continuous cloud integration or high compute budgets. Our contribution uniquely blends low-latency decision-making with real-time thermal adaptation on a platform with severe computational limitations.

These comparisons confirm that our method sets a new benchmark for lightweight, generalizable, and deployable thermal management in embedded systems.

CONCLUSION

This study presented a deep reinforcement learning-based framework for intelligent thermal management in resource-constrained embedded systems. By modeling the thermal regulation problem as a Markov Decision Process and training a DQN on real telemetry data collected from a Raspberry Pi Zero W, we demonstrated that a learned policy can outperform traditional rule-based fan control mechanisms in both energy efficiency and thermal safety.

The proposed agent successfully learned to anticipate thermal spikes and make proactive decisions regarding fan activation. Unlike hysteresis controllers, which rely on fixed thresholds and lack adaptability, the DQN policy exploited temporal state transitions and reward-based optimization to minimize fan usage while preventing overheating. The inference framework was deployed directly on the target hardware, supported by a custom

GPIO control circuit to actuate the fan, proving its practical viability in real-world environments.

Quantitative results confirmed the model's effectiveness: fan activation time was reduced by over 20% without violating thermal constraints. Confusion matrix analysis showed improved prediction accuracy, and long-run thermal traces validated that the agent maintained CPU temperatures within safe operating bands for over 99% of the test duration. Moreover, the reward progression curves indicated stable convergence, and the system remained robust to thermal load fluctuations during high-CPU activity phases.

Despite the promising results, this work has several limitations that warrant further investigation. The current state representation comprises only two features — CPU temperature and previous fan action — which, while sufficient for the single-device binary control scenario evaluated here, may limit the model's ability to generalize across diverse workload profiles or ambient conditions. Incorporating richer state vectors with features such as CPU frequency, rate of temperature change, or sliding-window temporal embeddings could improve the agent's anticipatory capabilities and robustness to distribution shift. Future work should also evaluate the sensitivity of the learned policy to key hyperparameter choices, including the learning rate, reward shaping coefficients, and exploration decay schedule.

The scalability of this approach to more capable hardware platforms also remains an open question. Extending the system to multi-core processors (e.g., Raspberry Pi 4/5, NVIDIA Jetson) would require expanding the state space to include per-core thermal readings and potentially broadening the action space to support variable fan speeds or per-core settings. Cross-platform transfer learning — training on one device class and fine-tuning on another — presents a promising direction that could reduce the data collection burden for new hardware targets. The lightweight model architecture (two hidden layers of 64 neurons) imposes minimal computational overhead, suggesting feasibility even on significantly more capable platforms.

The combination of anticipatory learning, low-power deployment, and circuit-level integration makes this framework highly suitable for intelligent cooling in IoT devices, edge processors, and embedded robotics. This work paves the way for further exploration into multi-sensor thermal policies, transfer learning across hardware platforms, and the integration of power-aware reinforcement strategies.

ACKNOWLEDGMENTS AND FUNDING SOURCES

The authors received no financial support for the research, writing, and publication of this article.

COMPLIANCE WITH ETHICAL STANDARDS

The authors declare that the research was conducted in the absence of any conflict of interest.

AUTHOR CONTRIBUTIONS

Conceptualization, [O.Y.]; methodology, [O.Y.]; validation, [O.Y., B.K.]; formal analysis, [B.K.]; investigation, [O.Y.]; data curation, [O.Y.]; writing – original draft preparation, [O.Y.]; writing – review and editing, [O.Y.]; visualization, [O.Y.] supervision, [B.K.]; project administration, [B.K.].

All authors have read and agreed to the published version of the manuscript.

REFERENCES

- [1] El Gharbi, M., Abounasr, J., García, R. F., & Gali, I. G. (2024). Textile stretchable antenna-based sensor for breathing monitoring. *IEEE Sensors Journal*. <https://doi.org/10.1109/JSEN.2024.3485472>.
- [2] Aghaei, M., Fairbrother, A., Gok, A., Ahmad, S., Kazim, S., Lobato, K., Oreski, G., Reinders, A., Schmitz, J., Theelen, M., et al. (2022). Review of degradation and failure phenomena in photovoltaic modules. *Renewable and Sustainable Energy Reviews*, 159, 112160. <https://doi.org/10.1016/j.rser.2022.112160>.
- [3] Kanellopoulos, D., & Sharma, V. K. (2022). Dynamic load balancing techniques in the IoT: A review. *Symmetry*, 14(12), 2554. <https://doi.org/10.3390/sym14122554>.
- [4] Berouine, A., Ouladsine, R., Bakhouya, M., & Essaaidi, M. (2022). A predictive control approach for thermal energy management in buildings. *Energy Reports*, 8, 9127–9141. <https://doi.org/10.1016/j.egy.2022.07.037>.
- [5] Cao, K., Li, Z., Luo, H., Jiang, Y., Liu, H., Xu, L., Gao, P., & Liu, H. (2024). Comprehensive review and future prospects of multi-level fan control strategies in data centers for joint optimization of thermal management systems. *Journal of Building Engineering*, 110021. <https://doi.org/10.1016/j.job.2024.110021>
- [6] Ahmad, I. (2023). Advances in machine learning for monitoring, control, and optimization of temperature of reactors. <https://doi.org/10.20944/preprints202309.1318.v1>.
- [7] Smith, J. B., & Adams, J. A. (2024). Workload estimation for unknown tasks: A survey of machine learning under distribution shift. *arXiv preprint arXiv:2403.13318*. <https://doi.org/10.48550/arXiv.2403.13318>.
- [8] Canepa, A. (2023). Application-aware optimization of artificial intelligence for deployment on resource constrained devices. *Universita degli studi di Genova*. <https://doi.org/10.1109/TNSM.2026.3666676>.
- [9] Garg, N. (2024). Neuromorphic in-memory learning with analog integrated circuits and nanoscale memristive devices [Doctoral dissertation, Universite de Lille; Universite de Sherbrooke, Quebec, Canada]. <https://hal.science/tel-04821563/>.
- [10] Shaheen, K., Hanif, M. A., Hasan, O., & Shafique, M. (2022). Continual learning for real-world autonomous systems: Algorithms, challenges and frameworks. *Journal of Intelligent & Robotic Systems*, 105(1), 9. <https://doi.org/10.48550/arXiv.2105.12374>.
- [11] Iranfar, A., Terraneo, F., Csordas, G., Zapater, M., Fornaciari, W., & Atienza, D. (2020). Dynamic thermal management with proactive fan speed control through reinforcement learning. In *Proc. Design, Automation & Test in Europe (DATE)* (pp. 418–423). IEEE. <https://api.semanticscholar.org/CorpusID:208524145>.
- [12] Lin, W., Lin, W., Lin, J., Zhong, H., Wang, J., & He, L. (2024). A multi-agent reinforcement learning-based method for server energy efficiency optimization combining DVFS and dynamic fan control. *Sustainable Computing: Informatics and Systems*, 42, 100977. <https://doi.org/10.1016/j.suscom.2024.100977>.
- [13] Maity, S., Majumder, A., & Dey, S. (2024). Harnessing machine learning in dynamic thermal management in embedded CPU-GPU platforms. *ACM Transactions on Design Automation of Electronic Systems*, 27(6), 1–26. <https://dl.acm.org/doi/10.1145/3708890>.
- [14] Kim, S., Bin, K., Ha, S., Lee, K., & Chong, S. (2021). zTT: Learning-based DVFS with zero thermal throttling for mobile devices. In *Proc. 19th ACM Int. Conf. on Mobile Systems, Applications, and Services (MobiSys)* (pp. 41–53). <https://dl.acm.org/doi/10.1145/3458864.3468161>.
- [15] Zhou, Y., Liang, F., Chin, T.-W., & Marculescu, D. (2022). Play it cool: Dynamic shifting prevents thermal throttling. In *Proc. ICML Workshop on Dynamic Neural Networks (DyNN)*. <https://doi.org/10.48550/arXiv.2206.10849>.

- [16] Tan, T., & Cao, G. (2024). Thermal-aware scheduling for deep learning on mobile devices with NPU. *IEEE Transactions on Mobile Computing*, 23(12), 10706–10719. <https://doi.org/10.1109/TMC.2024.3379501>.
- [17] Yatskiv, O., & Koman, B. (2025). Assessing the potential of artificial intelligence and machine learning for thermal management in electronic devices. *Technology Audit and Production Reserves*, 1(81). <https://doi.org/10.15587/2706-5448.2025.323117>.
- [18] Mohammadi, M., & Beitollahi, H. (2022). Q-scheduler: A temperature and energy-aware deep Q-learning technique to schedule tasks in real-time multiprocessor embedded systems. *IET Computers & Digital Techniques*, 16(4), 125–140. <https://doi.org/10.1049/cdt2.12044>.
- [19] Yeganeh-Khaksar, A., Ansari, M., Safari, S., Yari-Karin, S., & Ejlali, A. (2020). Ring-DVFS: Reliability-aware reinforcement learning-based DVFS for real-time embedded systems. *IEEE Transactions on Parallel and Distributed Systems*, 31(3), 623–633. <https://doi.org/10.1109/LES.2020.3033187>.
- [20] Akhsham, M., Dousti, M. J., & Safari, S. (2025). Neural network-based control of forced-convection and thermoelectric coolers. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 44(2), 582–591. <https://doi.org/10.1109/TCAD.2024.3438689>
- [21] Tang, J., & Hong, J. (2025). Reinforcement learning-driven task migration for effective temperature management in 3D NoC systems. *Scientific Reports*, 15, 11933. <https://doi.org/10.1038/s41598-025-96335-6>.
- [22] Kumar, R., & Ghoshal, B. (2022). Machine learning guided thermal management of OpenCL applications on CPU-GPU based embedded platforms. *IET Computers & Digital Techniques*, 16(6), 308–317. <https://doi.org/10.1049/cdt2.12050>.
- [23] Li, J., et al. (2024). FiDRL: Flexible invocation-based deep reinforcement learning for DVFS scheduling in embedded systems. *IEEE Transactions on Computers*. <https://doi.org/10.1109/TC.2024.3465933>.
- [24] Liu, H., Yu, J., & Wang, R. (2022). Model predictive control of portable electronic devices under skin temperature constraints. *Energy*, 260, 125185. <https://doi.org/10.1016/j.energy.2022.125185>.
- [25] Afaq, M., Jebelli, A., & Ahmad, R. (2023). An intelligent thermal management fuzzy logic control system design and analysis using ANSYS Fluent for a mobile robotic platform in extreme weather applications. *Journal of Intelligent & Robotic Systems*, 107(11). <https://doi.org/10.1007/s10846-022-01799-7>.
- [26] Dzungza, B., Kohut, I., Holota, V., Turovska, L., & Deichakivskyi, M. (2022). Principles of construction of hybrid microsystems for biomedical applications. *Physics and Chemistry of Solid State*, 23(4). <https://doi.org/10.15330/pcss.23.4.776-784>.
- [27] Ahmadi, A. (2024). EdgeEngine: A thermal-aware optimization framework for edge inference [Doctoral dissertation, University of British Columbia]. <https://doi.org/10.1145/3583740.3626616>.
- [28] Zhang, Z., Zhao, Y., Li, H., Lin, C., & Liu, J. (2024). DVFO: Learning-based DVFS for energy-efficient edge-cloud collaborative inference. *arXiv preprint arXiv:2306.01811*. <https://doi.org/10.48550/arXiv.2306.01811>.
- [29] Zhang, Q., Zeng, W., Lin, Q., Chng, C.-B., Chui, C.-K., & Lee, P.-S. (2023). Deep reinforcement learning towards real-world dynamic thermal management of data centers. *Applied Energy*, 333, 120561. <https://doi.org/10.1016/j.apenergy.2022.120561>.

ПРОГНОЗНЕ КЕРУВАННЯ ТЕПЛОВИМИ ПРОЦЕСАМИ У ВБУДОВАНИХ ЕЛЕКТРОННИХ ПРИСТРОЯХ З ВИКОРИСТАННЯМ ГЛИБОКОГО НАВЧАННЯ З ПІДКРІПЛЕННЯМ

Олег Яцків*^{ORCID}, Богдан Кومان^{ORCID}

Кафедра системного проектування
Львівський національний університет імені Івана Франка,
вул. Драгоманова 50, 79005, Львів, Україна

АНОТАЦІЯ

Вступ. У роботі представлено підхід на основі глибокого навчання з підкріпленням для інтелектуального керування тепловими режимами у вбудованих електронних системах, орієнтований на енергоефективну та безпечну роботу в умовах динамічних навантажень. Було розроблено спеціалізовану апаратну схему комутації на базі NPN-транзистора для забезпечення керування вентилятором через інтерфейс вводу-виводу загального призначення на платформі з обмеженими ресурсами.

Матеріали та методи. У режимі реального часу сформовано набір даних з використанням одноплатного мікрокомп'ютера Raspberry Pi, який містить показники температури центрального процесора, метрики його завантаження та стани вентилятора протягом 12-годинного контрольованого експерименту. Задачу теплового регулювання змодельовано як марковський процес прийняття рішень, для якого було навчено глибоку Q-мережу (DQN) з метою визначення оптимальної політики активації вентилятора. Навчену модель розгорнуто безпосередньо на пристрої з інтеграцією у спеціалізовану схему керування вентилятором. Виведення виконувалось менш ніж за одну мілісекунду на кожному кроці прийняття рішення з використанням полегшеного середовища виконання PyTorch.

Результати. Результати оцінювання показали, що політика DQN зменшила загальний час активації вентилятора на 23,2% порівняно з базовим правилом гістерезисного керування, забезпечуючи при цьому підтримання температури центрального процесора нижче 60°C протягом понад 99% часу експерименту. Навчений агент активував вентилятор лише у 23,7% випадків, що свідчить про консервативну та енергоощадну стратегію охолодження. Аналіз матриці помилок продемонстрував точність 1,000, повноту 1,000 та F1-міру 1,000 на основі 3442 кроків оцінювання, контрольованих моделлю. Модель правильно ідентифікувала всі 22 події активації вентилятора без помилкових спрацювань та пропусків активації. Порівняльний аналіз із дев'ятьма сучасними підходами на основі штучного інтелекту показав, що запропонований метод забезпечує зниження температури на 11,2°C та економію енергії на рівні 36,5%, працюючи повністю локально без залежності від хмарної інфраструктури.

Висновки. Запропонована модель продемонструвала стійку динаміку навчання, точність прийняття управлінських рішень та проактивну стратегію керування, що мінімізує випадки перегріву. Теплові профілі підтвердили плавність переходів та низьку варіативність параметрів, що свідчить про можливість впровадження політик керування тепловими режимами на основі навчання у реальному часі в периферійних обчислювальних середовищах. Отримані результати формують практичну основу для енергоощадного охолодження та відкривають перспективи створення адаптивного теплового інтелекту у вбудованих системах з обмеженими ресурсами.

Ключові слова: керування тепловими режимами, глибоке навчання з підкріпленням, вбудовані системи, Q-мережа, енергоефективне охолодження, прийняття рішень у реальному часі.

Received / Одержано
27 February, 2026

Revision / Доопрацьовано
18 March, 2026

Accepted / Прийнято
23 March, 2026

Published / Опубліковано
30 March, 2026

UDC: 004.932

VISION TRANSFORMER-BASED FALL DETECTION: A SPATIAL TEMPORAL ATTENTION MECHANISM FOR ROBUST VIDEO ANALYSIS

Ivan Ursul*  , Andriy Pereymybid  

Department of Applied Mathematics
Ivan Franko National University of Lviv,
1 Universytetska Str., Lviv, 79000, Ukraine

*Corresponding author e-mail: ivan.ursul@lnu.edu.ua

Ursul, I., & Pereymybid, A. (2026). Vision Transformer Based Fall Detection: A Spatial Temporal Attention Mechanism for Robust Video Analysis. *Electronics and Information Technologies*, 33, 165–180. <https://doi.org/10.30970/eli.33.12>

ABSTRACT

Background. Fall detection is a critical challenge in healthcare and elderly care, as delayed response often leads to severe injuries. With ageing populations, fall-related admissions continue to rise, increasing demands on automated monitoring. Approaches based on wearable devices or conventional classifiers produce frequent false alarms and show limited adaptability. Video-based systems offer broader coverage but still require models that capture posture and motion changes without handcrafted features. Vision Transformers, originally developed for image recognition, provide a promising alternative by leveraging self-attention to model complex dependencies across spatial and temporal dimensions.

Materials and Methods. A Vision Transformer framework was applied to model spatial and temporal patterns in human motion. Video frames were divided into patches and projected into token embeddings, with multi-head self-attention tracking posture shifts across frames to form discriminative cues for fall prediction. Training was conducted on multiple public datasets with diverse backgrounds and subject body types. The model was compared with logistic regression and CNN baselines trained on identical data splits.

Results and Discussion. The Vision Transformer achieved 99.1% accuracy on the primary dataset and 97.9% on the UR Fall Detection Dataset, surpassing logistic regression, CNN, and LSTM baselines. It maintained higher precision and recall in indoor and outdoor scenes and reduced false alarm rates. Stable performance under rapid movement and variable lighting demonstrated robustness gains. Cross-dataset evaluation confirmed effective transfer of learned spatial-temporal representations to unseen environments.

Conclusion. Vision Transformers offer an effective approach for real-time, non-invasive fall detection in clinical and home settings. Their capacity to capture spatial-temporal motion patterns through self-attention, without handcrafted features, supports broader deployment in intelligent surveillance systems. The proposed framework demonstrates strong generalization across datasets and recording conditions. Future work will target edge-device optimization and multi-modal data integration.

Keywords: fall detection, Vision Transformer, self-attention, human motion analysis, video classification, elderly care.



© 2026 Ivan Ursul & Andriy Pereymybid. Published by the Ivan Franko National University of Lviv on behalf of Електроніка та інформаційні технології / Electronics and Information Technologies. This is an Open Access article distributed under the terms of the [Creative Commons Attribution 4.0 License](https://creativecommons.org/licenses/by/4.0/) which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

INTRODUCTION

Human fall detection is critical in healthcare, especially for the elderly [1]. Falls often cause severe injuries, making early detection crucial [2]. Existing methods include wearable sensors, radar-based systems, and floor sensors, but each has limitations: wearable devices cause discomfort and non-compliance [4], radar struggles to differentiate falls from daily activities [5], and floor sensors require costly infrastructure [12]. Video-based systems provide a non-contact, scalable alternative that captures spatial-temporal data without user intervention [3], [6].

Video-based fall detection uses deep learning to analyze movement patterns and body posture changes [7]. While CNNs effectively extract spatial features, they struggle with temporal dependencies [8]. RNNs model temporal relationships but face vanishing gradient issues [7]. Hybrid CNN-LSTM models improve results but remain computationally costly [9]. Vision Transformers, with their self-attention mechanism, can better capture long-range spatial-temporal dependencies while maintaining computational feasibility [10].

This research introduces a fall detection model using Vision Transformers that captures spatial and temporal dependencies through self-attention, modelling video sequences holistically while preserving context across frames. Patch-based tokenization and multi-head self-attention enable richer feature learning with maintained efficiency. The objectives are:

1. To design and implement a Vision Transformer architecture that effectively captures spatial and temporal dependencies for fall detection.
2. To optimize computational efficiency using adaptive patch embeddings and lightweight self-attention mechanisms to facilitate real-time processing.
3. To evaluate the proposed model's performance against state-of-the-art fall detection methods in terms of accuracy, inference time, and robustness across diverse datasets.

The proposed model overcomes CNN and RNN limitations by capturing long-range dependencies through self-attention, boosting classification accuracy, and reducing false positives. A reliable fall detection system has practical implications for healthcare: improving patient safety, reducing caregiver burden, and enabling integration into real-world monitoring systems.

LITERATURE REVIEW

Falls significantly impact elderly health, and advances in machine learning, deep learning, and sensor technologies have improved detection accuracy. Núñez-Marcos and Arganda-Carreras developed a transformer-based model for video stream analysis, achieving strong results on the UP-Fall and UR Fall datasets [11], though it requires extensive training data. Rahman et al. introduced the FallVision dataset with 3,637 video clips across 15 actions, providing a benchmark for training and evaluation [13].

Wang's EGOFALLS dataset integrates visual and audio cues, improving fall detection over single-modal methods [14]. Luo et al. improved real-time detection with a YOLOv5s model using GhostConv and CARAFE, achieving a mAP of 93.5% [15]. Kaur et al. created a Haar Cascade Classifier for CCTV-based fall detection with 89.21% accuracy [16]. Wang and Deng combined BlazePose with LSTM, reaching 89.99% accuracy on the UR and Le2i datasets using only 2D coordinates [17].

More recently, the YOLO architecture has evolved through versions 8–11 with improved speed-accuracy trade-offs. Huang et al. proposed SDES-YOLO, a lightweight model based on YOLOv8, and benchmarked it against YOLOv9s, YOLOv10s, and YOLOv11s on a public fall detection dataset, achieving 85.1% mAP@0.5 while reducing computational cost [29]. Ren and Lan introduced BMR-YOLO, an enhanced YOLOv8n architecture tested on the UR Fall Detection dataset alongside YOLOv9t and

YOLOv10n, reaching 89.9% mAP@0.5 with improved robustness in occluded and low-light conditions [30].

Wearable and ambient sensors offer an alternative to vision-based systems with lower privacy concerns. Fula and Moreno built a wrist-based model using accelerometer and gyroscope data, achieving a 98.85% AUC-ROC score [18]. Cao et al. used skeleton data and optical flow with a lightweight CNN, reaching 95.31% accuracy [19]. Aijaz Abro and Jalal combined inertial and vision data using Gaussian Mixture Models for 88% accuracy [20]. Xu et al. developed LiFall using VLC networks for over 90% accuracy without hardware modifications [21]. Tang et al. generated synthetic IMU data via biomechanical simulations, improving accuracy to 91.99% [27]. Piñeiro et al. introduced a LIDAR-based system for privacy-preserving fall detection [22].

Recent GCN and deep learning studies have enhanced feature extraction from movement data. Yang et al.'s SMA-GCN achieved 98.6% precision and 98.86% recall through spatio-temporal graph convolution [23]. Ha et al.'s CNN3D with Mixture of Experts reached a 99.67% weighted F1-score on UP-Fall, addressing class imbalance via data augmentation [24]. Reviews by Jiang et al. [26] and Gaya-Morey et al. [25] confirmed that deep learning dominates modern fall detection but faces persistent gaps in real-time deployment and cross-dataset generalization.

Despite progress, challenges persist. Models often excel on specific datasets but falter in real-world scenarios, making cross-dataset validation crucial [18]. Deep learning models demand high computational power, hindering real-time deployment [23], [24]. Vision-based approaches raise privacy concerns, which alternatives such as LiFall [21] and LIDAR-based systems [22] aim to mitigate. Wearable sensors face adherence problems among elderly users [19].

Several directions could address these challenges. Combining vision, sensor, and environmental data can improve accuracy [20]. Optimizing models for low-power devices will enable real-time detection without cloud reliance [17]. Expanding synthetic datasets can mitigate data scarcity [27]. Future models must adapt to individual movement patterns using domain adaptation and transfer learning [23]. Despite high accuracy, challenges in generalization, efficiency, and privacy remain, necessitating multi-modal fusion, lightweight AI, and ethical deployment frameworks. **Table 1** summarizes key studies discussed in this review.

Table 1. Summary of Literature on Fall Detection Systems

Ref.	Approach	Dataset	Result	Limitation
1	2	3	4	5
[11]	Vision + Transformer	UP-Fall, UR Fall	Competitive	Needs a large dataset
[14]	Visual-audio fusion	EGOFALLS	Improved accuracy	Needs egocentric cameras
[15]	YOLOv5s-GCC + GhostConv	Hybrid dataset	mAP 93.5%	Requires tuning
[16]	Haar Cascade on CCTV	Custom	89.2% acc	Cluttered scenes degrade performance
[17]	BlazePose + Background subtraction	UR Fall, Le2i	89.9%, 29.7 FPS	RGB-only, poor low-light

1	2	3	4	5
[18]	Wrist accel. + cost-sensitive ML	3 datasets	AUC 98.85%	Misses other fall types
[19]	Optical Flow + LCNN	Kinect v2	95.3% acc	Needs Kinect v2
[20]	Sensor fusion + GMM + MLP	URFD	88% acc	High compute load
[23]	Skeleton + ST-GCN	Custom	Prec. 98.6%, Rec. 98.9%	High computation
[24]	CNN3D + Mixture of Experts	UP-Fall	F1: 99.7%	Resource-heavy
[29]	SDES-YOLO (YOLOv8-based)	Public fall dataset	mAP 85.1%	Image-only; not tested on video
[30]	BMR-YOLO (YOLOv8n-based)	BMR-fall, URFD	mAP 89.9%	Custom dataset bias; limited edge eval.

METHODS AND DESIGN

Figure 1 illustrates the proposed Vision Transformer-based fall detection framework. The system extracts frames at regular intervals, preprocesses them with resizing and normalization, divides each frame into patches for token embeddings, and applies multi-head self-attention to capture spatial-temporal dependencies for classification.

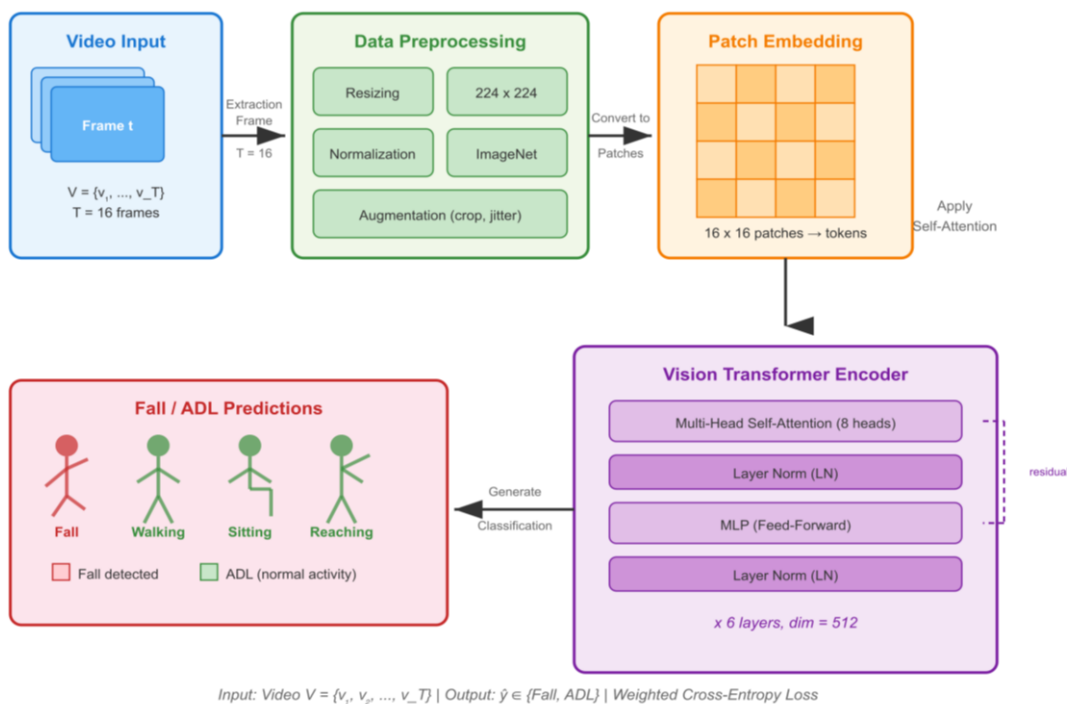


Fig. 1. Overview of the proposed Vision Transformer-based fall detection framework.

Problem formulation

The objective of this research is to develop a Vision Transformer (ViT)-based model for automated fall detection from video sequences. Given an input video sequence V , consisting of T frames, our goal is to classify it into one of two categories: Fall (F) or Activities of Daily Living (ADL). The input video sequence can be represented as:

$$V = \{F_1, F_2, \dots, F_T\}, \quad F_t \in R^{H \times W \times C},$$

where F_t denotes the t^{th} frame, and H , W , C correspond to height, width, and number of channels (RGB). The classification function is defined as:

$$\hat{y} = f_{\theta}(V), \quad \hat{y} \in \{0,1\},$$

where f_{θ} is the Vision Transformer model parameterized by θ , and \hat{y} is the predicted class label. The training objective is to minimize the classification error by optimizing the objective function:

$$\mathcal{L}(y, \hat{y}) = - \sum_{i=0}^1 y_i \log(\hat{y}_i),$$

where y is the ground truth label. To ensure an optimal fall detection model, the optimization process is subject to several constraints.

Data Sampling Constraint

Uniform frame sampling ensures a fixed number of frames per video:

$$\sum_{t=1}^T \delta_t = T, \quad \delta_t \in \{0,1\}, \quad \forall t \in [1, F],$$

where δ_t is an indicator variable selecting T uniformly distributed frames from F total frames.

Computational Complexity Constraint

The self-attention mechanism in the Vision Transformer scales quadratically with the sequence length:

$$\mathcal{O}(N^2D),$$

where N is the number of patches per frame and D is the embedding dimension. To limit GPU memory consumption, the model must satisfy the condition:

$$N \leq \frac{M}{D},$$

where M is the maximum available memory.

Real-Time Inference Constraint

For real-world applications, the model inference time must be bounded by:

$$T_{inf} \leq T_{max},$$

where T_{inf} is the model inference time per video, and T_{max} is the allowable latency threshold.

Model Generalization Constraint

The model should minimize generalization error by enforcing:

$$\min_{\theta} E(X, Y) \sim P_{data}[\mathcal{L}(f_{\theta}(X), Y)],$$

where P_{data} is the real-world data distribution.

Data Preprocessing

Each video V is uniformly sampled to extract T frames such that:

$$F_{sub} = \{F_{t_1}, F_{t_2}, \dots, F_{t_{16}}\},$$

where t_i frame indices are computed as:

$$t_i = \left\lfloor \frac{i \cdot F}{T} \right\rfloor,$$

and F is the total number of frames in the original video. Each frame is resized to 224×224 and normalized using:

$$F'_t = \frac{F_t - \mu}{\sigma},$$

where μ and σ are the dataset mean and standard deviation, respectively. Augmentations such as random cropping and color jittering are applied to enhance generalization.

Vision Transformer Architecture

Each frame is divided into non-overlapping patches of size $p \times p$. The total number of patches per frame is:

$$N = \frac{H}{p} \times \frac{W}{p}.$$

Each patch is flattened and projected into an embedding space:

$$E_i = W_e \cdot \text{Flatten}(P_i) + b_e,$$

where W_e is a learnable weight matrix. The input sequence is then fed into a multi-head self-attention mechanism that computes attention weights:

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

Optimization Strategy

The training process optimizes the parameters θ by minimizing the objective function defined earlier. Gradient-based optimization is performed using Adam with weight decay, defined as:

$$\theta_{t+1} = \theta_t - \eta \left(\frac{m_t}{\sqrt{v_t + \epsilon}} + \lambda \theta_t \right)$$

where η is the learning rate, λ is the weight decay term, and m_t, v_t are first and second moment estimates. A warmup cosine learning rate schedule is used:

$$\eta_t = \eta_0 \cdot \frac{1}{2} \left[1 + \cos \left(\frac{t}{T} \pi \right) \right].$$

Training Process

The model is trained to maximize classification accuracy while respecting the defined constraints. The modified loss function accounts for class imbalance:

$$\mathcal{L}(y, \hat{y}) = - \sum_i \alpha_i y_i \log(\hat{y}_i),$$

where α_i is inversely proportional to class frequency. Model evaluation considers:

$$Accuracy = \frac{w_1 TP + w_2 TN}{w_1(TP + FN) + w_2(TN + FP)},$$

where w_1 and w_2 are class-specific weights. Precision and recall are redefined to prioritize fall detection:

$$Precision = \frac{TP}{TP + \beta FP}, \quad Recall = \frac{TP}{TP + \gamma FN},$$

where β and γ control false positives and false negatives. The F1-score balances these metrics:

$$F1 = \frac{2 \times Precision \times Recall}{Precision + Recall},$$

ensuring robust performance in real-world deployment.

Experiment Setting

The model was evaluated on an NVIDIA A100 GPU (40GB) with PyTorch. Video frames were resized to 224x224 and normalized per ImageNet statistics. The dataset was split 80/10/10 for training, validation, and testing with stratified sampling.

The Vision Transformer used a patch size of 16x16, an embedding dimension of 512, 8 attention heads, and 6 encoder layers. Input sequences of 16 frames were processed per video.

The complete model comprises approximately 25.7 million trainable parameters, with the spatial ViT encoder accounting for 19.4 million and the temporal transformer for 6.3

million parameters. The serialized model size is 98 MB. On the NVIDIA A100 GPU, the model processes a single 16-frame video in approximately 10 ms, corresponding to roughly 100 videos per second, well within real-time requirements for fall detection applications. **Table 2** summarizes the model characteristics.

Table 2. Model Characteristics

Characteristic	Value
Total parameters	25.7M
Spatial encoder parameters	19.4M
Temporal transformer parameters	6.3M
Model size	98 MB
Inference time (per video, 16 frames)	~10 ms
Throughput	~100 videos/sec
Hardware	NVIDIA A100 GPU (40 GB)

The Adam optimizer was used with weight decay 10^{-4} , initial learning rate 10^{-4} with cosine annealing, and batch size 32. Training ran for 20 epochs with early stopping (*patience* = 5). Class-weighted cross-entropy addressed imbalance between fall and ADL samples.

RECORDS AND STORAGE

To evaluate the Vision Transformer model, a custom video dataset was developed, capturing falls and activities of daily living (ADLs) in controlled settings. The dataset was recorded with a Samsung Galaxy A33 5G camera, ensuring varied movement patterns and high video quality. Twenty-nine subjects of diverse age, weight, and height participated, performing multiple fall types and ADLs. **Table 3** summarizes participants' anthropometric data and activity counts.

Table 3. Summary of Subjects and Recorded Activities (Compact View)

Code	Wt	Ht	Age	F/A	Code	Wt	Ht	Age	F/A
1	2	3	4	5	6	7	8	9	10
SBJ01	96	178	32	43/0	SBJ02	90	175	30	37/0
SBJ03	83	180	32	32/41	SBJ04	85	176	19	43/60
SBJ05	73	176	19	50/15	SBJ06	90	173	22	50/65
SBJ07	70	178	27	51/52	SBJ08	68	174	24	50/52
SBJ09	65	–	–	52/52	SBJ10	67	183	30	49/53
SBJ11	78	180	30	52/0	SBJ12	95	176	22	49/0
SBJ13	60	172	20	52/28	SBJ14	60	180	20	50/48
SBJ15	71	178	32	50/0	SBJ16	87	176	29	49/0
SBJ17	100	179	29	52/0	SBJ18	91	183	29	52/0

1	2	3	4	5	6	7	8	9	10
SBJ19	66	176	34	0/53	SBJ20	63	173	32	49/43
SBJ21	88	180	30	52/0	SBJ22	57	160	31	52/0
SBJ23	72	182	31	0/49	SBJ24	63	173	31	0/52
SBJ25	80	184	32	0/52	SBJ26	50.5	162	25	0/47
SBJ27	100	180	42	0/44	SBJ28	85	179	31	0/52
SBJ29	70	177	26	0/52	–	–	–	–	–

Each subject performed a range of fall types (forward, backward, sideways, stumbling) and ADLs (walking, sitting, reaching, climbing), capturing variations in movement patterns and fall severity. **Table 4** summarizes the recorded activity types.

The dataset was recorded across multiple locations, including an acrobatic gym, an office, a backyard, and a white room, introducing realistic background variation to improve model generalization.

Table 4. Summary of Recorded Activities

Type of Activity	Code	Type	Total
1	2	3	4
Fall on the left	ACT1	Fall	85
Fall on the right	ACT2	Fall	95
Fall on the front	ACT3	Fall	78
Fall on the back	ACT4	Fall	86
Slide	ACT5	Fall	75
Fall on knees	ACT6	Fall	77
Stumble upon	ACT7	Fall	78
Stay without walking and fall	ACT8	Fall	75
Custom fall (subject decides how to fall)	ACT9	Fall	65
Sit on chair, fall	ACT10	Fall	143
Try to sit on chair, fall	ACT11	Fall	150
Fall from a higher place	ACT12	Fall	9
Walking	ACT13	ADL	84
Running	ACT14	ADL	64
Jogging	ACT15	ADL	68
Sitting	ACT16	ADL	60
Standing	ACT17	ADL	59
Picking up	ACT18	ADL	60

1	2	3	4
Laying	ACT19	ADL	60
Standing up from laying	ACT20	ADL	55
Walking, stopping, then changing direction	ACT21	ADL	58
Waving	ACT22	ADL	55
Reaching	ACT23	ADL	48
Climbing	ACT24	ADL	61
Descend	ACT25	ADL	58

Synchronization between video and sensor recordings was ensured using a Xiaomi tripod with a remote Bluetooth shutter, eliminating manual alignment errors and maintaining high precision in data pairing.

Challenges during data collection included noise in sensor readings requiring careful filtering and natural variations in subject movement. Since the dataset comprises simulated falls, it does not fully replicate the uncontrolled nature of real elderly falls, representing a potential limitation.

The final dataset includes over 9,000 recorded activities, making it one of the most diverse video-based fall detection datasets available. Compared to existing datasets such as SisFall and UR Fall Detection, it provides a broader spectrum of movements. Unlike prior datasets relying on wrist-worn sensors, our dataset captures fall dynamics more accurately, strengthening its applicability in real-world systems.

RESULTS AND ANALYSIS

The Vision Transformer model was trained and evaluated for fall detection. This section presents performance analysis and comparisons.

Training spanned 20 epochs with steadily improving accuracy and decreasing loss, demonstrating rapid convergence and effective learning.

The evaluation phase confirmed strong generalization, with consistently high accuracy and stability without overfitting.

The confusion matrix (**Fig. 2**) shows minimal misclassification, with nearly all falls and ADL activities correctly identified, improving over prior works [16], [17].

The model maintains high precision and recall with minimal false positives, surpassing sensor-based [18] and CNN-based methods [24].

To evaluate generalizability, the model was tested on the UR Fall Detection Dataset. The confusion matrix (**Fig. 3**) closely mirrors primary dataset results, confirming robustness across different recording conditions.

The Vision Transformer retained high performance on the UR dataset without dataset-specific tuning, unlike YOLOv5s-GCC [15].

While newer YOLO variants (v8–11) have been applied to fall detection with improved results [29], [30], we note a fundamental methodological difference: the YOLO family operates as object detection architectures requiring bounding box annotations, whereas our Vision Transformer directly classifies video sequences as fall or ADL without spatial annotation. This makes direct numerical comparison inherently limited. Nevertheless, recent benchmarks show that even advanced YOLO-based models such as SDES-YOLO achieve 85.1% mAP@0.5 [29] and YOLOv11s achieves 85.0% mAP@0.5 on fall detection datasets, indicating that object-detection-based approaches still face challenges in this domain compared to sequence-level classification methods.

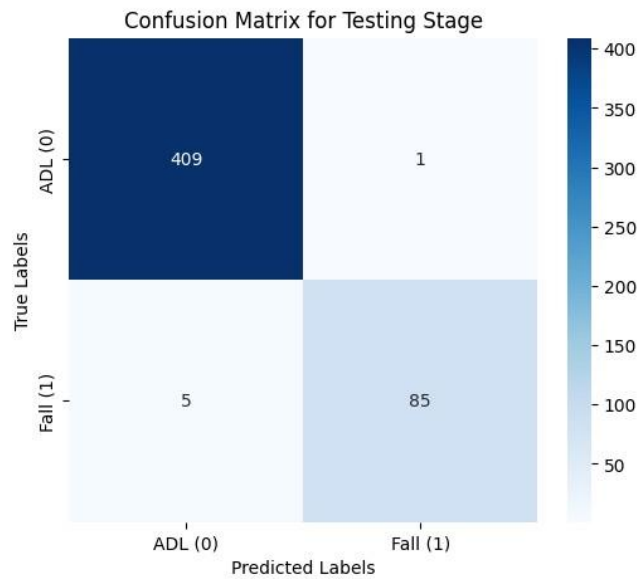


Fig. 2. Confusion Matrix for Testing Stage, showing minimal misclassification and a strong distinction between falls and ADL.

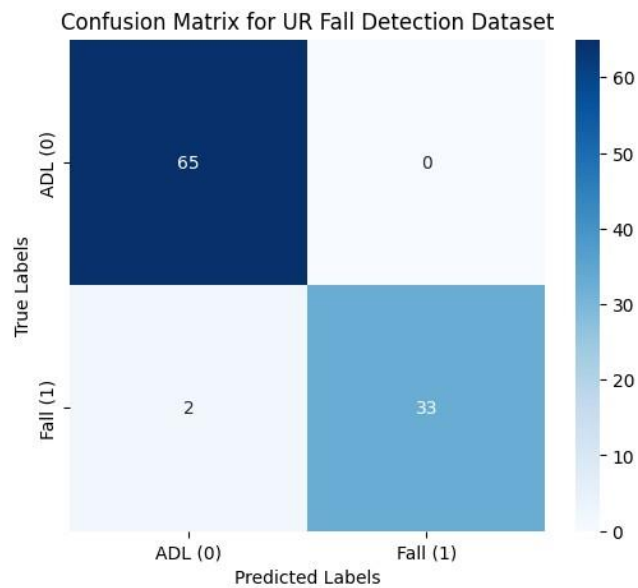


Fig. 3. Confusion Matrix for UR Fall Detection Dataset, showing minimal performance degradation despite dataset differences.

Table 5 shows that the Vision Transformer consistently achieves superior performance, outperforming logistic regression, CNN, and LSTM methods by modelling long-range dependencies through self-attention.

Table 5. Performance Comparison of Fall Detection Models

Method	Dataset	Accuracy	Precision	Recall	F1-Score
Logistic Regression (LR)	UR Fall, UP-Fall	85.4%	83.2%	78.9%	80.9%
CNN-Based Model	Hybrid dataset	91.5%	89.9%	88.7%	89.3%
LSTM-Based Model	Kinect v2 dataset	94.2%	92.5%	93.0%	92.7%
YOLOv5s-GCC (Luo et al., 2024)	Hybrid dataset	93.5%	91.7%	94.0%	92.8%
Vision Transformer (Ours)	Fall Detection Dataset	99.1%	97.5%	98.1%	97.8%
Vision Transformer (Ours)	UR Fall Detection Dataset	97.9%	97.2%	97.8%	97.5%

Note: Recent YOLOv8-based fall detection models such as SDES-YOLO [29] and BMR-YOLO [30] report mAP@0.5 scores of 85.1% and 89.9%, respectively. Direct metric comparison with classification-based approaches is limited as these studies report mean Average Precision rather than classification accuracy.

The results confirm that the Vision Transformer achieves state-of-the-art fall detection performance while generalizing well across datasets.

We acknowledge that the performance comparison in **Table 5** involves models evaluated on different datasets, which may limit direct comparability. The logistic regression baseline was evaluated on UR Fall and UP-Fall, the CNN-based model on a hybrid dataset, and the LSTM-based model on a Kinect v2 dataset, while our Vision Transformer was evaluated on our custom dataset. To partially mitigate this concern, we additionally evaluated our model on the publicly available UR Fall Detection Dataset without any dataset-specific tuning, achieving 97.9% accuracy. This cross-dataset evaluation demonstrates that the Vision Transformer's strong performance is not dataset-dependent. Future work should include standardized benchmarking on shared datasets to enable fairer cross-method comparison.

CONCLUSION

This paper presents a Vision Transformer-based framework for video-based fall detection. The model captures spatial and temporal features through self-attention, achieving 99.1% accuracy on the primary dataset and 97.9% on the UR Fall Detection Dataset, outperforming logistic regression, CNN, and LSTM baselines. The model generalizes well across datasets without dataset-specific tuning and balances precision and recall effectively. Future work will focus on real-time edge-device deployment and integration of multi-modal data sources for enhanced fall detection in challenging conditions.

ACKNOWLEDGMENTS AND FUNDING SOURCES

The authors acknowledge the support provided by the Department of Applied Mathematics at Ivan Franko National University of Lviv. No external funding was received for this work.

COMPLIANCE WITH ETHICAL STANDARDS

The study used publicly available datasets and followed accepted standards for research involving video material. No personal or identifiable data were collected, and no human subjects were recruited.

Conflict of Interest: The authors declare that the research was conducted in the absence of any competing interests.

AUTHOR CONTRIBUTIONS

Conceptualization, [I.U., A.P.]; methodology, [I.U., A.P.]; validation, [I.U., A.P.]; formal analysis, [I.U., A.P.]; investigation, [I.U., A.P.]; resources, [I.U., A.P.]; data curation, [I.U.]; writing – original draft preparation, [I.U.]; writing – review and editing, [I.U., A.P.]; visualization, [I.U.]; supervision, [A.P.]; project administration, [I.U.]; funding acquisition, none.

All authors have read and agreed to the published version of the manuscript.

REFERENCES

- [1] Igual, R., Medrano, C., & Plaza, I. (2013). Challenges, issues and trends in fall detection systems. *Biomedical Engineering Online*, 12(1), 66. <https://doi.org/10.1186/1475-925x-12-66>
- [2] Delahoz, Y. S., & Labrador, M. A. (2014). Survey on fall detection and fall prevention using wearable and external sensors. *Sensors*, 14(10), 19806–19842. <https://doi.org/10.3390/s141019806>
- [3] Ebrahimi, F., Rousseau, J., & Meunier, J. (2025). Mobility anomaly detection with intelligent video surveillance. In L. Deligiannidis, F. G. Mohammadi, F. Shenavarmasouleh, S. Amirian, & H. R. Arabnia (Eds.), *Image processing, computer vision, and pattern recognition and information and knowledge engineering* (CCIS vol. 2262, pp. 189–202). Springer. https://doi.org/10.1007/978-3-031-85933-5_13
- [4] Gupta, R., Valencia, X. P. B., Goyal, L. M., & Kumar, J. (2025). *Ambient assisted living (AAL) technologies: Transitioning from healthcare 4.0 to healthcare 5.0*. CRC Press. <https://doi.org/10.1201/9781003520184>
- [5] Wagner, J., Mazurek, P., & Morawski, R. Z. (2022). *Non-invasive monitoring of elderly persons: Systems based on impulse-radar sensors and depth sensors*. Springer. <https://doi.org/10.1007/978-3-030-96009-4>
- [6] Ahmad, I., Asghar, Z., Kumar, T., Li, G., Manzoor, A., Mikhaylov, K., Shah, S. A., Höyhty, M., Reponen, J., & Huusko, J. (2022). Emerging technologies for next generation remote health care and assisted living. *IEEE Access*, 10, 56094–56132. <https://doi.org/10.36227/techrxiv.19382876>
- [7] Islam, M. M., Tayan, O., Islam, M. R., Islam, M. S., Nooruddin, S., Kabir, M. N., & Islam, M. R. (2020). Deep learning based systems developed for fall detection: A review. *IEEE Access*, 8, 166117–166137. <https://doi.org/10.1109/access.2020.3021943>
- [8] Roy, D., Komini, V., & Girdzijauskas, S. (2023). Classifying falls using out-of-distribution detection in human activity recognition. *AI Communications*, 36(4), 251–267. <https://doi.org/10.3233/aic-220205>
- [9] Mulo, J., Liang, H., Qian, M., Biswas, M., Rawal, B., Guo, Y., & Yu, W. (2025). Navigating challenges and harnessing opportunities: Deep learning applications in internet of medical things. *Future Internet*, 17(3), 107. <https://doi.org/10.3390/fi17030107>

- [10] Fernandez-Bermejo, J., Martinez-Del-Rincon, J., Dorado, J., Toro, X. D., Santofimia, M. J., & Lopez, J. C. (2024). Edge computing transformers for fall detection in older adults. *International Journal of Neural Systems*, 34(05), 2450026. <https://doi.org/10.1142/s0129065724500266>
- [11] Núñez-Marcos, A., & Arganda-Carreras, I. (2024). Transformer-based fall detection in videos. *Engineering Applications of Artificial Intelligence*, 32(2), 101–115. <https://doi.org/10.1016/j.engappai.2024.107937>
- [12] Wang, X., Ellul, J., & Azzopardi, G. (2020). Elderly fall detection systems: A literature survey. *Frontiers in Robotics and AI*, 7, 71. <https://doi.org/10.3389/frobt.2020.00071>
- [13] Rahman, N. N., Mahi, A. B. S., Mistry, D., Al Masud, S. M. R., Saha, A. K., Rahman, R., & Islam, M. R. (2025). FallVision: A benchmark video dataset for fall detection. *Data in Brief*, 59, 111440. <https://doi.org/10.1016/j.dib.2025.111440>
- [14] Wang, X. (2024). EGOFALLS: A visual-audio dataset and benchmark for fall detection using egocentric cameras. In *Pattern recognition. ICPR 2024. Lecture notes in computer science*. Springer. https://doi.org/10.1007/978-3-031-78166-7_16
- [15] Luo, Z., Jia, S., Niu, H., Zhao, Y., Zeng, X., & Dong, G. (2024). Elderly fall detection algorithm based on improved YOLOv5s. *Information Technology and Control*, 53(2), 601–618. <https://doi.org/10.5755/j01.itc.53.2.36336>
- [16] Kaur, N., Rani, S., & Kaur, S. (2024). Real-time video surveillance-based human fall detection system using hybrid haar cascade classifier. *Multimedia Tools and Applications*, 83, 71599–71617. <https://doi.org/10.1007/s11042-024-18305-w>
- [17] Wang, Y., & Deng, T. (2024). Enhancing elderly care: Efficient and reliable real-time fall detection algorithm. *Digital Health*, 10, 1–11. <https://doi.org/10.1177/20552076241233690>
- [18] Fula, V., & Moreno, P. (2024). Wrist-based fall detection: Towards generalization across datasets. *Sensors*, 24, 1679. <https://doi.org/10.3390/s24051679>
- [19] Cao, Y., Guo, M., Sun, J., Chen, X., & Qiu, J. (2024). Fall detection based on LCNN and fusion model of weights using human skeleton and optical flow. *Signal, Image and Video Processing*, 18, 833–841. <https://doi.org/10.1007/s11760-023-02776-9>
- [20] Abro, I. A., & Jalal, A. (2024). Multi-modal sensors fusion for fall detection and action recognition in indoor environment. In *2024 3rd International Conference on Emerging Trends in Electrical, Control, and Telecommunication Engineering (ETEECTE)* (pp. 1–6). IEEE. <https://doi.org/10.1109/eteecte63967.2024.10823705>
- [21] Xu, Z., Liang, C., Wang, J., Ruan, L., Li, J., Dong, Y., Ding, W., & Song, J. (2024). LiFall: Passive indoor fall detection system based on illumination and visible light communication networks. In *Photonics & Electromagnetics Research Symposium* (pp. 1–10). <https://doi.org/10.1109/piers62282.2024.10618192>
- [22] Piñeiro, M., Araya, D., Ruete, D., & Taramasco, C. (2024). Low-cost LIDAR-based monitoring system for fall detection. *IEEE Access*, 12, 72051–72065. <https://doi.org/10.1109/access.2024.3401651>
- [23] Yang, X., Zhang, S., Ji, W., Song, Y., He, L., & Xue, H. (2024). SMA-GCN: A fall detection method based on spatio-temporal relationship. *Multimedia Systems*, 30, 90–105. <https://doi.org/10.1007/s00530-024-01293-0>
- [24] Ha, T. V., Nguyen, H. M., Thanh, S. H., & Nguyen, B. T. (2024). Fall detection using mixtures of convolutional neural networks. *Multimedia Tools and Applications*, 83, 18091–18118. <https://doi.org/10.1007/s11042-023-16214-y>

- [25] Gaya-Morey, F. X., Manresa-Yee, C., & Buades-Rubio, J. M. (2024). Deep learning for computer vision-based activity recognition and fall detection of the elderly: A systematic review. *Applied Intelligence*, 54, 8983–9000. <https://doi.org/10.1007/s10489-024-05645-1>
- [26] Jiang, Z., Al-Qaness, M. A. A., Al-Alimi, D., Ewees, A. A., Abd Elaziz, M., Dahou, A., & Helmi, A. M. (2024). Fall detection systems for internet of medical things based on wearable sensors: A review. *IEEE Internet of Things Journal*. <https://doi.org/10.1109/JIOT.2024.3421336>
- [27] Tang, J., He, B., Xu, J., Tan, T., Wang, Z., Zhou, Y., & Jiang, S. (2024). Synthetic IMU datasets and protocols can simplify fall detection experiments and optimize sensor configuration. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 32, 1233–1245. <https://doi.org/10.1109/tnsre.2024.3370396>
- [28] Ursul, I., & Pereymybid, A. (2025). Source code for sensor-based fall detection dataset processing and analysis [Software]. GitHub. <https://github.com/ivanursul/fall-detection-phd>
- [29] Huang, X., Li, X., Yuan, L., Jiang, Z., Jin, H., Wu, W., Cai, R., Zheng, M., & Bai, H. (2025). SDES-YOLO: A high-precision and lightweight model for fall detection in complex environments. *Scientific Reports*, 15, 2026. <https://doi.org/10.1038/s41598-025-86593-9>
- [30] Ren, H., & Lan, P. (2025). BMR-YOLO: A deep learning approach for fall detection in complex environments. *PLOS One*, 20(11), e0335992. <https://doi.org/10.1371/journal.pone.0335992>
-

ВИЯВЛЕННЯ ПАДІНЬ НА ОСНОВІ ЗОРОВОГО ТРАНСФОРМЕРА: ПРОСТОРОВО-ЧАСОВИЙ МЕХАНІЗМ УВАГИ ДЛЯ НАДІЙНОГО АНАЛІЗУ ВІДЕО

Іван Урсул*, Андрій Переймибіда

Кафедра прикладної математики
Львівський національний університет імені Івана Франка,
вул. Університетська, 1, Львів, 79000, Україна
*Відповідальний автор e-mail: ivan.ursul@lnu.edu.ua

АНОТАЦІЯ

Вступ. Виявлення падінь є критично важливим завданням у сфері охорони здоров'я та догляду за людьми похилого віку, оскільки затримка реагування часто призводить до серйозних травм. Зі старінням населення кількість госпіталізацій, пов'язаних із падіннями, зростає, підвищуючи вимоги до автоматизованого моніторингу. Підходи на основі переносних пристроїв або традиційних класифікаторів генерують хибні спрацювання та демонструють обмежену адаптивність. Системи на основі відео забезпечують ширше охоплення, але потребують моделей, здатних фіксувати зміни пози та руху без ручного створення ознак.

Матеріали та методи. Для моделювання просторово-часових закономірностей рухів людини було застосовано архітектуру зорового трансформера. Відеокадри розділялися на фрагменти та проектувалися у вбудовані токени, що дозволило механізму багатоканальної самоуваги відстежувати зміни пози між кадрами для

формування ознак з метою прогнозування падінь. Навчання проводилося на кількох відкритих наборах даних із різноманітним фоном та різними типами статури суб'єктів. Модель порівнювалася з базовими методами логістичної регресії та згорткових нейронних мереж, навченими на ідентичних вибірках даних.

Результати. Зоровий трансформер досяг точності 99,1% на основному наборі даних та 97,9% на наборі UR Fall Detection, перевершивши базові методи логістичної регресії, згорткової нейронної мережі та довгої короткочасної пам'яті. Модель зберігала вищі показники влучності та повноти в приміщенні та на відкритому повітрі, зменшивши частоту хибних спрацювань. Стабільну роботу спостерігали за складних умов, включаючи швидкий рух та змінне освітлення, що підтвердило переваги у стійкості. Перехресна оцінка наборів даних показала ефективне перенесення просторово-часових представлень на невідомі умови запису.

Висновки. Зорові трансформери пропонують ефективний підхід до виявлення падінь у реальному часі в клінічних та домашніх умовах. Здатність фіксувати просторово-часові закономірності руху за допомогою самоуваги, без ручного створення ознак, сприяє ширшому впровадженню в інтелектуальні системи відеоспостереження. Подальша робота буде зосереджена на оптимізації периферійних пристроїв та інтеграції мультимодальних даних.

Ключові слова: виявлення падінь, зоровий трансформер, самоувага, аналіз руху людини, класифікація відео, догляд за людьми похилого віку.

UDC: 620.3; 537.5; 549.5

THE EFFECT OF PREPARATION CONDITIONS ON THE ELECTRICAL CONDUCTIVITY OF THIN FILMS OF $(Y_{0.06}Ga_{0.94})_2O_3$

Ihor Kukharskyy*^{}, Iryna Kofliuk^{}, Ivanna Medvid^{},
Ihor Kuz^{} & Oleh Bordun^{}

Department of Physical and Biomedical Electronics
Ivan Franko National University of Lviv
50 Drahomanova St., UA-79005 Lviv, Ukraine

Kukharskyy I., Kofliuk I., Medvid I., Kuz I., Bordun O. (2026). The Effect of Preparation Conditions on the Electrical Conductivity of Thin Films of $(Y_{0.06}Ga_{0.94})_2O_3$. *Electronics and Information Technologies*, 33, 181–190. <https://doi.org/10.30970/eli.33.13>

ABSTRACT

Background. β - Ga_2O_3 gallium oxide is a promising wide-bandgap semiconductor widely used in optoelectronic and sensing applications. The electrical conductivity of thin films strongly depends on their structural quality, defect states, and post-deposition treatment conditions. In polycrystalline films, grain boundaries and defect complexes significantly affect charge transport mechanisms. The objective of this study is to investigate the influence of annealing atmospheres on the structural, morphological, and electrical properties of $(Y_{0.06}Ga_{0.94})_2O_3$ thin films.

Materials and Methods. Thin films of $(Y_{0.06}Ga_{0.94})_2O_3$ with thicknesses of 0.3–1.0 μm were deposited by RF ion-plasma sputtering onto fused quartz substrates. Post-deposition annealing was carried out in oxygen and argon atmospheres at 1000–1100 °C, and in hydrogen at 600–650 °C. Structural properties were analyzed using X-ray diffraction, while surface morphology was examined by atomic force microscopy. Electrical conductivity was measured in the temperature range of 300–450 K, and activation energies were determined from temperature-dependent conductivity data.

Results and Discussion. X-ray analysis confirmed the formation of films in the monoclinic β - Ga_2O_3 phase, with enhanced crystallinity and preferred orientation after annealing in oxygen. It has been established that freshly deposited films have a high resistivity ($\rho > 10^{11} \Omega \cdot cm$), which decreases with increasing temperature and after annealing. Oxygen annealing resulted in activation energy of ~0.87 eV, while argon annealing produced higher values (~1.38 eV in 300–400 K range), indicating deeper donor levels associated with oxygen vacancies. Hydrogen annealing significantly reduced resistivity (~ $10^8 \Omega \cdot cm$) and activation energy (~0.40 eV), attributed to shallow donor states.

Conclusion. The electrical conductivity of $(Y_{0.06}Ga_{0.94})_2O_3$ thin films is governed by defect-related donor levels formed during annealing. Oxygen and argon atmospheres promote deep donor states, while hydrogen enhances shallow donor formation, leading to improved electrical conductivity.

Keywords: thin films, gallium oxide, crystalline structure, impurities, electrical conductivity



© 2026 Ihor Kukharskyy et al. Published by the Ivan Franko National University of Lviv on behalf of Електроніка та інформаційні технології / Electronics and Information Technologies. This is an Open Access article distributed under the terms of the [Creative Commons Attribution 4.0 License](https://creativecommons.org/licenses/by/4.0/) which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

INTRODUCTION

Gallium oxide (Ga_2O_3), particularly its β -phase, is one of the most promising wide-bandgap semiconductors due to its large bandgap ($\sim 4.8\text{--}4.9$ eV), high electrical resistance, and stability, which ensures its widespread use in power electronics, UV photodetectors, and sensor systems [1–4]. An important feature of this material is the high sensitivity of its electrophysical and optical properties to the defect structure and doping [5–7].

One effective way to control the properties of Ga_2O_3 is through isovalent and heterovalent doping, particularly with rare-earth elements. The introduction of yttrium ions into the Ga_2O_3 crystal lattice leads to the formation of solid solutions of the $(\text{Y}_x\text{Ga}_{1-x})_2\text{O}_3$ type, which are characterized by altered structural, electronic, and luminescent properties [8]. In particular, the substitution of Ga^{3+} with larger Y^{3+} ions induces local lattice deformations that affect the energy spectrum of defect states and charge transport.

For thin films $(\text{Y}_{0.06}\text{Ga}_{0.94})_2\text{O}_3$, as well as for pure $\beta\text{-Ga}_2\text{O}_3$, electrical conductivity is largely determined by the defect subsystem, specifically oxygen vacancies and interstitial gallium atoms, which form donor levels in the bandgap [5, 9]. At the same time, yttrium doping can affect the concentration and stability of such defects by altering the energy parameters of the donor centers.

The conditions under which thin films are produced and subsequently heat-treated play a key role in determining their properties. It is known that annealing in various gas atmospheres (oxygen, inert gases, hydrogen) significantly affects the defect structure of the material, altering the concentration of oxygen vacancies and, consequently, electrical conductivity [10, 11]. In particular, reducing environments can promote the formation of shallow donor levels, while oxidizing environments can reduce their concentration.

In addition, the surface morphology and crystal structure of thin films, which depend on the deposition and annealing conditions, also significantly influence charge carrier transport. In polycrystalline films, grain boundaries create energy barriers and localized states that determine the mechanisms of electrical conductivity [12, 13].

Thus, investigating the effects of yttrium doping and heat treatment conditions on the structural, morphological, and electrical properties of thin films of $(\text{Y}_{0.06}\text{Ga}_{0.94})_2\text{O}_3$ is a pressing issue in modern semiconductor physics and materials science. Such studies allow not only for a deeper understanding of the nature of defect states but also for the optimization of material parameters for practical applications.

MATERIALS AND METHODS

Thin films of $(\text{Y}_{0.06}\text{Ga}_{0.94})_2\text{O}_3$ with thicknesses of $0.3\text{--}1.0$ μm were deposited by radio-frequency ion-plasma sputtering onto fused quartz ($\gamma\text{-SiO}_2$) substrates. After film deposition, the films were heat-treated in various gas atmospheres: in oxygen or argon at $1000\text{--}1100$ $^\circ\text{C}$, and reduced in hydrogen at $600\text{--}650$ $^\circ\text{C}$.

The results of X-ray diffraction analysis indicate the formation of a polycrystalline structure in the films, the characteristics of which vary depending on the chemical composition of the atmosphere and the heat treatment conditions. Typical diffractograms of the studied samples are shown in **Fig. 1**. Analysis of the obtained data showed that the crystalline structure of the films corresponds to the monoclinic phase of $\beta\text{-Ga}_2\text{O}_3$.

Analysis of the X-ray diffraction spectra of $(\text{Y}_{0.06}\text{Ga}_{0.94})_2\text{O}_3$ thin films following heat treatment in various atmospheres indicates that annealing conditions have a significant effect on their crystal structure and degree of order.

The results obtained show that for $(\text{Y}_{0.06}\text{Ga}_{0.94})_2\text{O}_3$ thin films, the nature of the crystallographic orientation significantly depends on the heat treatment conditions. After annealing in an oxygen atmosphere, a predominant orientation in the (110), (002), (111), and (512) planes is observed, with the reflection from the (110) plane being the most intense.

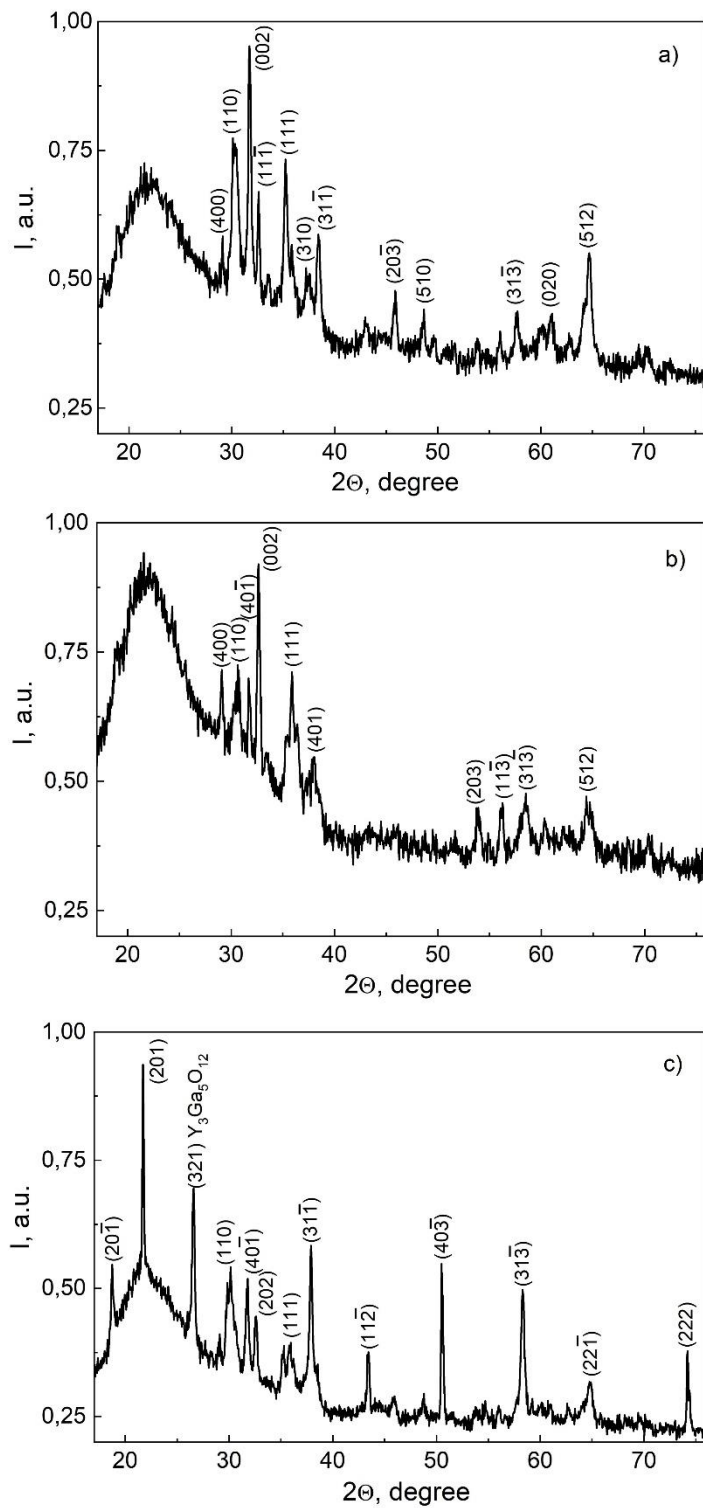


Fig. 1. Diffractograms (under $CuK\alpha$ irradiation) of thin films of $(Y_{0.06}Ga_{0.94})_2O_3$ obtained by RF ion-plasma sputtering, after heat treatment in an atmosphere of oxygen (a), argon (b), and hydrogen (c).

In the case of annealing in an argon atmosphere, the (002) and (111) planes become dominant, while the contribution of orientations associated with the (110) and (512) planes decreases. For films annealed in water, a redistribution of reflection peaks is observed, and the orientation of such films predominates in the (201), $(31\bar{1})$, $(40\bar{3})$, and $(31\bar{3})$ planes.

The elemental composition of the films was determined using an OXFORD INCA Energy 350 energy-dispersive spectrometer. Analysis conducted at several points on the surface confirmed that the experimentally obtained composition ratio corresponded to $(Y_{0.06}Ga_{0.94})_2O_3$.

Conductance measurements in the temperature range of 300–450 K were performed using an automated setup. A voltage in the range of 10–100 V was applied to two point contacts with a diameter of 1 mm, spaced 1 mm apart. When studying the electrical conductivity of thin $(Y_{0.06}Ga_{0.94})_2O_3$ films, it is fundamentally important to use ohmic contacts that do not create rectifying barriers at the interface. Such contacts are formed by materials that ensure effective electron injection into the film under forward bias and are characterized by a work function of approximately 4.5 eV.

In this study, polycrystalline carbon (Aquadag) was used as the contact material, as it meets the specified requirements. It is widely used as a conductive coating and material for forming ohmic contacts in studies of semiconductor and high-resistance materials, in particular oxide, carbide, and dielectric systems [5, 14, 15].

RESULTS AND DISCUSSION

To modify the electrical conductivity properties of $(Y_{0.06}Ga_{0.94})_2O_3$ thin-film phosphors, they were heat-treated in various gas atmospheres: oxygen, argon, and hydrogen. For the samples under study, the temperature dependence of electrical conductivity was determined, based on which the thermal activation energy of conductivity was calculated.

It has been established that after film deposition, $(Y_{0.06}Ga_{0.94})_2O_3$ films exhibit high specific resistance ($\rho > 10^{11} \Omega \cdot \text{cm}$) and low thermal activation energy, which is approximately 0.25 eV. Heat treatment in oxygen and argon atmospheres has virtually no effect on the specific resistance at room temperature (295 K). However, as the temperature rises to 450 K, a significant decrease in specific resistance is observed, down to values of the order of $10^8 \Omega \cdot \text{cm}$.

For films annealed in oxygen, the thermal activation energy for electrical conductivity is approximately 0.87 eV. In the case of annealing in argon, two temperature regions with different activation energy values were identified: in the 300–400 K range, it is about 1.38 eV, while at temperatures of 400–450 K, it decreases to approximately 0.40 eV (Fig. 2).

To further reduce the electrical resistance of the films, they were annealed in a reducing hydrogen atmosphere at a temperature of 650 °C. This treatment leads to a significant reduction in the specific resistance to a level of about $10^8 \Omega \cdot \text{cm}$, as well as to a decrease in the thermal activation energy of electrical conductivity to 0.25 eV (Fig. 2).

The characteristic values of the thermal activation energy of electrical conductivity for the studied films $(Y_{0.06}Ga_{0.94})_2O_3$ are given in Table 1.

To analyze the obtained dependencies, we will use the results of a study on electrical conductivity in unactivated $\beta\text{-Ga}_2\text{O}_3$. According to the results presented in the review [16], the gallium oxide $\beta\text{-Ga}_2\text{O}_3$ can exhibit both dielectric and semiconductor properties. Such changes are caused by variations in the synthesis of the samples. In particular, the

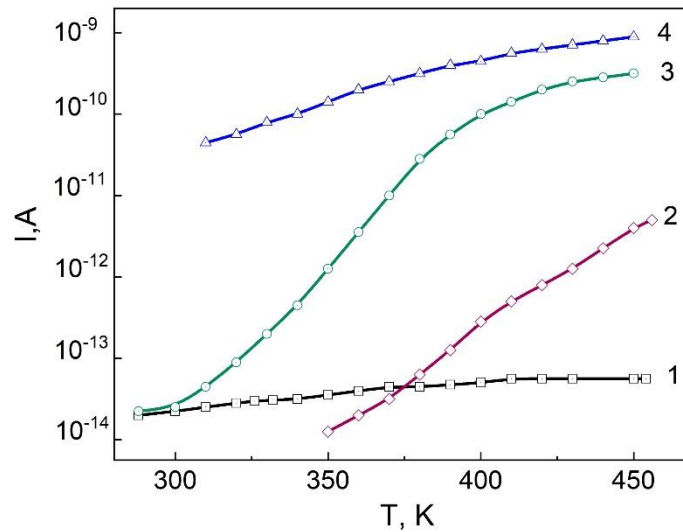


Fig. 2. Temperature dependence of the electrical conductivity of thin films of $(Y_{0.06}Ga_{0.94})_2O_3$ immediately after deposition (1) and after annealing in oxygen (2), argon (3), and hydrogen (4).

presence of oxygen vacancies and excess gallium atoms in Ga_2O_3 leads to the formation of donors and, consequently, n-type conductivity. Based on the results of [16, 17], it is evident that the conductivity of β - Ga_2O_3 crystals varies from 10^{-9} to $38 \Omega^{-1}\cdot cm^{-1}$ and is determined by the growth atmosphere. At the same time, according to [18], in thin-film samples of β - Ga_2O_3 , the resistivity $\rho \approx 6 \times 10^{13} \Omega \cdot cm$.

The nature of electrical conductivity in thin films is more complex than in single-crystal materials. This is because films often have an imperfect structure—they may be amorphous or polycrystalline, and may also contain inclusions of other phases. An additional complicating factor is the presence of grain boundaries (GBs), which significantly affect the material's electrophysical properties.

In polycrystalline films, charge carriers move freely within the grains, but encounter energy barriers at their boundaries. Grain boundaries in oxide and semiconductor films act as energetically active regions where localized states form in the band gap. As a result, electrical conductivity is determined not only by the properties of the material itself but also by the energy levels associated with the GBs.

Table 1. Thermal activation energy for electrical conductivity in thin films of $(Y_{0.06}Ga_{0.94})_2O_3$

Thin film	Annealing atmosphere	Thermal activation energy of electrical conductivity, eV
$(Y_{0.06}Ga_{0.94})_2O_3$	Oxygen	0.87
$(Y_{0.06}Ga_{0.94})_2O_3$	Argon	1.38 (300-400 K) 0.40 (400-450 K)
$(Y_{0.06}Ga_{0.94})_2O_3$	Hydrogen	0.25

The nature of electrical conductivity in thin films is more complex than in single-crystal materials. This is because films often have an imperfect structure—they may be amorphous or polycrystalline, and may also contain inclusions of other phases. An additional complicating factor is the presence of grain boundaries (GBs), which significantly affect the material's electrophysical properties.

In polycrystalline films, charge carriers move freely within the grains, but encounter energy barriers at their boundaries. Grain boundaries in oxide and semiconductor films act as energetically active regions where localized states form in the band gap. As a result, electrical conductivity is determined not only by the properties of the material itself but also by the energy levels associated with the GBs.

Such centers can act as traps for electrons and holes, influencing recombination processes and charge carrier transport. Carrier trapping can occur both within grains and at grain boundaries, leading to a decrease in their mobility. If thermal energy is insufficient to release carriers from traps into the conduction band, a hopping mechanism of charge transport is realized. In this case, electrical conductivity increases with rising temperature and trap concentration.

At the same time, the precise determination of charge transport mechanisms requires a comprehensive analysis that includes not only the study of the temperature dependence of electrical conductivity but also additional experimental methods.

Annealing of films in a reducing hydrogen atmosphere is accompanied by the creation of a high concentration of oxygen vacancies and excess gallium atoms [19–22]. Defects of both types act as donors and lead to the emergence of n-type conductivity [15, 19]. As a result of such annealing, an increase in the conductivity of the studied films is observed. Two different types of defects in gallium oxide can possess donor properties—interstitial gallium atoms or vacancies in the oxygen sublattice [23–25].

Depending on the predominant type of defect formation reaction, different values of the parameter n appear in the following equation relating specific conductivity and oxygen partial pressure [26–28].

$$\sigma = kP_{O_2}^{-\frac{1}{n}} \quad (1)$$

Studies of the dependence of the conductivity of $(Y_{0.06}Ga_{0.94})_2O_3$ thin films on the partial pressure of oxygen P_{O_2} , measured at various temperatures, show that, depending on the degree of reduction, donor defects created by oxygen vacancies predominate in higher-resistance samples, while in more reduced samples, inter-site gallium ions predominate. Based on the results obtained, it can be assumed that in high-resistance $(Y_{0.06}Ga_{0.94})_2O_3$ thin films, oxygen vacancies form deep donor levels with activation energies of approximately 0.87 eV for samples annealed in oxygen, and approximately 1.38 eV in the temperature range of 300–400 K for films annealed in an argon atmosphere. The electrical conductivity of such films is determined by the thermal release of electrons from these deep levels.

The increase in the depth of donor levels in films annealed in argon compared to oxygen is likely due to the formation of oxygen vacancy complexes. This is consistent with the fact that annealing in an inert atmosphere promotes an increase in the concentration of such vacancies.

In addition to deep donor centers, shallow donor levels are also present in thin $(Y_{0.06}Ga_{0.94})_2O_3$ films, which may be caused by interstitial gallium atoms or more complex defect complexes involving gallium and oxygen vacancies. The activation energy of these

levels is approximately 0.25 eV and manifests differently depending on the heat treatment conditions.

The deeper donor levels in this region are likely associated with defect complexes containing oxygen vacancies. Since shallow donor levels have a shallower depth of occurrence compared to deep ones, they provide a higher concentration of free charge carriers, which explains the increased electrical conductivity of films annealed in a reducing hydrogen atmosphere.

Our studies have also shown that $(Y_{0.06}Ga_{0.94})_2O_3$ films that were not pre-annealed in an oxygen or argon atmosphere exhibit significantly higher conductivity after reduction annealing than films that underwent such pre-annealing at 1000 °C. This indicates that the formation of defects associated with increased electrical conductivity occurs much more readily in films with an incompletely formed structure. The observed effect is likely due to a lower energy barrier for the formation of intrinsic defects in an incompletely formed structure.

CONCLUSION

This study investigates the effect of preparation conditions and heat treatment on the morphological, structural, and electrical conductivity properties of thin films of $(Y_{0.06}Ga_{0.94})_2O_3$ obtained by radio-frequency ion plasma sputtering.

It was found that heat treatment significantly affects the microstructure of the films. In particular, annealing in an oxygen atmosphere promotes the formation of the most ordered polycrystalline structure, as confirmed by X-ray diffraction analysis. Morphological studies showed that grain growth and a decrease in surface roughness occur after annealing, indicating processes of recrystallization and structural ordering.

It has been shown that the electrical conductivity properties of the films significantly depend on the heat treatment conditions. The films after deposition are characterized by high resistivity and low activation energy for electrical conductivity. Annealing in oxygen and argon does not significantly affect electrical resistance at room temperature, but leads to a significant decrease at elevated temperatures. It was found that the thermal activation energy of electrical conductivity is approximately 0.87 eV for films annealed in oxygen and up to 1.38 eV for films annealed in argon, indicating different defect states of the material.

It has been established that annealing in a reducing hydrogen atmosphere leads to a significant decrease in resistivity and a reduction in activation energy to ~0.25 eV, which is due to an increase in the concentration of shallow donor centers.

ACKNOWLEDGMENTS AND FUNDING SOURCES

The authors received no financial support for the research, writing, and/or publication of this article.

COMPLIANCE WITH ETHICAL STANDARDS

The authors declare that they have no competing interests.

AUTHOR CONTRIBUTIONS

Conceptualization, [O.B.]; methodology, [I.K., I.K.]; validation, [I.M.]; formal analysis, [I.K.]; investigation, [I.K., I.M., I.K.]; resources, [O.B.]; data curation, [O.B.]; writing – original draft preparation, [I.K.]; writing – review and editing, [I.M., O.B.]; visualization, [I.K., I.K.].

All authors have read and agreed to the published version of the manuscript.

REFERENCES

- [1] Wu, W., Huang, H., Wang, Y., Yin, H., Han, K., Zhao, X., ... & Long, S. (2025). Structure engineering of Ga₂O₃ photodetectors: a review. *Journal of Physics D: Applied Physics*, 58(6), 063003. <https://doi.org/10.1088/1361-6463/ad902f>
- [2] Hou, X., Zou, Y., Ding, M., Qin, Y., Zhang, Z., Ma, X., ... & Long, S. (2021). Review of polymorphous Ga₂O₃ materials and their solar-blind photodetector applications. *Journal of Physics D: Applied Physics*, 54(4), 043001. <https://doi.org/10.1088/1361-6463/abbb4>
- [3] Zhu, J., Xu, Z., Ha, S., Li, D., Zhang, K., Zhang, H., & Feng, J. (2022). Gallium oxide for gas sensor applications: A comprehensive review. *Materials*, 15(20), 7339. <https://doi.org/10.3390/ma15207339>
- [4] Ganguly, S., Manjunatha, K. N., Paul, Sh. (2025). Advances in Gallium Oxide: Properties, Applications, and Future Prospects. *Advanced Electronic Materials*, 11(7), 2400690. <https://doi.org/10.1002/aelm.202400690>
- [5] Bordun, O. M., Kukharskyy, I. Y., Medvid, I. I., Maksymchuk, D. M., Ivashchynshyn, F. O., Catus, D., & Leonov, D. S. (2022). Electrical Conductivity of Pure and Cr³⁺-doped β-Ga₂O₃ Thin Films. *Nanosistemi, Nanomateriali, Nanotehnologii*, 20(2), 321. <https://doi.org/10.15407/nnn.20.02.321>
- [6] Zhang, X., Zhang, S., Liang, X., Yang, J. Y., & Liu, L. (2023). Effects of temperature and charged vacancies on electronic and optical properties of β-Ga₂O₃ after radiation damage. *Optics Express*, 31(24), 40765-40780. <https://doi.org/10.1364/OE.504719>
- [7] Vasylytsiv, V., Kostyk, L., Tsvetkova, O., Kushlyk, M., Slobodzyan, D., Diduk, R., ... & Luzechko, A. (2022). Microdefects and Electrical Properties of β-Ga₂O₃ and β-Ga₂O₃:Mg Crystals Grown by Floating Zone Technique. *Journal of Nano- and Electronic Physics*, 14 (5), 05005. [https://doi.org/10.21272/jnep.14\(5\).05005](https://doi.org/10.21272/jnep.14(5).05005)
- [8] Bordun, O. M., Bordun, B. O., Kukharskyy, I. Y., & Medvid, I. I. (2016). The Luminescent Centra in thin Films of β-Ga₂O₃ and (Y_{0.06}Ga_{0.94})₂O₃. *Physics and Chemistry of Solid State*, 17(1), 53-59. <https://doi.org/10.15330/pcss.17.1.53-59>
- [9] Sharma, R., Law, M. E., Ren, F., Polyakov, A. Y., & Pearton, S. J. (2021). Diffusion of dopants and impurities in β-Ga₂O₃. *Journal of Vacuum Science & Technology A*, 39(6), 060801. <https://doi.org/10.1116/6.0001307>
- [10] Chen, Z. H., Wang, Y. S., Zhang, N., Zhou, B., Gao, J., Wu, Y. X., ... & Yu, S. W. (2023). Effects of preparation parameters on growth and properties of β-Ga₂O₃ film. *Chinese Physics B*, 32(1), 017301. <https://doi.org/10.1088/1674-1056/ac728c>
- [11] Islam, M. M., Liedke, M. O., Winarski, D., Butterling, M., Wagner, A., Hosemann, P., ... & Selim, F. A. (2020). Chemical manipulation of hydrogen induced high p-type and n-type conductivity in Ga₂O₃. *Scientific reports*, 10(1), 6134.. <https://doi.org/10.1038/s41598-020-62948-2>
- [12] Min, Y., Salvatore, G., Ruth, M., Theo, S., & Reece Michael, J. (2017). Review of flash sintering: materials, mechanisms and modeling. *Adv Appl Ceram*, 116, 24-60. <http://dx.doi.org/10.1080/17436753.2016.1251051>
- [13] Huang, X., Han, S., Huang, W., & Liu, X. (2013). Enhancing solar cell efficiency: the search for luminescent materials as spectral converters. *Chemical society reviews*, 42(1), 173-201. <https://doi.org/10.1039/C2CS35288E>
- [14] Ilyasly, T., Gahramanova, G., Abbasova, R., Veysova, S., & Ismailov, Z. (2021). Investigation of the electrical properties of glasses of Tm-As-S and Tm-As-Se systems. *New Materials, Compounds and Applications*, 5(3), 227-233.

https://jomardpublishing.com/UploadFiles/Files/journals/NMCA/V5N3/Ilyasly_et%20a_l.pdf

- [15] Kageura, T., Kato, K., Yamano, H., Suaebah, E., Kajiya, M., Kawai, S., ... & Kawarada, H. (2017). Effect of a radical exposure nitridation surface on the charge stability of shallow nitrogen-vacancy centers in diamond. *Applied Physics Express*, 10(5), 055503. <https://doi.org/10.7567/APEX.10.055503>
- [16] Pearton, S. J., Yang, J., Cary, P. H., Ren, F., Kim, J., Tadjer, M. J., & Mastro, M. A. (2018). A review of Ga₂O₃ materials, processing, and devices. *Appl. Phys. Rev.*, 5(1), 011301. <https://doi.org/10.1063/1.5006941>
- [17] Ueda, N., Hosono, H., Waseda, R., & Kawazoe, H. (1997). Synthesis and control of conductivity of ultraviolet transmitting β-Ga₂O₃ single crystals. *Appl. Phys. Lett.*, 70(26), 3561–3565. <https://doi.org/10.1063/1.119233>
- [18] Passlack, M., Hunt, N. E. J., Schubert, E. F., Zydzik, G. J., Hong, M., Mannaerts, J. P., ... & Fischer, R. J. (1994). Dielectric properties of electron-beam deposited Ga₂O₃ films. *Appl. Phys. Lett.*, 64(20), 2715–2717. <https://doi.org/10.1063/1.111452>
- [19] Cha, S. Y., Ahn, B. G., Kang, H. C., Lee, S. Y., & Noh, D. Y. (2018). Direct conversion of β-Ga₂O₃ thin films to β-Ga₂O₃ nanowires by annealing in a hydrogen atmosphere. *Ceramics International*, 44(14), 16470-16474. <https://doi.org/10.1016/j.ceramint.2018.06.062>
- [20] Polyakov, A. Y., Lee, I. H., Smirnov, N. B., Yakimov, E. B., Shchemerov, I. V., Chernykh, A. V., ... & Pearton, S. J. (2019). Effects of Hydrogen Plasma Treatment Condition on Electrical Properties of β-Ga₂O₃. *ECS Journal of Solid State Science and Technology*, 8(11) P661-P666. <https://doi.org/10.1149/2.0041911jss>
- [21] Kim, S., Kim, S. J., Kim, K. H., Kim, H. D., & Kim, T. G. (2014). Improved performance of Ga₂O₃/ITO based transparent conductive oxide films using hydrogen annealing for near-ultraviolet light-emitting diodes. *Phys. Status Solidi A*, 211(11), 2569–2573. <https://doi.org/10.1002/pssa.201431278>
- [22] Villafuerte, J., Chaix-Pluchery, O., Kioseoglou, J., Donatini, F., Sarigiannidou, E., Pernot, J., & Consonni, V. (2021). Engineering nitrogen- and hydrogen-related defects in ZnO nanowires using thermal annealing. *Phys. Rev. Materials*, 5, 056001. <https://doi.org/10.1103/PhysRevMaterials.5.056001>
- [23] Gregori G., Merkle R., Maier J., (2017). Ion conduction and redistribution at grain boundaries in oxide systems. *Progress in Materials Science*, 89, 252-305. <https://doi.org/10.1016/j.pmatsci.2017.04.009>
- [24] Golz, C., Galazka, Z., Lähnemann, J., Hortelano, V., Hatami, F., Masselink, W. T., & Bierwagen, O. (2019). The electrical conductivity tensor of Ga₂O₃ analyzed by van der Pauw measurements: Inherent anisotropy, off-diagonal element, and the impact of grain boundaries. *Phys. Rev. Materials*, 3, 124604. <https://doi.org/10.1103/PhysRevMaterials.3.124604>
- [25] Nguyen, V. H., Gottlieb, U., Valla, A., Muñoz, D., Bellet, D., & Muñoz-Rojas, D. (2018). Electron tunneling through grain boundaries in transparent conductive oxides and implications for electrical conductivity: the case of ZnO:Al thin films. *Mater. Horiz.*, 5(4), 715-726. <https://doi.org/10.1039/C8MH00402A>
-

ВПЛИВ УМОВ ОДЕРЖАННЯ НА ЕЛЕКТРОПРОВІДНІ ВЛАСТИВОСТІ ТОНКИХ ПЛІВОК $(Y_{0.06}Ga_{0.94})_2O_3$

*Ігор Кухарський**^{ORCID}, *Ірина Кофлюк*^{ORCID}, *Іванна Медвідь*^{ORCID},
Ігор Кузь^{ORCID} & *Олег Бордун*^{ORCID}

*Кафедра фізичної та біомедичної електроніки
Львівський національний університет імені Івана Франка
вул. Драгоманова, 50, 79005 Львів, Україна*

АНОТАЦІЯ

Вступ. Оксид галію $\beta\text{-Ga}_2\text{O}_3$ є перспективним широкозонним напівпровідниковим матеріалом, який широко застосовується в оптоелектроніці та сенсорних системах. При цьому електропровідність тонких плівок значною мірою визначається їх структурною досконалістю, дефектною підсистемою та умовами післяростової обробки. У полікристалічних плівках міжзернові границі та дефектні комплекси суттєво впливають на механізми перенесення заряду. Метою роботи є дослідження впливу атмосфери відпалу на структурні, морфологічні та електропровідні властивості тонких плівок $(Y_{0.06}Ga_{0.94})_2O_3$.

Матеріали та методи. Тонкі плівки $(Y_{0.06}Ga_{0.94})_2O_3$ товщиною 0,3–1,0 мкм були отримані методом ВЧ іонно-плазмового розпилення на підкладках із плавненого кварцу. Після осадження проводився відпал у кисні та аргоні при температурах 1000–1100 °С, а також у водні при 600–650 °С. Структурні властивості досліджували методом рентгенівської дифракції, морфологію поверхні — за допомогою атомно-силової мікроскопії. Електропровідність вимірювали в температурному діапазоні 300–450 К, а енергію активації визначали з температурних залежностей провідності.

Результати. Рентгеноструктурний аналіз підтвердив формування плівок у моноклінній фазі $\beta\text{-Ga}_2\text{O}_3$ та підвищення ступеня кристалічності після відпалу в кисні. Встановлено, що свіжнанесені плівки мають високий питомий опір ($\rho > 10^{11}$ Ом·см), який зменшується при підвищенні температури та після відпалу. Для плівок, відпалених у кисні, енергія активації становить близько 0,87 еВ, тоді як для відпалу в аргоні вона досягає ~1,38 еВ (у діапазоні 300–400 К), що свідчить про формування глибших донорних рівнів, пов'язаних із кисневими вакансіями. Відпал у водні призводить до значного зниження питомого опору (~ 10^8 Ом·см) та енергії активації (~0,25 еВ), що зумовлено появою мілких донорних рівнів.

Висновок. Електропровідність тонких плівок $(Y_{0.06}Ga_{0.94})_2O_3$ визначається донорними рівнями дефектної природи, які формуються залежно від умов відпалу. Кисневе та аргонне середовище сприяють утворенню глибоких донорних центрів, тоді як відпал у водні призводить до формування мілких донорів і підвищення електропровідності.

Ключові слова: тонкі плівки, оксид галію, кристалічна структура, домішки, електропровідність.

UDC: 537.312, 621.383

PHOTODETECTORS BASED ON FIELD EFFECT IN POROUS SILICON – REDUCED GRAPHENE OXIDE STRUCTURES

Igor Olenych*  & Andrii Kozak 

Department of Radioelectronic and Computer Systems
Ivan Franko National University of Lviv
50 Dragomanov Street, 79005 Lviv, Ukraine

Olenych, I. and Kozak, A. (2026). Photodetectors Based on Field Effect in Porous Silicon – Reduced Graphene Oxide Structures. *Electronics and Information Technologies*, 33, 191–200. <https://doi.org/10.30970/eli.33.14>

ABSTRACT

Background. Graphene field-effect transistors (FETs) have high potential for application in sensor electronics as detectors of electromagnetic radiation in a wide spectral range due to the high sensitivity of graphene ambipolar conductivity to local changes in the electric field. The use of a reduced graphene oxide (RGO) film provides cost reduction of photodetectors based on graphene FETs. On the other hand, an additional porous silicon light-absorbing layer can increase their sensitivity due to an increase in surface area.

Materials and methods. Graphene field-effect photodetectors were created by drying a film-forming RGO suspension deposited on the surface of the porous silicon on a silicon substrate, which served as the gate of the FET. Electrical source and drain contacts were formed on the surface of the obtained RGO film. To improve the insulating properties of porous silicon, it was electrochemically oxidized, and an additional layer of Al_2O_3 was deposited. The electrical and photoelectric properties of the created field-effect photodetectors were investigated in DC and AC modes using a white LED and standard optical equipment.

Results. An increase in the conductivity and capacitance of the RGO channel of the FETs was detected under the influence of white light irradiation. Based on the analysis of the drain current dependencies on the gate voltage, it has been established that the efficiency and photosensitivity of the FETs based on the porous silicon and RGO film are increased by the deposition of an additional Al_2O_3 layer on the surface of electrochemically oxidized porous silicon. The maximum sensitivity of the created photodetectors is in the spectral range of 800–900 nm. The response time to white light pulses is about 0.5 ms. Passivation of the porous silicon surface with the oxide film and the Al_2O_3 layer causes an increase in the photosignal relaxation time.

Conclusions. The features of using FETs based on porous silicon structures and RGO film as visible radiation detectors have been investigated. The electrical, spectral, and time characteristics of the created field-effect photodetectors were determined.

Keywords: Graphene field-effect transistor, reduced graphene oxide, porous silicon, photosensitivity.

INTRODUCTION

The unique combination of ultra-high charge carrier mobility in graphene, sensitivity of its electrical conductivity to local changes in the electric field, high transparency, and flexibility provides exceptional prospects for graphene and materials based on sp^2 -bonded



© 2026 Igor Olenych & Andrii Kozak. Published by the Ivan Franko National University of Lviv on behalf of Електроніка та інформаційні технології / Electronics and Information Technologies. This is an Open Access article distributed under the terms of the [Creative Commons Attribution 4.0 License](https://creativecommons.org/licenses/by/4.0/) which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

carbon atoms in sensor electronics for creating photodetectors, touch displays, and other next-generation optoelectronic devices [1-3]. Due to the absence of a band gap, graphene absorbs photons of almost any energy, which opens up the possibility of creating photodetectors with a wide spectral range, namely from UV to terahertz radiation. These properties make graphene especially attractive for high-speed telecommunications systems [4, 5], contactless scanning systems [6], and optical sensors in spectroscopic methods of medical diagnostics [7, 8], where both response time and energy efficiency are critical.

Despite the attractiveness of graphene for sensor applications, challenges can be noted that hinder the development of graphene-based devices and their mass industrial application. One obstacle remains scaling up the production of high-quality graphene and ensuring its uniformity over large areas. A solution to this problem may be the use of reduced graphene oxide (RGO), since its manufacturing technology is quite simple and low-cost [9,10]. Field-effect transistors (FETs) based on the RGO film demonstrate high sensitivity to electromagnetic and ionizing radiation [11].

On the other hand, pure graphene demonstrates relatively low quantum efficiency due to its limited absorption capacity (~2.3 % per layer). Therefore, the basic approach to graphene applications in photoelectronics is focused on engineering hybrid structures that combine the charge-sensitive carbon monolayer and a photosensitive material for charge carrier generation. In particular, combining graphene with strong light absorbers, such as perovskites or nanostructured semiconductors and metals, provides a significant increase in photosensitivity due to charge transfer effects or enhancement of local electric fields [12-14]. According to this concept, porous silicon (PS) has high potential for use as a substrate for graphene deposition due to its attractive antireflective properties [15]. The large ratio of surface area to volume of the porous layer provides not only effective absorption of light quanta, but also promotes the deposition of various nature nanoparticles, including RGO [16]. In addition, the nanostructured PS can be used as a supporting layer for the RGO channel of the FETs, owing to its low electrical conductivity [17]. Further improvement of the porous layer dielectric properties increases the efficiency of FETs based on the PS–RGO structures. Therefore, the work aimed to study the relationship between the structural features of the PS-based supporting layer and the electrical and photosensitive properties of the FETs with the RGO film channel.

MATERIALS AND METHODS

Slightly doped silicon wafers with a thickness of 400 μm were used as the substrate and gate of the RGO-based FETs. A thin gold film used as an electrical contact for the gate was thermally deposited on the back surface of the wafers and annealed at a temperature of 600°C for 30 min. On the opposite side of the wafers, the PS layer was formed by the electrochemical method. An ethanolic solution of hydrofluoric acid with a component ratio of $\text{HF}:\text{C}_2\text{H}_5\text{OH} = 1:1$ was used as the electrolyte. The current density and duration of anodic etching were 30 mA/cm^2 and 5 min, respectively. The wafer working surface was illuminated with a Feron HB1 J118 500 W incandescent lamp throughout the entire electrochemical etching process to generate positive charge carriers necessary for the silicon etching chemical reactions. After washing with distilled water, some of the samples were subjected to electrochemical oxidation of the porous layer in a H_2O_2 solution at a current density of 15 mA/cm^2 for 10 min to stabilize the surface and form a thin dielectric film. To improve the dielectric properties of the PS used as an insulating layer between the gate and the RGO channel of the FETs, a 150 nm thick Al_2O_3 film was additionally deposited on the surface of the anodically oxidized porous layer of several samples by RF magnetron sputtering under the conditions described in [11].

The conductive channel of the FETs was formed by air-drying the RGO film-forming suspension deposited on the surface of the insulating layer, obtained by reducing graphene oxide (in the form of an aqueous suspension from Sigma-Aldrich) with hydrazine monohydrate as described in [16]. Silver contacts of source and drain were thermally deposited onto the surface of the formed RGO film at a distance of 1 mm from each other, as shown in Fig. 1.

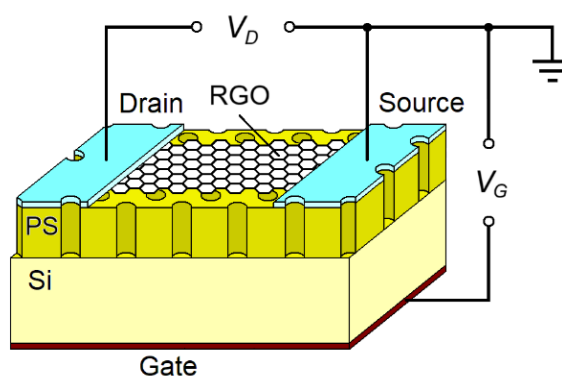


Fig. 1. Schematic representation of the FET based on the PS and RGO film.

The electrical and photoelectric characteristics of the created FETs were studied in DC and AC modes. In particular, the dependencies of the drain current I_D on the bias voltage V_D and the gate voltage V_G were measured using a Siglent SDM 3045 multimeter. The capacitive and resistive characteristics of the RGO channel were investigated using a Hantek 1833C RLC meter in the frequency range of 10^2 – 10^5 Hz. A white LED FYLP–1W–UWB–A with a power of 1 W and a luminous flux of 76 lumens was used to study photoelectric phenomena. The LED emission spectrum has a broad band in the range of 470–600 nm and an intense band with a maximum at about 450 nm. The photoresponse spectra were measured on standard optical equipment and corrected to account for the spectral characteristics of the equipment.

RESULTS AND DISCUSSION

The dependencies of the drain current I_D on the potential difference between the source and drain contacts V_D of the created FETs in the absence of gate voltage ($V_G = 0$) are nonlinear, as shown in Fig. 2. The lowest conductivity was found in the RGO film deposited on the surface of the PS (Si–PS–RGO structure). The RGO film formed on the surface of an anodically oxidized porous layer (Si–PS_{ox}–RGO structure) demonstrates slightly higher conductivity. The highest conductivity of the RGO film is observed in the structure with an additional Al₂O₃ layer (Si–PS_{ox}–Al₂O₃–RGO). Irradiation with white light from the LED FYLP–1W–UWB–A causes an increase in the drain current I_D of the RGO-based FETs. The observed increase in the conductivity of the RGO film is likely caused by the electric field of charge carriers photogenerated in the PS and silicon substrate.

Although the different conductivity values are most likely caused by the different number of carbon nanosheet layers in the RGO film, the nonlinear nature of the I_D – V_D dependencies may have several reasons. Considering that graphene is characterized by ohmic contact resistance with metals [18], it can be assumed that the found nonlinearity is due to both the heterogeneity of the formed films and the influence of the PS layer on the

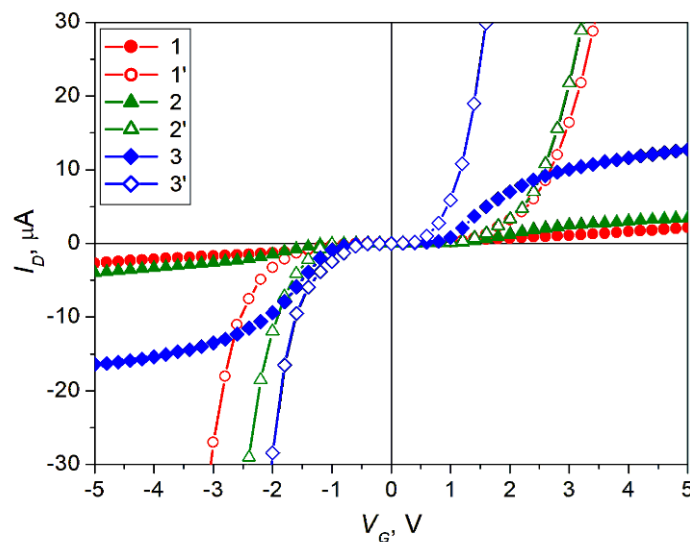
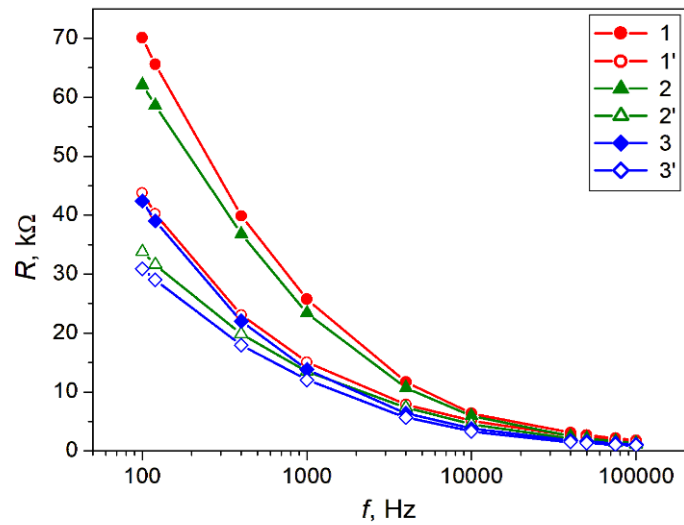


Fig. 2. Dependencies of drain current I_D on bias voltage V_D of the FETs based on the Si-PS-RGO (1,1'), Si-PS_{ox}-RGO (2,2'), and Si-PS_{ox}-Al₂O₃-RGO (3,3') structures measured in the dark (1,2,3) and under irradiation with white light (1',2',3') at the gate voltage $V_G = 0$.

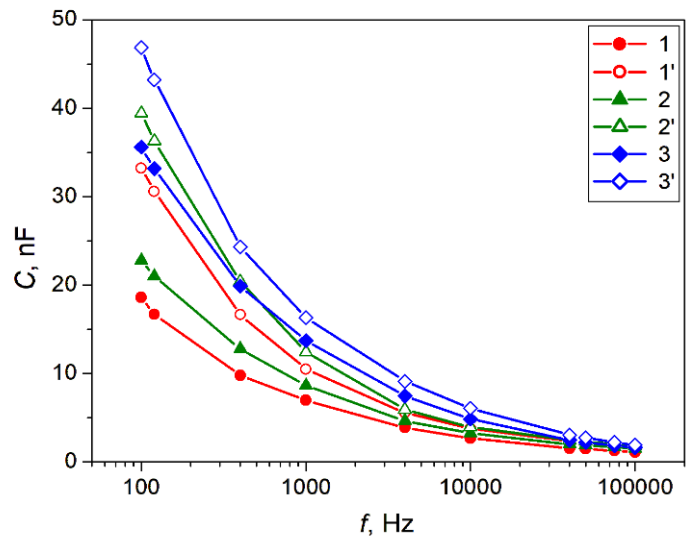
RGO conductivity. First of all, the electrical properties of the RGO film depend not only on the conductivity of the 2D carbon nanoparticles but also on the electrical barriers between them, which are likely formed due to the presence of a surfactant in the film-forming suspension [16]. In addition, since porous silicon is not a perfect insulator [15], charge carriers injected from the PS should not be ignored. Finally, due to the structural imperfection of the PS dielectric coating, the electrical conductivity of the RGO film can also be affected by charge carriers localized on electrically active defects and the interface with the RGO.

The capacitive and resistive properties of the created FETs were investigated in the AC mode to obtain additional information about the charge transfer processes in the RGO film. **Fig. 3** shows the frequency dependencies of the internal resistance and capacitance of the RGO film, measured between the source and drain contacts. A decrease in electrical characteristics was detected with increasing frequency from 100 Hz to 100 kHz. A decrease in resistance and an increase in capacitance of the RGO film were also found under irradiation of the working surface of the field-effect photodetectors. The observed dispersion of the electrical characteristics of the RGO channel of the FETs may further support the assumption of inhomogeneities in the RGO film formed from carbon nanosheets.

Dependencies of the drain current I_D on the gate voltage V_G were measured to evaluate the efficiency of the created FETs based on the RGO film. The obtained I_D - V_G curves at $V_D = \pm 1.5$ V are shown in **Fig. 4**. The drain current increases linearly with the change in gate voltage from approximately 0 to -4 V for the bias voltage of $V_D = 1.5$ V. Similarly, the conductivity of the RGO channel of the FETs increases almost linearly with increasing V_G from 1.5 to 4 V for $V_D = -1.5$ V. The minimum conductivity of the RGO-based FETs at about 1 V gate voltage is caused by the features of the graphene band structure in the form of Dirac cones and is associated with the charge neutrality point, which divides the conductivity profile into hole and electron components [19].



a)



b)

Fig 3. Frequency dependencies of the internal resistance (a) and capacitance (b) of the RGO channel of the FETs based on the Si-PS-RGO (1,1'), Si-PS_{ox}-RGO (2,2'), and Si-PS_{ox}-Al₂O₃-RGO (3,3') structures measured in the dark (1,2,3) and under irradiation with white light (1',2',3') at the gate voltage of $V_G = 0$.

Analysis of the measured I_D-V_G dependencies shows that the FET based on the Si-PS_{ox}-Al₂O₃-RGO structure, which demonstrates the largest range of change in drain current when the gate voltage changes, has the highest efficiency. The revealed features of the electrical characteristics of the FETs based on the PS and RGO film may be due to the different quality of the insulating layers, since the performance of graphene field-effect devices depends significantly on defects in the supporting dielectric layer [20].

As a result of irradiation of the FETs with white light, an increase in the drain current I_D was observed at a bias voltage of $V_D = \pm 1.5$ V. The highest photosensitivity to white light is demonstrated by the structure with the additional Al₂O₃ layer, probably due to the higher FET efficiency. The small thickness of the transparent to visible light Al₂O₃ layer does not

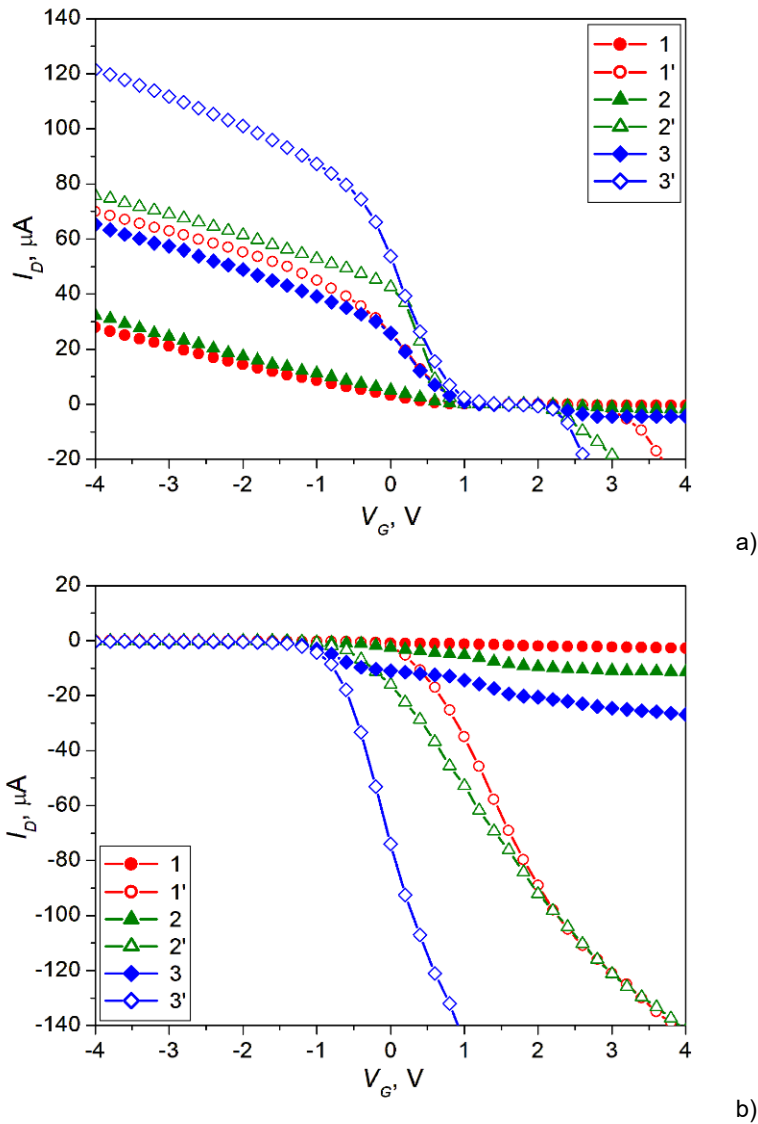


Fig 4. Dependencies of drain current I_D on gate voltage V_G of the FETs based on the Si-PS-RGO (1,1'), Si-PS_{ox}-RGO (2,2'), and Si-PS_{ox}-Al₂O₃-RGO (3,3') structures measured in the dark (1,2,3) and under irradiation with white light (1',2',3') at the bias voltages $V_D = 1.5$ V (a) and $V_D = -1.5$ V (b).

significantly affect the sensitivity of the proposed field-effect photodetector. In addition, improved passivation of the PS surface with the Al₂O₃ layer can reduce the rate of photogenerated charge carrier recombination. An additional argument in favor of this hypothesis is the increase in the relaxation time of the photoresponse to white light pulses, as shown in the inset of Fig. 5. In general, the analysis of the photoresponse kinetics of the created photodetectors revealed almost the same rate of increase in the photosignal and different times of its decay when excited by light pulses with a duration of 1 ms.

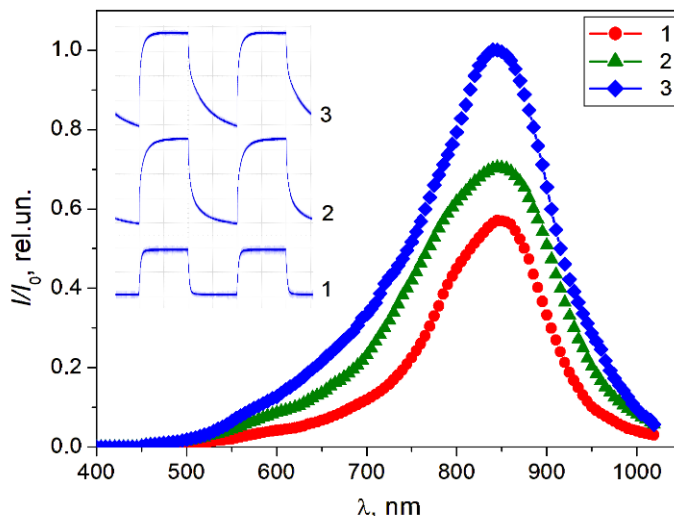


Fig 5. Photoresponse spectra of the FETs based on the Si-PS-RGO (1), Si-PS_{ox}-RGO (2), and Si-PS_{ox}-Al₂O₃-RGO (3) structures. Inset: response kinetics of the photodetectors to white light pulses.

The spectral characteristic of the field-effect photodetectors depends on the light-absorbing properties of the PS and silicon substrate. The photosensitivity spectra of the RGO-based FETs are characterized by a broad band with a maximum in the 800–900 nm range, similar to a silicon photodiode (see Fig. 5).

CONCLUSION

As a result of studying the electrical and photoelectric characteristics of the FETs based on the PS and RGO film, the features of their use as photodetectors have been established. It was revealed that irradiation of FETs with white light causes an increase in the conductivity of the RGO channel in both DC and AC modes. In addition, the efficiency and photosensitivity of the created FETs were increased due to the deposition of an additional Al₂O₃ layer on the surface of the anodically oxidized PS. Moreover, the improved passivation of the PS surface with the Al₂O₃ layer increases the relaxation time of the photoresponse to white light pulses. The response time of the field-effect photodetectors is about 0.5 ms. It has been established that the photosensitivity of FETs based on the RGO film is maximum in the range of 800–900 nm.

ACKNOWLEDGMENTS AND FUNDING SOURCES

This work was supported by the Ministry of Education and Science of Ukraine [0124U000982].

COMPLIANCE WITH ETHICAL STANDARDS

The authors declare that the research was conducted in the absence of any conflict of interest.

AUTHOR CONTRIBUTIONS

Conceptualization, [IO]; methodology, [IO, AK]; investigation, [IO, AK]; writing – original draft preparation, [IO, AK]; writing – review and editing, [IO]; visualization, [IO].

All authors have read and agreed to the published version of the manuscript.

REFERENCES

- [1] Zhang, X., Liu, X., Huang, Y., Sun, B., Liu, Z., Liao, G., & Shi, T. (2023). Review on flexible perovskite photodetector: processing and applications. *Front. Mech. Eng.*, 18, 33. <https://doi.org/10.1007/s11465-023-0749-z>
- [2] Esteghamat, A. & Akhavan, O. (2023). Graphene as the ultra-transparent conductive layer in developing the nanotechnology-based flexible smart touchscreens. *Microelectronic Engineering*, 267-268, 111899. <https://doi.org/10.1016/j.mee.2022.111899>
- [3] Ishida, S., Anno, Y., Takeuchi, M., Matsuoka, M., Takei, K., Arie, T., & Akita, S. (2015). Highly photosensitive graphene field-effect transistor with optical memory function. *Sci. Rep.*, 5, 15491. <https://doi.org/10.1038/srep15491>
- [4] Romagnoli, M., Sorianello, V., Midrio, M., Koppens, F.H.L., Huyghebaert, C., Neumaier, D., Galli, P., Templ, W., D'Errico, A., & Ferrari, A.C. (2018). Graphene-based integrated photonics for next-generation datacom and telecom. *Nature Reviews Materials*, 3, 392-414. <https://doi.org/10.1038/s41578-018-0040-9>
- [5] Ding, Y., Cheng, Z., Zhu, X., Yvind, K., Dong, J., Galili, M., Hu, H., Mortensen, N., Xiao, S., & Oxenløwe, L. (2020). Ultra-compact integrated graphene plasmonic photodetector with bandwidth above 110 GHz. *Nanophotonics*, 9(2), 317-325. <https://doi.org/10.1515/nanoph-2019-0167>
- [6] Liu, J., Li, X., Jiang, R., Yang, K., Zhao, J., Khan, S. A., He, J., Liu, P., Zhu, J., & Zeng, B. (2021). Recent Progress in the Development of Graphene Detector for Terahertz Detection. *Sensors*, 21(15), 4987. <https://doi.org/10.3390/s21154987>
- [7] Suvamphaet, P., & Pechprasarn, S. (2017). Graphene-Based Materials for Biosensors: A Review. *Sensors*, 17(10), 2161. <https://doi.org/10.3390/s17102161>
- [8] Baruah, A., Newar, R., Das, S., Kalita, N., Nath, M., Ghosh, P., Chinnam, S., Sarma, H., & Narayan, M. (2024). Biomedical applications of graphene-based nanomaterials: recent progress, challenges, and prospects in highly sensitive biosensors. *Discover Nano*, 19, 103. <https://doi.org/10.1186/s11671-024-04032-6>
- [9] Pei, S. & Cheng, H.M. (2012). The reduction of graphene oxide. *Carbon*, 50, 3210–3228. <https://doi.org/10.1016/j.carbon.2011.11.010>
- [10] Abakumov, A.A., Bychko, I.B., Voitsihovska, O.O., Rudenko, R.M., & Strizhak, P.E. (2024). Tuning the surface area of reduced graphene oxide by modulating graphene oxide concentration during hydrazine reduction. *Materials Letters*, 354, 135417. <https://doi.org/10.1016/j.matlet.2023.135417>
- [11] Olenych, I.B., Monastyrskii, L.S., Sokolovskii, B.S., Turko, B.I., & Dzendzelyuk, O.S. (2025). Field-effect transistors based on reduced graphene oxide film for photo and radiation detectors. *Sensor Electronics and Microsystem Technologies*, 22(2), 19–26, (in Ukrainian). <https://doi.org/10.18524/1815-7459.2025.2.333193>
- [12] Abbas, K., Ji, P., Ullah, N., Shafique, S., Zhang, Z., Ameer, M.F., Qin, S., & Yang, S. (2024). Graphene photodetectors integrated with silicon and perovskite quantum dots. *Microsyst. Nanoeng.*, 10, 81, <https://doi.org/10.1038/s41378-024-00722-4>
- [13] Thai, K.Y., Park, I.J., Kim, B.J., Hoang, A.T., Na, Y., Park, C.U., Chae, Y., & Ahn, J.-H. (2021). MoS₂/Graphene Photodetector Array with Strain-Modulated Photoresponse up to the Near-Infrared Regime. *ACS Nano*, 15(8), 12836–12846. <https://doi.org/10.1021/acsnano.1c04678>

- [14] Ye, M., Zha, J., Tan, C., & Crozier, K.B. (2021). Graphene-based mid-infrared photodetectors using metamaterials and related concepts. *Appl. Phys. Rev.* 8, 031303. <https://doi.org/10.1063/5.0049633>
- [15] Bisi, O., Ossicini, S., & Pavesi, L. (2000). Porous silicon: a quantum sponge structure for silicon based optoelectronics. *Surf. Sci. Rep.*, 38, 1–126. [https://doi.org/10.1016/S0167-5729\(99\)00012-6](https://doi.org/10.1016/S0167-5729(99)00012-6)
- [16] Olenych, I. B. & Horbenko, Yu. Yu. (2024). Electrical and photoelectric properties of hybrid structures based on reduced graphene oxide and Pd-doped porous silicon. *Mol. Cryst. Liq. Cryst.*, 768, 135-144. <https://doi.org/10.1080/15421406.2023.2235191>
- [17] Olenych, I.B., Horbenko, Y.Y., & Sokolovskii, B.S. (2024). Effect of supporting layer on electrical characteristics of field-effect transistor based on reduced graphene oxide film. *Mol. Cryst. Liq. Cryst.*, 768(11), 426–435. <https://doi.org/10.1080/15421406.2024.2353960>
- [18] Xia, F., Perebeinos, V., Lin, Y.-M., Wu, Y., & Avouris P. (2011). The origins and limits of metal graphene junction resistance. *Nature Nanotechnology*, 6(3), 179–184. <https://doi.org/10.1038/nnano.2011.6>
- [19] Zhan, D., Yan, J., Lai, L., Ni, Z., Liu, L., & Shen, Z. (2012). Engineering the Electronic Structure of Graphene. *Adv. Mater.*, 24, 4055–4069. <https://doi.org/10.1002/adma.201200011>
- [20] Imamura, G. & Saiki, K. (2015). Modification of Graphene/SiO₂ Interface by UV-Irradiation: Effect on Electrical Characteristics. *ACS Appl. Mater. Interfaces*, 7, 2439–2443. <https://doi.org/10.1021/am5071464>
-

ФОТОДЕТЕКТОРИ НА ОСНОВІ ЕФЕКТУ ПОЛЯ У СТРУКТУРАХ ПОРУВАТИЙ КРЕМНІЙ – ВІДНОВЛЕНИЙ ОКСИД ГРАФЕНУ

Ігор Оленич* , Андрій Козак 

Кафедра радіоелектронних і комп'ютерних систем
Львівський національний університет імені Івана Франка,
вул. Драгоманова 50, 79005 м. Львів, Україна

АНОТАЦІЯ

Вступ. Завдяки високій чутливості біполярної провідності графену до локальної зміни електричного поля, графенові польові транзистори мають високий потенціал для застосування у сенсорній електроніці, зокрема як детектори електромагнітного випромінювання у широкому спектральному діапазоні. Використання плівки відновленого оксиду графену (ВОГ) дає змогу знизити вартість фотодетекторів на основі графенових польових транзисторів, тоді як додатковий світлопоглинаючий шар поруватого кремнію може забезпечити підвищення їхньої чутливості завдяки збільшенню площі поверхні.

Матеріали та методи. Фоточутливі графенові польові транзистори отримано висушуванням плівкоутворювальної суспензії ВОГ, нанесеної на поверхню поруватого кремнію на кремнієвій підкладці, що слугувала затвором польового транзистора. На поверхні утвореної плівки ВОГ сформовано електричні контакти витоку та стоку. Для покращення ізоляційних властивостей поруватого кремнію його було електрохімічно окиснено та осаджено додатковий шар Al₂O₃. Електричні та фотоелектричні властивості створених польових транзисторів досліджено у режимах постійного та

змінного струму з використанням світловипромінювального діода білого світла і стандартного оптичного обладнання.





Результати. Виявлено збільшення провідності та ємності ВОГ-каналу польових транзисторів за впливу опромінення білим світлом. На основі аналізу залежностей струму стоку від напруги затвора встановлено, що ефективність і фоточутливість польових транзисторів на основі поруватого кремнію та плівки ВОГ збільшується завдяки осадженню додаткового шару Al_2O_3 на поверхню електрохімічно окисненого поруватого кремнію. Максимум чутливості створених фотодетекторів знаходиться у спектральному діапазоні 800–900 нм. Час відгуку на імпульси білого світла становить близько 0,5 мс. Пасивація поверхні поруватого кремнію оксидною плівкою та шаром Al_2O_3 зумовлює збільшення часу релаксації фотосигналу.

Висновки. Досліджено особливості використання польових транзисторів на основі поруватого кремнію та плівки ВОГ як детекторів видимого випромінювання. Визначено електричні, спектральні та часові характеристики створених фотодетекторів.

Ключові слова: Графеновий польовий транзистор, відновлений оксид графену, поруватий кремній, фоточутливість.

UDC 004.89

INTEGRATION OF DECENTRALIZED PERFORMANCE VERIFICATION IN HYBRID ARCHITECTURES EDGE-FOG-CLOUD TO INCREASE IoT SYSTEMS RELIABILITY

Roman Diachok*  , Halyna Klym  
Lviv Polytechnic National University,
12 Stepan Bandera St., Lviv 79013, Ukraine

Roman Diachok & Halyna Klym (2026). Integration of Decentralized Performance Verification in Hybrid Architectures Edge-Fog-Cloud to Increase IoT Systems' Reliability. *Electronics and Information Technologies*, 33, 201–208. <https://doi.org/10.30970/eli.33.15>

ABSTRACT

Background. The rapid growth of Internet of Things (IoT) systems has increased the demand for scalable and low-latency data processing architectures. Traditional cloud-centric approaches often suffer from high communication delays and bandwidth limitations. Edge–Fog–Cloud computing introduces a multi-tier model that distributes computational tasks closer to data sources. However, evaluating computational methods in such heterogeneous environments requires systematic performance analysis and architectural optimization. In this context, integrating mathematically stable and computationally efficient methods, such as harmonic potential field–based approaches, is essential to ensure reliable real-time operation, scalability, and system resilience across distributed layers.

Methods. This study evaluates the Laplace artificial potential field method implemented within a multi-tier Edge–Fog–Cloud architecture. The experimental framework includes distributed simulation, real-time processing scenarios, and comparative benchmarking. Performance metrics such as latency, computational load, and system stability were analyzed. The proposed approach was tested under variable workload conditions to assess scalability and efficiency across architectural layers.

Results and Discussion. Experimental results demonstrate reduced end-to-end latency and improved task distribution across edge and fog layers. Compared to centralized processing, the proposed architecture maintains stability under increased workload. The Laplace-based computational model ensures efficient obstacle handling and balanced resource utilization. These findings confirm that multi-tier orchestration enhances system responsiveness while preserving acceptable computational overhead in dynamic IoT environments.

Conclusion. Integrating the Laplace artificial potential field method within an Edge–Fog–Cloud architecture significantly improves distributed system performance. The proposed framework increases scalability, reliability, and computational efficiency in real-time IoT applications, providing a solid foundation for further optimization of resource management and intelligent task allocation in heterogeneous distributed environments.

Keywords: Edge–Fog–Cloud; IoT; Laplace artificial potential field; distributed computing; real-time processing; latency; runtime verification.

INTRODUCTION

The goal of the research is to theoretically substantiate and practically develop a distributed framework that combines the low latency of edge computing with the power of cloud analytical platforms. The primary objective is to create a dynamic orchestration



© 2026 Roman Diachok & Halyna Klym. Published by the Ivan Franko National University of Lviv on behalf of Електроніка та інформаційні технології / Electronics and information technologies. This is an Open Access article distributed under the terms of the [Creative Commons Attribution 4.0 License](https://creativecommons.org/licenses/by/4.0/) which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

mechanism that automatically distributes processing and verification tasks based on their computational complexity and response-time requirements.

Traditional cloud-centric IoT architectures suffer from performance bottlenecks caused by limited bandwidth and increased latency when large volumes of raw data are transmitted to centralized cloud infrastructures. This limitation becomes particularly critical in real-time and mission-critical systems. To mitigate these constraints, research has increasingly focused on distributed computing paradigms that move computation closer to data sources. Edge computing significantly reduces network traffic and improves system responsiveness by enabling local data processing, while fog computing introduces an intermediate layer that supports regional aggregation and context-aware analytics [1,2].

The effectiveness of heterogeneous edge–fog–cloud collaboration has been further demonstrated through architectural implementations and benchmarking studies, which confirm improvements in scalability, energy efficiency, and response time under dynamic workloads [3,4]. However, while orchestration and resource management techniques such as NSGA-III and MEC concentrate on traffic optimization and computational efficiency, they often lack formal guarantees of reliability and strict compliance with real-time requirements.

In dynamic IoT environments, where network topology and system states continuously evolve, runtime verification (RV) emerges as a critical mechanism for ensuring correctness. Edge-based runtime verification frameworks enable monitors to operate directly on distributed nodes, reducing reliance on centralized analysis and improving responsiveness [5]. The application of Metric First-Order Temporal Logic (MFOTL) for real-time policy enforcement provides a rigorous formal foundation for specifying and verifying temporal constraints over event streams [2]. By employing monitoring tools such as MonPoly, systems can transition from exhaustive design-time verification to lightweight real-time event monitoring, which is more suitable for evolving IoT deployments.

Formally verified monitoring frameworks, such as VeriMon, further enhance reliability by ensuring correctness of the monitoring algorithms themselves [6]. More recent advancements in anticipatory runtime verification extend this concept by enabling predictive reasoning over temporal properties, thereby allowing systems to detect potential violations before they fully manifest [7]. Additional runtime verification tools targeting context-aware distributed systems further reinforce the feasibility of applying formal monitoring in heterogeneous IoT ecosystems [8].

Collectively, these works indicate a clear evolution from purely performance-oriented optimization strategies toward integrated architectures that combine distributed edge–fog–cloud processing with formal runtime verification mechanisms. Such integration is essential for achieving scalability, real-time responsiveness, and dependable operation in modern IoT systems.

The rapid expansion of Internet of Things (IoT) systems has significantly increased the demand for scalable, low-latency, and computationally efficient distributed architectures. Traditional centralized cloud-based solutions often suffer from bandwidth limitations and high end-to-end latency, particularly in real-time applications. Multi-tier Edge–Fog–Cloud architectures enable task distribution closer to data sources, thereby reducing communication overhead and improving responsiveness. However, the integration of computational intelligence methods within such heterogeneous environments requires systematic evaluation of performance and scalability [9-10].

This study aims to evaluate the computational efficiency and scalability of the Laplace artificial potential field method implemented within a multi-tier Edge–Fog–Cloud architecture for real-time IoT applications.

METHODS

Recent studies demonstrate the effectiveness of Edge–Fog–Cloud architectures in reducing latency and improving workload distribution in IoT systems. At the same time,

research on runtime verification and distributed monitoring highlights the importance of formal reliability enforcement in dynamic environments. Despite significant progress, many existing works emphasize performance optimization without a comprehensive evaluation of computational models under heterogeneous deployment conditions. Therefore, there remains a need for integrated approaches that combine architectural scalability with computational efficiency assessment.

To ensure the validity of the comparative analysis between Raspberry Pi Zero, Raspberry Pi 3, and the x86 VM, a standardized software environment was maintained across all platforms. Each node operated on a Debian-based Linux distribution (Raspberry Pi OS for ARM-based nodes and Debian 12 for the cloud VM) using an identical version of the Python 3.10 runtime environment. This approach minimizes performance variances that could arise from operating system overhead or differing software stacks, ensuring that the measured throughput primarily reflects the computational capabilities of the underlying hardware.

The methodology is based on a three-level hierarchical data processing model that encompasses six stages of the information life cycle. Data from IoT sensors undergo filtering to eliminate noise and anomalies. Low-pass and high-pass filters, as well as Kalman filters, are used to reduce random fluctuations. The mathematical model for preprocessing is defined as (1):

$$PD = f(RD, FP, NF), \quad (1)$$

where PD is *Preprocessed Data*, RD – *Raw Data*, FP – *Filtering Parameters*, and NF – *Normalization Factors*.

The tiered computation distribution is organized as follows. The Edge Layer (L0) performs simple operations (filtering, threshold alerts) and local RV. The processing time is defined as (2):

$$EPT = \frac{DV \times PC}{ECR}, \quad (2)$$

where EPT is *Edge Processing Time*, DV – *Data Volume*, PC – *Processing Complexity*, and ECR – *Edge Computing Resources*.

The Fog Layer (L1) aggregates data from multiple Edge nodes. Processing efficiency at this level is calculated using the formula (3):

$$EDP = \frac{FNP}{LaFN + TD} \quad (3)$$

where EDP is *Effective Data Processed*, FNP – *Fog Node Performance*, $LaFN$ – *Latency at Fog Node*, and TD – *Transmission Delay*.

The Cloud Layer (L2) is designed for executing complex machine learning tasks and long-term storage (4). The processing time is defined as:

$$CPT = \frac{TC \times DV}{CR}, \quad (4)$$

where CPT is *Cloud Processing Time*, DV – *Data Volume*, TC – *Task Complexity*, and CR – *Cloud Resources*.

For dynamic orchestration and life cycle management, the system automatically classifies tasks by complexity (low, medium, high) and resource intensity. The orchestration decision is made based on maximizing performance relative to latency (5):

$$OD = \operatorname{argmax}_{\text{layer}} \left(\frac{PC}{L} \right), \quad (5)$$

where OD is *Orhestration Decision*, PC – *Processing Capability*, L – *Latency*.

The data retention period is dynamically optimized depending on its priority (6):

$$RP = \frac{DIF \times DT}{AS}, \quad (6)$$

where RP is *Retention Period*, DIF – *Data Importance Factor*, DT – *Data Type*, AS is *Available Storage*.

The Laplace artificial potential field method was implemented within the proposed architecture and evaluated under controlled simulation conditions. The experimental setup included distributed workload scenarios with varying event rates to assess latency dynamics and computational load balancing. Performance metrics such as end-to-end latency, processing efficiency, and system stability were measured. Comparative analysis with centralized deployment was conducted to validate the effectiveness of distributed task allocation.

RESULTS AND DISCUSSION

The proposed system architecture follows a hierarchical multi-tier model consisting of edge, fog, and cloud layers. The edge layer performs initial data acquisition and preprocessing, minimizing raw data transmission. The fog layer is responsible for intermediate aggregation and distributed computational coordination. The cloud layer provides centralized analytics and long-term data storage. Such hierarchical structuring enables dynamic task redistribution depending on workload conditions and enhances modular scalability of the system.

The obtained results demonstrate that the multi-tier architecture significantly reduces end-to-end latency compared to centralized processing models. The Laplace-based computational approach ensures stable system behavior even under increased workload conditions. Efficient distribution of computational tasks between edge and fog layers contributes to improved responsiveness and balanced resource utilization. These findings confirm that distributed orchestration enhances overall system performance without introducing excessive computational overhead.

The primary dependencies obtained during experiments on a test bench (Raspberry Pi 3, x86 Cloud VM) are presented below in **Table 1**.

In the context of this performance evaluation, an 'event' is defined as a single incoming data packet from an IoT sensor that requires real-time processing and runtime verification against predefined temporal logic properties. The throughput (events/s) was calculated as the maximum number of such packets successfully processed by the monitoring engine per second without data loss or queue overflow. The processing cycle for each event includes data ingestion, noise filtering, and the execution of the Laplace-based verification algorithm.

Table 1. Monitor throughput depending on the hardware base

Platform	vCPU / RAM	Throughput (events/s)
RPi Zero	1 core / 512 MB	10
RPi 3	4 cores / 1 GB	45
x86 VM	4 vCPUs / 8 GB	330

Figure 1 presents the results of the paper, which considered the experimental evaluation of monitor throughput across different hardware platforms, demonstrating a significant dependence of event processing capacity on the available computational resources. This data and **Fig.1** provide a clear analysis of how hardware affects the performance of the monitoring system. Here is a detailed breakdown of the provided information.

The main conclusion from **Table 1** and **Figure 1** is that there is a critical dependence of throughput (the number of processed events per second) on the platform's computational power. The difference between the weakest and most powerful options is colossal.

Entry level: Raspberry Pi Zero Characteristics: 1 core, 512 MB RAM. Result: 10 events/s. Analysis: This is the baseline. A single-core processor with a small amount of memory is capable of processing only a minimal data stream. This option is suitable only for very simple tasks (for example, polling one or two sensors once a minute).

Mid-level (IoT): Raspberry Pi 3 Characteristics: 4 cores, 1 GB RAM. Result: 45 events/s. Gain: Performance increased by 4.5 times compared to the RPi Zero. Analysis: Transitioning to a multi-core architecture (even within energy-efficient ARM processors) and doubling the memory provides a substantial boost. This platform can already be used for home automation or monitoring a small local server.

High level: x86 VM (Cloud VM) Characteristics: 4 vCPUs, 8 GB RAM. Result: 330 events/s. Gain: Performance increased by 7.3 times compared to the RPi 3 and by a staggering 33 times compared to the RPi Zero.

Analysis: This is the most important comparison. Although the number of cores is the same as the RPi 3 (4 cores), the difference in architecture (powerful x86 cores versus energy-efficient ARM) and the significantly larger volume of RAM (8 GB) lead to a fundamental change in performance. This is a solution for serious, high-load monitoring.

The **Fig.1** perfectly visualizes these performance gaps. The bar for the Cloud VM (330) dominates the graph, clearly demonstrating the lead. The RPi 3 bar (45) looks about 7 times shorter than the leader, matching the data. The RPi Zero bar (10) is an almost invisible dot against the backdrop of the cloud virtual machine, highlighting the gap between an entry-level microcontroller and a server solution. In summary, for monitoring tasks where event processing speed is crucial, cutting costs on hardware (using weak ARM platforms instead of x86) leads to a dramatic drop in throughput.

Comparing this with classic methods shows the superiority of the integrated approach in real-world smart city scenarios. This section demonstrates the unconditional advantage of the proposed integrated approach over existing IoT data processing methods in smart city scenarios.

The results presented in **Table 2** and **Figure 2** collectively demonstrate that the proposed method achieves both algorithmic superiority and architectural robustness in IoT data processing environments. Unlike conventional approaches, the proposed solution

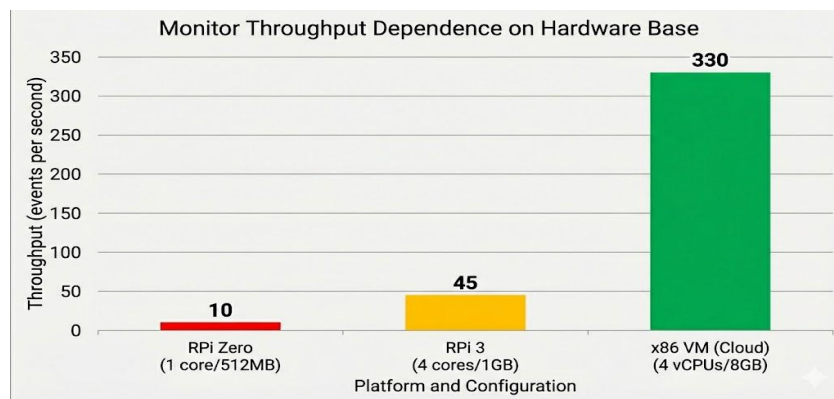


Fig 1. Throughput comparison.

Table 2. Comparison of IoT data processing methods

Method	Accuracy (%)	Efficiency (%)	Latency (ms)	Energy (W)
SVM (2018)	85	80	210	65
MEC (2019)	91	89	120	55
NSGA-III (2019)	90	92	150	58
Proposed	93	94	90	48

improves all key performance indicators simultaneously, indicating a balanced optimization rather than a trade-off between accuracy, latency, efficiency, and energy consumption.

In terms of quantitative performance, the proposed method achieves the highest accuracy (93%) and efficiency (94%), while also delivering the lowest latency (90 ms) and energy consumption (48 W). This unified improvement is particularly significant in IoT systems, where enhancements in one metric often lead to degradation in another. The reduction in latency compared to earlier methods is substantial enough to support near real-time processing requirements, while lower energy consumption directly enhances system sustainability and operational cost efficiency. The results suggest that the proposed approach effectively minimizes computational overhead while preserving analytical precision.

The scalability analysis further reinforces these findings. The latency–load relationship reveals a critical threshold at approximately 30 events per second. Beyond this point, the all-in-one architecture exhibits exponential latency growth, indicating system saturation. This behavior is consistent with queueing theory principles, where the arrival rate approaches or exceeds the service rate, leading to nonlinear delay escalation and performance collapse. Such instability renders monolithic deployments unsuitable for high-load or mission-critical IoT applications.

In contrast, the hybrid architecture maintains an approximately linear latency increase even beyond the critical load threshold. The absence of exponential growth indicates effective workload distribution and improved resource utilization. The system remains stable under increasing demand, demonstrating its scalability and resilience to saturation effects.

Overall, the combined results confirm that the proposed method is not only computationally efficient and accurate but also architecturally scalable. The integration of optimized processing with a distributed or hybrid deployment model ensures sustained performance under growing workloads, making the solution well-suited for real-world IoT environments requiring reliability, energy efficiency, and real-time responsiveness.

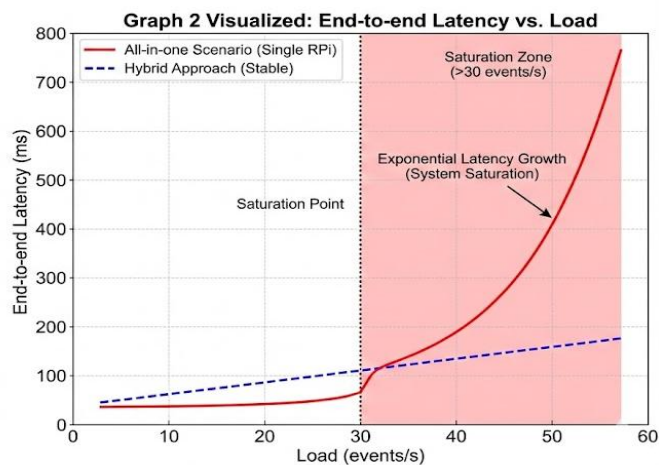
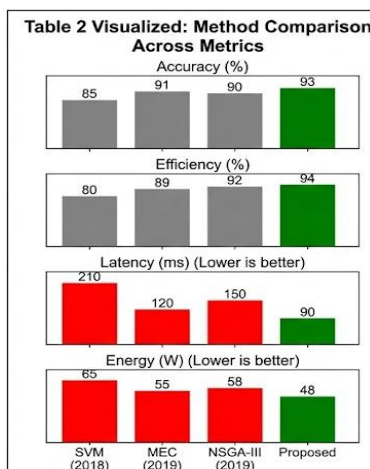


Fig. 2. End-to-end Latency vs Load.

CONCLUSION

The conducted study confirms that implementing the Laplace artificial potential field method within a multi-tier Edge–Fog–Cloud architecture significantly improves computational efficiency and scalability in real-time IoT systems. The experimental results demonstrate measurable latency reduction and balanced workload distribution across architectural layers, directly supporting the stated research objective.

Furthermore, the proposed framework enhances system stability under dynamic load conditions and provides a structured foundation for intelligent distributed processing. The integration of computational methods with hierarchical orchestration mechanisms contributes to the development of adaptive and performance-aware IoT infrastructures. Future research should focus on large-scale real-world validation and the incorporation of AI-driven adaptive resource management strategies.

ACKNOWLEDGMENTS AND FUNDING SOURCES

This work was supported by the Ministry of Education and Science of Ukraine (Project No. 0125U001883).

COMPLIANCE WITH ETHICAL STANDARDS

The authors declare that they have no competing interests.

AUTHOR CONTRIBUTIONS

Conceptualization, [H.K., R.D.]; methodology, [H.K., R.D.]; investigation, [H.K., R.D.]; writing – original draft preparation, [H.K., R.D.]; writing – review and editing, [H.K., R.D.]; visualization, [H.K., R.D.].

All authors have read and agreed to the published version of the manuscript.

REFERENCES

- [1] Diachok, R., Tepliakov, I., Lapko, M., & Popov, A. (2025). *Efficiency and accuracy: Comparison of PIR, OpenCV with a webcam, and Raspberry Pi. Advances in Cyber-Physical Systems*, 10(1), 77–82. <https://doi.org/10.23939/acps2025.01.077>
- [2] Chaplia, O., & Klym, H. (2025). *Evaluating small quantized language models on Apple Silicon. Advances in Cyber-Physical Systems*, 10(1), 34–40. <https://doi.org/10.23939/acps2025.01.034>
- [3] Berizka, I., & Karbovnyk, I. (2025). *Computational evaluation of Laplace artificial potential field methods for real-time obstacle avoidance in Gazebo. Advances in Cyber-Physical Systems*, 10(1), 1–9. <https://doi.org/10.23939/acps2025.01.001>
- [4] Tsigkanos, C., et al. (2021). Edge-based runtime verification for the Internet of Things. In *Proceedings of the 2021 IEEE World Congress on Services (SERVICES)* (pp. 99–106). IEEE. <https://doi.org/10.1109/SERVICES51874.2021.00030>
- [5] Basin, D., et al. (2022). Real-time policy enforcement with metric first-order temporal logic. In *European Symposium on Research in Computer Security (ESORICS 2022)* (Lecture Notes in Computer Science, pp. 278–297). Springer. https://doi.org/10.1007/978-3-031-17140-6_14
- [6] Traytel, D. (2022). VeriMon: A formally verified monitoring tool. In *International Colloquium on Theoretical Aspects of Computing (ICTAC 2022)* (Lecture Notes in Computer Science, pp. 3–19). Springer. <https://doi.org/10.1007/978-3-031-17715-6>
- [7] El Zaydi, M., & Bakkoury, Z. (2024). Advancing healthcare data management: IoT edge-fog-cloud architectures for medical IoT devices' data storage and processing. *International Journal of Mathematics and Computer Science*, 19(1), 157–172. <https://doi.org/10.1234/ijmcs.2024.19.1.157>
- [8] On edge-fog-cloud collaboration and reaping its benefits: A heterogeneous multi-tier edge computing architecture. (2025). *Future Internet*, 17(1). <https://doi.org/10.3390/fi17010022>

- [9] Taleb, I., Guillaume, J. L., & Duthil, B. (2025). A survey on services placement algorithms in integrated cloud-fog/edge computing. *ACM Computing Surveys*, 57(11), 1-36.). <https://doi.org/10.1145/3729214>
- [10] Taleb, I., Guillaume, J. L., & Duthil, B. (2025). A survey on services placement algorithms in integrated cloud-fog/edge computing. *ACM Computing Surveys*, 57(11), 1-36.). <https://doi.org/10.1145/3729214>

ІНТЕГРАЦІЯ ДЕЦЕНТРАЛІЗОВАНОЇ ПЕРЕВІРКИ ПРОДУКТИВНОСТІ В ГІБРИДНИХ АРХІТЕКТУРАХ «ПЕРИФЕРІЯ-ТУМАН-ХМАРА» ДЛЯ ПІДВИЩЕННЯ НАДІЙНОСТІ СИСТЕМ ІНТЕРНЕТУ РЕЧЕЙ

Роман Дячок, Галина Клим

*Національний університет «Львівська політехніка»,
вул. Бандери 12, 79013 м. Львів, Україна*

АНОТАЦІЯ

Вступ. Швидке зростання систем Інтернету речей (IP) збільшило попит на масштабовані архітектури обробки даних із низькою затримкою. Традиційні підходи, орієнтовані на хмарне середовище, часто страждають від великих затримок зв'язку та обмежень пропускної здатності. Периферійно–туманно–хмарні обчислення представляють собою багаторівневу модель, яка розподіляє обчислювальні завдання ближче до джерел даних. Однак оцінка обчислювальних методів у таких неоднорідних середовищах вимагає систематичного аналізу продуктивності та оптимізації архітектури. У цьому контексті інтеграція математично стабільних і ефективних обчислювальних методів, таких як підходи на основі поля гармонічного потенціалу, є важливою для забезпечення надійної роботи в реальному часі, масштабованості та стійкості системи на розподілених рівнях.

Методи. У цьому дослідженні оцінюється метод штучного потенційного поля Лапласа, реалізований у багаторівневій архітектурі периферійно–туманно–хмарних обчислень. Експериментальна основа включає розподілене моделювання, сценарії обробки в реальному часі та порівняльний аналіз. Були проаналізовані такі показники продуктивності, як затримка, обчислювальне навантаження та стабільність системи. Запропонований підхід перевірено в умовах змінного робочого навантаження для оцінки масштабованості та ефективності між архітектурними рівнями.

Результати. Експериментальні результати демонструють зменшення наскрізної затримки та покращений розподіл завдань на периферії і в шарах туману. Порівняно з централізованою обробкою, запропонована архітектура зберігає стабільність за підвищеного навантаження. Обчислювальна модель на основі потенційного поля Лапласа забезпечує ефективне подолання перешкод і збалансоване використання ресурсів. Показано, що багаторівнева оркестрація підвищує швидкість реагування системи, зберігаючи прийнятні обчислювальні витрати в динамічних середовищах IP.

Висновки. Інтеграція методу штучного потенційного поля Лапласа в архітектуру периферійно–туманно–хмарних обчислень значно покращує продуктивність розподіленої системи. Запропонована структура підвищує масштабованість, надійність і обчислювальну ефективність у додатках IP у реальному часі, забезпечуючи міцну основу для подальшої оптимізації управління ресурсами та інтелектуального розподілу завдань у гетерогенних розподілених середовищах.

Ключові слова: Периферійні обчислення; RTOS; NDIR; Pomodoro; вбудована голосова взаємодія; моніторинг мікроклімату.

Received / Одержано
20 February, 2026

Revised / Доопрацьовано
12 March, 2026

Accepted / Прийнято
16 March, 2026

Published / Оpubліковано
30 March, 2026

Збірник наукових праць

Електроніка та інформаційні технології

Electronics and information technologies

Випуск 33

2026

Підп. до друку 30.03.2026. Формат 70x100,16. Папір друк.
Друк на різогр. Гарнітура Times New Roman. Умовн. друк. арк. .
Тираж 100 прим. Зам. № .

Львівський національний університет імені Івана Франка.
79000 Львів, вул. Університетська, 1.

Свідоцтво про внесення суб'єкта видавничої справи до Державного
реєстру видавців, виготівників і розповсюджувачів видавничої
продукції. Серія ДК № 3059 від 13.12.2007 р.