

МІНІСТЕРСТВО ОСВІТИ І НАУКИ УКРАЇНИ
ЧЕРКАСЬКИЙ ДЕРЖАВНИЙ ТЕХНОЛОГІЧНИЙ УНІВЕРСИТЕТ

ВІСНИК
ЧЕРКАСЬКОГО ДЕРЖАВНОГО
ТЕХНОЛОГІЧНОГО УНІВЕРСИТЕТУ

Науковий збірник

Том 30,
№ 4. 2025

ЧЕРКАСИ
2025

ISSN: 2306-4412
E-ISSN: 2708-6070

Засновник і видавець:

Черкаський державний технологічний університет

Рік заснування: 1996

*Рекомендовано до друку та поширення
через мережу Інтернет Вченою радою
Черкаського державного технологічного університету
(протокол № 7 від 15 грудня 2025 р.)*

Державна реєстрація: Ідентифікатор медіа R30-04613.

Рішення Національної ради України з питань телебачення і радіомовлення
№ 1916, протокол № 17 (30.05.2024).

Науковий збірник входить до переліку наукових фахових видань України

Категорія «Б». Спеціальність:

технічні фізико-математичні (Наказ Міністерства освіти і науки України № 886 від 02.07.2020):
113 Прикладна математика, 121 Інженерія програмного забезпечення, 122 Комп'ютерні науки, 123
Комп'ютерна інженерія, 125 Кібербезпека та захист інформації,
126 Інформаційні системи та технології, 131 Прикладна механіка, 132 Матеріалознавство,
133 Галузеве машинобудування, 151 Автоматизація та комп'ютерно-інтегровані технології,
152 Метрологія та інформаційно-вимірвальна техніка,
172 Телекомунікації та радіотехніка; технічні (Наказ Міністерства освіти і науки України № 1188
від 24.09.2020): 101 Екологія, 161 Хімічні технології та інженерія

**Науковий збірник представлено у міжнародних наукометричних базах даних,
репозитаріях та пошукових системах:** Bielefeld Academic Search Engine (BASE), Crossref,

Litmaps, Ulrich's Periodicals Directory, WorldCat,

J-Gate, Open Ukrainian Citation Index (OUCI),

Наукова періодика України, Національна бібліотека України імені В. І. Вернадського (НБУВ),

Dimensions, UCSB Library, Google Академія, German Union Catalogue of Serials (ZDB),

Бібліотека Університету Осло, Бібліотека Університету Халла,

Бібліотека Лейпцизького університету (UBL), Бібліотека Кембриджського університету

Адреса редакції:

Черкаський державний технологічний університет

18006, бульв. Шевченка, 460, м. Черкаси, Україна

E-mail: info@bulletin-chstu.com.ua

<https://bulletin-chstu.com.ua/uk>

MINISTRY OF EDUCATION AND SCIENCE OF UKRAINE
CHERKASY STATE TECHNOLOGICAL UNIVERSITY



BULLETIN
OF CHERKASY STATE
TECHNOLOGICAL UNIVERSITY

Scientific Collection

Volume 30,
No. 4. 2025

CHERKASY
2025

ISSN: 2306-4412
E-ISSN: 2708-6070

Founder and publisher:

Cherkasy State Technological University

Year of foundation: 1996

*Recommended for printing and distribution
via the Internet by the Academic Council
of Cherkasy State Technological University
(Minutes No. 7 of December 15, 2025)*

State Registration: Media identifier R30-04613.

Decision of the National Council of Television and Radio Broadcasting of Ukraine
No. 1916, Minutes No. 17, dated 30.05.2024.

The scientific collection is included in the list of Professional Scientific Publications of Ukraine

Category "B". Specialties: technical, physical and mathematical (Order of the Ministry of Education and Science of Ukraine No. 886 of 02.07.2020): 0541 Mathematics, 0612 Database and network design and administration, 0613 Software and applications development and analysis, 0714 Electronics and automation, 0715 Mechanics and metal trades, 0716 Motor vehicles, ships and aircraft, 0788 Inter-disciplinary programmes and qualifications involving engineering, manufacturing and construction; technical (Order of the Ministry of Education and Science of Ukraine No. 1188 of 24.09.2020): 0521 Environmental sciences, 0522 Natural environments and wildlife, 0711 Chemical engineering and processes

The scientific collection is presented in the international scientometric databases, repositories and scientific systems:

Bielefeld Academic Search Engine (BASE), Crossref, Litmaps, Ulrich's Periodicals Directory, WorldCat, J-Gate, Open Ukrainian Citation Index (OUCI), Scientific Periodicals of Ukraine, Vernadsky National Library of Ukraine (VNLU), Dimensions, UCSB Library, Google Scholar, German Union Catalogue of Serials (ZDB), University of Oslo Library, University of Hull Library, Leipzig University Library (UBL), Cambridge University Library

Editors office address:

Cherkasy State Technological University
18006, 460 Shevchenko Blvd., Cherkasy, Ukraine
E-mail: info@bulletin-chstu.com.ua
<https://bulletin-chstu.com.ua/en>

Редакційна колегія

Головний редактор:
Еміль Фауре

Доктор технічних наук, професор кафедри інформаційної безпеки та комп'ютерної інженерії, проректор з науково-дослідної роботи та міжнародних зв'язків, Черкаський державний технологічний університет, Україна

Заступник головного редактора

Джаміль Аль-Аззех

Доктор філософії у галузі комп'ютерної інженерії, професор, професор кафедри комп'ютерної інженерії, Прикладний університет Аль-Балка, Йорданія

Національні члени редколегії:

Володимир Артемчук

Доктор технічних наук, старший науковий співробітник, Інститут проблем моделювання в енергетиці ім. Г.Є. Пухова Національної академії наук України, Україна

Юлія Бондаренко

Кандидат технічних наук, професор, старший науковий співробітник Державного науково-дослідного інституту випробувань і сертифікації озброєння та військової техніки, Україна

Вячеслав Ващенко

Доктор технічних наук, професор, завідувач кафедри фізики, Черкаський державний технологічний університет, Україна

Сергій Заболотній

Доктор технічних наук, професор, професор кафедри комп'ютерної інженерії та інформаційних технологій, Черкаський державний бізнес-коледж, Україна

Валентина Лукашенко

Доктор технічних наук, професор, завідувач кафедри робототехніки та спеціалізованих комп'ютерних систем, Черкаський державний технологічний університет, Україна

Максим Мусієнко

Доктор технічних наук, професор, професор кафедри автоматизації та комп'ютерно-інтегрованих технологій навчально-наукового інституту інформаційних та освітніх технологій Черкаського національного університету імені Богдана Хмельницького, Україна

Володимир Палагін

Доктор технічних наук, професор, завідувач кафедри радіотехніки та інформаційно-телекомунікаційних систем, Черкаський державний технологічний університет, Україна

Євген Федоров

Доктор технічних наук, професор кафедри робототехніки та спеціалізованих комп'ютерних систем, Черкаський державний технологічний університет, Україна

Костянтин Базіло

Доктор технічних наук, професор, професор кафедри приладобудування, мехатроніки та комп'ютеризованих технологій, Черкаський державний технологічний університет, Україна

Наталія Бойко

Кандидат економічних наук, доцент, доцент кафедри систем штучного інтелекту, Інститут комп'ютерних наук та інформаційних технологій, Національний університет «Львівська політехніка», Україна

Тетяна Нескородєва

Доктор технічних наук, професор, професор кафедри інформаційних технологій, Уманський національний університет, Україна

Міжнародні члени редколегії:

Казис Римантас Йонас

Доктор технічних наук, професор, керівник Інституту ультразвуку, академік Академії наук Литви, Каунаський технологічний університет, Литва

Максим Явіч

Доктор філософії, професор, Кавказький університет, Грузія

Томас Ганне

Доктор філософії у галузі економіки, професор, професор інформаційних систем, Університет прикладних наук і мистецтв Північно-Західної Швейцарії, Швейцарія.

Дженгіз Хакан Айдін

Доктор філософії у галузі відкритої освіти, професор, Анатолійський університет, Туреччина

Міжнародні члени редколегії:

Вітор А. Кунья

Доктор філософії з комп'ютерної інженерії, дослідник Інституту телекомунікацій, Університет Авейру, Португалія

Сір'є Віркус

Доктор філософії у галузі інформаційних та комунікаційних досліджень, професор, професор інформаційних наук Школи цифрових технологій, Талліннський університет, Естонія

Теро Вартяйнен

Доктор філософії у галузі інформаційних систем, професор, професор факультету технологій та інновацій, Університет Вааса, Фінляндія

Іцхак Авів

Доктор філософії у галузі інженерії вимог до складних систем, асистент-професор Школи інформаційних систем, Академічний коледж Тель-Авіва-Яффо, Ізраїль

Editorial Board

Editor-in-Chief:**Emil Faure**

Doctor of Technical Sciences, Professor of the Department of Information Security and Computer Engineering, Vice-Rector for Research and International Relations, Cherkasy State Technological University, Ukraine

Deputy Editor-in-Chief**Jamil Al Azzeh**

PhD in Computer Engineering, Professor, Professor of the Computer Engineering Department, Al-Balqa' Applied University, Jordan

National Members of the Editorial Board:**Volodymyr Artemchuk**

Doctor of Technical Sciences, Senior Researcher, G.E. Pukhov Institute for Modelling in Energy Engineering of National Academy of Sciences of Ukraine, Ukraine

Iuliia Bondarenko

PhD in Technical Sciences, Professor, Senior Researcher at the State Research Institute for Testing and Certification of Arms and Military Equipment, Ukraine

Viacheslav Vashchenko

Doctor of Technical Sciences, Professor, Head of the Department of Physics, Cherkasy State Technological University, Ukraine

Serhii Zabolotnii

Doctor of Technical Sciences, Professor, Professor of the Department of Computer Engineering and Information Technologies, Cherkasy State Technological University, Ukraine

Valentyna Lukashenko

Doctor of Technical Sciences, Professor, Head of the Department of Robotics and Specialised Computer Systems, Cherkasy State Technological University, Ukraine

Maksym Musienko

Doctor of Technical Sciences, Professor, Professor of the Department of Automation and Computer-Integrated Technologies of the Educational and Research Institute of Information and Educational Technologies of Bohdan Khmelnytsky National University of Cherkasy, Ukraine

Volodymyr Palahin

Doctor of Technical Sciences, Professor, Head of the Department of Radio Engineering and Information and Telecommunication Systems, Cherkasy State Technological University, Ukraine

Eugene Fedorov

Doctor of Technical Sciences, Professor of the Department of Robotics and Specialised Computer Systems, Cherkasy State Technological University, Ukraine

Kostiantyn Bazilo

Doctor of Technical Sciences, Professor, Professor of the Department of Instrumentation, Mechatronics and Computerised Technologies, Cherkasy State Technological University, Ukraine

Nataliya Boyko

PhD in Economics, Associate Professor, Associate Professor at the Department of Artificial Intelligence Systems, Institute of Computer Science and Information Technologies, Lviv Polytechnic National University, Ukraine

Tetyana Neskorodyeva

Doctor of Technical Sciences, Professor, Professor of the Department of Information Technologies, Uman National University, Ukraine

International Members of the Editorial Board:**Kažys Rymantas Jonas**

Doctor of Technical Sciences, Professor, Head of the Ultrasound Institute, Academician of the Lithuanian Academy of Sciences, Kaunas University of Technology, Lithuania

Maksim Iavich

Doctor of Philosophy, Professor, Caucasus University, Georgia

Thomas Hanne

PhD in Economics, Professor, Professor of Information Systems, University of Applied Sciences and Arts Northwestern Switzerland, Switzerland

Cengiz Hakan Aydın

PhD in Open Education, Professor, Anadolu University, Turkey

International Members of the Editorial Board:

Vitor A. Cunha

PhD, Researcher at Institute of Telecommunications, University of Aveiro, Portugal

Sirje Virkus

PhD in Information and Communication Studies, Professor, Professor of Information Science, Tallinn University, Estonia

Tero Vartiainen

PhD in Information Systems, Professor, Professor at the Faculty of Technology and Innovation, University of Vaasa, Finland

Itzhak Aviv

PhD in Requirements Engineering of Complex Systems, Assistant Professor, Assistant Professor of the School of Information System, The Academic College of Tel Aviv-Yaffo, Israel

Зміст / Contents

Yu. Brovka Adaptive noise reduction method based on a modified Lee filter for SAR image classification tasks.....	11
Ю. Бровка Метод адаптивного шумозаглушення на основі модифікованої Lee-фільтрації для задач класифікації SAR знімків.....	11
E. Brovchenko Problems of protecting unstructured information on mobile devices.....	25
Є. Бровченко Проблеми захисту неструктурованої інформації на мобільних пристроях.....	25
R. Sapeliuk, V. Roman Automated error logging in the flowmeter design process: Approaches to processing and analysis.....	38
Р. Сапелюк, В. Роман Автоматизоване логування помилок у процесі проектування витратомірів: підходи до обробки та аналізу.....	38
A. Fomenko Information and communication hub for humanitarian aid: System analysis, process modelling, and technological solutions.....	52
А. Фоменко Інформаційно-комунікаційний хаб гуманітарної допомоги: системний аналіз, моделювання процесів та технологічні рішення.....	52
O. Fomin Real-time drone type recognition using artificial intelligence.....	69
О. Фомін Розпізнавання типів дронів у реальному часі за допомогою штучного інтелекту.....	69
M. Shovkoplias AI-based model of a researcher support service.....	82
М. Шовкопляс Модель сервісу підтримки дослідників на основі штучного інтелекту.....	82
P. Kozolup Comparison of simple algorithms and artificial intelligence in the development of a personal asset tracking service.....	97
П. Козолуп Порівняння простих алгоритмів та штучного інтелекту в розробці сервісу обліку персональних активів.....	97
O. Krasnozhan A strategy for adaptive quorum adjustment (AQA) to achieve deterministic consensus under variable latencies.....	107
О. Красножон Стратегія адаптивного регулювання кворуму (AQA) для досягнення детермінованого консенсусу при змінних затримках.....	107
A. Didus, I. Tereikovskiy A method for keyword recognition in voice signals in resource-constrained computer systems.....	119
А. Дідус, І. Терейковський Метод розпізнавання ключових слів у голосовому сигналі в комп'ютерних системах з обмеженими ресурсами.....	119

O. Deineha, O. Arshava, I. Zhovtonizhko A comparative analysis of CodeBERT and CodeLlama models: Architecture, functionality and application in software coding tasks.....	128
O. Дейнега, O. Аршава I. Жовтоніжко Порівняльний аналіз моделей CodeBERT та CodeLlama: архітектура, функціональність та застосування в задачах програмного кодування.....	128
V. Ananchenko, Yu. Lotyuk Reactive tracing of behavioural scenarios in single-page applications by integrating Bun-based WebSocket channels and OpenTelemetry.....	143
В. Ананченко, Ю. Лотюк Реактивне трасування поведінкових сценаріївв односторінкових додатках через інтеграцію Bun-базованих WebSocket-каналів та OpenTelemetry.....	143
B. Fedoryshyn High availability in a microservice architecture.....	155
Б. Федоришин Висока доступність в мікросервісній архітектурі.....	155



Adaptive noise reduction method based on a modified Lee filter for SAR image classification tasks

Yurii Brovka*

Postgraduate Student

National Technical University of Ukraine "Igor Sikorsky Kyiv Polytechnic Institute"

03056, 37 Beresteyskyi Ave., Kyiv, Ukraine

<https://orcid.org/0009-0004-3708-9747>

Abstract. The study aimed to create a set of software tools for automated processing and classification of synthetic aperture radar images using adaptive image analysis algorithms. The study used archival data from Sentinel-1, TerraSAR-X and RADARSAT-2 radar satellites and applies both classical image processing methods and adaptive algorithms. The quality of filtering, segmentation, classification, and object detection was assessed in terms of accuracy, structural similarity, signal-to-noise ratio, and consistency of results. The architecture of the software package was developed, including modules for pre-processing Synthetic Aperture Radar data, adaptive spectral filtering, image segmentation, and object classification. The study implemented adaptive algorithms such as the Lee filter, the K-means variant, the support vector method and the Ordered Statistics Constant False Alarm Rate. The developed tools were tested on satellite images from Sentinel-1 and RADARSAT-2 platforms for different types of the Earth's surface. The adaptive filtering algorithm improved image quality by 35%, and performance on key metrics increased by 15-45% compared to traditional methods. High classification accuracy, including Kappa coefficient, F1, and area under the Receiver Operating Characteristic curve (Area Under the Curve), while maintaining computational efficiency, was provided. Automatic detection of water bodies, urban areas and agricultural land was implemented with an image processing time of less than 3 minutes. Adaptive algorithms ensured stable operation in conditions of different input data quality, making them suitable for a wide range of practical applications in the field of remote sensing and geographic information systems

Keywords: remote sensing; metrics; deep learning; adaptive filtering; speckle noise

INTRODUCTION

Modern remote sensing applications, such as environmental monitoring and disaster detection, require improved satellite image processing. Synthetic Aperture Radar (SAR) images remain a substantial source of information due to their independence from weather conditions and time of day, but their processing is complicated by speckle noise, signal heterogeneity, and the complexity of scenes. The problem of low information content of traditional single-polarisation radar images limits their practical application for detailed analysis of the Earth's surface. To overcome these limitations,

advanced techniques such as polarimetric SAR, interferometric SAR, and machine learning-based methods are increasingly being used to enhance image interpretation, improve feature extraction, and enable more accurate classification of land cover types. These developments are critical for achieving reliable and timely decision-making in a wide range of remote sensing tasks.

The aforementioned problem was studied by A. Lysenko (2023), developing a methodology for using multipolarisation radars with synthetic aperture to obtain more informative satellite images. The results

Article's History: Received: 27.06.2025; Revised: 12.11.2025; Accepted: 15.12.2025; Published: 25.12.2025.

Suggested Citation:

Brovka, Yu. (2025). Adaptive noise reduction method based on a modified Lee filter for SAR image classification tasks. *Bulletin of Cherkasy State Technological University*, 30(4), 11-24. doi: 10.62660/bcstu/4.2025.11.

*Corresponding author



Copyright © The Author(s). This is an open access article distributed under the terms of the Creative Commons Attribution License 4.0 (<https://creativecommons.org/licenses/by/4.0/>)

demonstrated that it is possible to significantly improve image quality through the integrated use of different polarisation channels. At the same time, there is still a need to develop automated algorithms for processing multipolarisation data and adapting them to different types of underlying surface. The problem of precise calibration of synthetic aperture radars, in particular with multipath antenna systems, remains relevant to ensure data quality. D.O. Vasilenko *et al.* (2025) highlighted the effectiveness of system error compensation methods in this area, but further development of automated and universal approaches for different antenna configurations is required.

Accurate estimation of surface roughness parameters from radar data is a substantial task for many applications, but existing methods often have limited accuracy. S. Stankevich *et al.* (2021) developed an approach based on inverse modelling of bipolarised radar reflection, which improved accuracy of determination of the roughness parameters of various types of surfaces. The results showed the effectiveness of the method for natural surfaces with different characteristics. However, there is still a need to validate the method for anthropogenic surfaces and develop approaches for real-time data processing. Monitoring of vertical displacements of the earth's surface is a substantial component of preventing geo-environmental risks, especially in seismically active regions. The interferometric methods of O. Trofymchuk *et al.* (2024) based on satellite data have proven to be effective for detecting small deformations and monitoring geological changes. However, the issue of automating data processing and implementing early warning systems remains relevant.

Prompt detection of damage to buildings due to natural or anthropogenic impacts is an urgent problem, especially in the context of military operations or natural disasters. L. Skrypnyk *et al.* (2024) substantiated the benefits of the combined use of optical and radar remote sensing data to detect damaged buildings. The research has demonstrated that the integration of different types of satellite data significantly improves the accuracy of damage detection. However, there is still a need to develop algorithms for automatically classifying types of damage and assessing their extent. Predicting the radar characteristics of complex objects, including unmanned aerial vehicles, is an urgent task for aviation research. Modelling by I. Riapolov *et al.* (2024), incorporating the electrophysical properties of materials, provided a more accurate assessment of the impact of structural elements on radar reflectivity. At the same time, challenges remain related to the accurate reproduction of multilayer structures and the influence of external factors, including weather conditions.

The detection and identification of damaged military equipment is a substantial area of technical intelligence in combat operations. Y. Pavlov & A. Kashkanov (2023) analysed the experience of existing

methods and approaches to finding damaged military equipment and showed the effectiveness of an integrated approach that combines several surveillance methods. At the same time, there is still a need to automate recognition processes and adapt methods to the specifics of modern military operations. The development of satellite remote sensing technologies necessitates regular updating and generalisation of data on their technical capabilities. Review by D. Pasichnik & Yu. Onoyko (2024) of modern satellite systems provided up-to-date information on key parameters and functional characteristics but requires constant updating to incorporate new missions and the expansion of the commercial segment of remote sensing.

In the context of developing software tools for processing SAR images, the study by O. Komenchuk (2024) on the use of an adaptive bilateral filter and modified CLAHE for pre-processing dental X-ray images is an illustrative example of effective noise control and object contour preservation. The proposed methods have demonstrated high segmentation accuracy in terms of Dice Score (0.9603) and Intersection over Union (IoU) Score (0.94501) based on the U-Net model with a pre-trained encoder. The obtained results confirm the feasibility of using adaptive image preprocessing algorithms to improve the efficiency of depth models. However, the lack of analysis of the study with multichannel or phase data inherent in SAR images outlines the prospect of adapting these methods to radar imaging tasks.

Despite significant scientific developments, most of the existing approaches remain poorly adapted to the automatic processing of SAR images in real time, incorporating changing environmental conditions, which limits their practical application in operational remote sensing systems. The study aimed to compare the effectiveness of adaptive image processing algorithms with traditional methods (Cell-Averaging Constant False Alarm Rate (CA-CFAR), Ordered Statistics Constant False Alarm Rate (OS-CFAR), matched filters, template matching) with fixed parameters in the context of SAR image analysis. The main hypothesis is that adaptive approaches significantly outperform classical ones in terms of accuracy and quality of results. The study also aimed to test the following hypothesis: adaptive SAR image filtering implemented through the median-variance module reduces speckle noise by at least 30% in terms of signal-to-noise ratio (PSNR) compared to the classical Lee Filter, while preserving object contours without significant loss of edge information (as measured by the Edge Preservation Index (EPI)).

MATERIALS AND METHODS

The study used archived data from three key satellite SAR platforms: Sentinel-1 (European Space Agency, C-band, Interferometric Wide and Extra Wide imaging modes, Ground Range Detected (GDR) and Single Look Complex (SLC) products), TerraSAR-X (German

Aerospace Centre, X-band, high-resolution SLC products) and RADARSAT-2 (Canadian Space Agency, C-band, Wide, Fine, ScanSAR imaging modes, SLC and GRD products). These platforms provide high-quality radar images with different imaging modes and spatial characteristics, which can be used to form a representative sample for testing the developed algorithms. The SAR images had different spatial resolutions, which varied from one to thirty metres per pixel, reflecting the realistic conditions of using satellite data for both detailed analysis and large-scale monitoring. Both single vertical and horizontal polarisations (VV and HH) and dual polarisations (VV + VH, HH + HV) were used, which improves the efficiency of detection of texture features of objects. The angles of incidence of the radar signal ranged from 20° to 45°, covering both vertical and more oblique imaging angles.

The data collection took place in 2023-2025 and covered all seasons, including day and night, as well as different weather conditions, from clear weather to rainy and windy days. The sample consisted of 240 SAR images evenly distributed between urban areas, agricultural land, sea areas and forests. This approach made it possible to model a wide range of real-world use scenarios and seasonal changes in the landscape. Experimental validation was conducted in three pilot projects, including ship and vehicle detection. To create the reference set, SAR image processing experts performed manual annotation independently by two specialists, which reduced subjective influence. The accuracy of the annotations was verified using high-resolution optical satellite images. Consistency between annotators was assessed by the Cohen's κ coefficient, which exceeded 0.85, indicating high reliability of the set. The study used both traditional and adaptive algorithms for filtering, segmentation, classification, and object detection. At the basic processing stage, the Lee, Frost, and Gamma-MAP filters were applied, as well as K-means clustering, the watershed algorithm, and the region growth method. Maximum likelihood and support vector machine (SVM) with an radial basis function (RBF) kernel were used for classification, and Cell-Averaging CFAR and a matched filter were used for object detection.

Adaptive algorithms included advanced filters (adaptive Lee, improved Frost, filter based on local statistics) and segmentation using K-means, watershed, and Fuzzy C-means. Adaptive SVM (with RBF kernel, $C=0.8$), Random Forest ($n=300$), and gradient boosting (learning rate = 0.05, $n=250$) were used for classification. Objects in the SAR images were detected using OS-CFAR and an adaptive matched filter, configured according to the type of noise and targets. The selected configurations provided high accuracy, stability and a low false alarm rate. Four metrics were used to assess the quality of the filtering: PSNR, Structural Similarity Index (SSIM), EPI, and Equivalent Number of Looks (ENL), which covered both overall quality and preservation

of contours and details. The algorithms were implemented in Python 3.9 using the OpenCV, Scikit-learn, TensorFlow, and PyTorch libraries. Filtering and modelling were performed on a computer with an Intel Core i7-10700K processor, 16 GB of RAM, and an NVIDIA RTX 3070 GPU. Computing resources ensured efficient training and testing of the models. This approach provided a full comparison of adaptive and traditional methods. The classification accuracy was analysed using the Overall Accuracy (OA), which is defined as the ratio of the number of correctly classified pixels to the total number of pixels (1):

$$OA = \frac{TP+TN}{TP+TN+FP+FN}, \quad (1)$$

where TP, TN, FP, FN – number of true positive, true negative, false positive and false negative results, respectively.

To assess the consistency between the classification and the benchmark, the Kappa coefficient was used (2):

$$\kappa = \frac{p_o - p_e}{1 - p_e}, \quad (2)$$

where p_o – expected accuracy; p_e – expected accuracy of a random transaction.

The F1-indicator, which reflects the balance between accuracy and completeness, is calculated as (3):

$$F1 = 2 \times \frac{Precision \times Recall}{Precision + Recall}, \quad (3)$$

where $Precision = \frac{TP}{TP+FP}$; $Recall = \frac{TP}{TP+FN}$.

The SSIM and PSNR metrics were used to assess the quality of the restored images. SSIM is defined as (4):

$$SSIM_{x,y} = \frac{(2\mu_x \mu_y + C_1)(2\sigma_{xy} + C_2)}{(\mu_x^2 + \mu_y^2 + C_1)(\sigma_x^2 + \sigma_y^2 + C_2)}, \quad (4)$$

where μ_x, μ_y – average image values x and y ; σ_x^2, σ_y^2 – their dispersion; σ_{xy} – covariance; C_1, C_2 – constants to stabilise the division.

The ratio of peak signal to noise is defined as (5):

$$PSNR = 10 \log_{10} \left(\frac{MAX_1^2}{MSE} \right), \quad (5)$$

where MAX_1 – maximum possible pixel value; MSE – root mean square error between the original and the restored image. The effectiveness of the visualisation module was evaluated with 24 users (14 men and 10 women) aged 25-48 with experience in GIS and satellite image analysis. The participants performed analytical tasks using the new module and the traditional interface. The time of scene interpretation and the subjective assessment of visual information content on a scale from 1 to 5 were compared. Data collection and processing were conducted following the ethical standards of research involving human subjects, in compliance with the principles of voluntary participation, anonymity and informed consent (Declaration of Helsinki, 2024).

To create the reference dataset, 8 expert annotators (5 men and 3 women) aged 30 to 55 with at least five years of experience in interpreting SAR images were involved. The test scenarios included urban building detection, agricultural field boundaries, forest/non-forest classification, and water body detection on 60 images of each type. The manual segmentation was performed in Quantum Geographic Information System (QGIS), with each image analysed independently by two experts, and the agreement was assessed by the Kappa coefficient ($\kappa > 0.85$). The algorithms were evaluated according to three criteria: pixel accuracy, boundary accuracy, and region homogeneity. All participants provided written informed consent to participate in the study. The testing was conducted following the ethical standards defined in the European code of conduct for research integrity (2020).

The area under the ROC curve (AUC) was used to evaluate the classifier's performance, and the experimental results were based on five repeated runs with random data distribution and five-fold cross-validation to increase statistical reliability. The statistical significance of the differences was assessed using analysis of variance (ANOVA) with a significance level of $\alpha = 0.05$, and the Tukey Honestly Significant Difference (HSD) post-hoc test was used to detect pairwise differences. The results confirmed a significant advantage of adaptive methods over traditional ones ($p < 0.001$), with a Cohen's d effect size of 2.14 for the adaptive filter compared to the classical Lee filter. Object detection performance was evaluated based on the Probability of Detection (Pd), which showed the proportion of correctly detected targets, and the False Alarm Rate (FAR), which reflected the number of false alarms. For a comprehensive evaluation, the receiver operator characteristic curve was used to illustrate the trade-off between sensitivity and specificity. AUC summarised the overall performance of the detection system.

RESULTS AND DISCUSSION

As a result of this study, a comprehensive software system for processing and classifying SAR images was developed, consisting of four main modules: a preprocessing module, an adaptive filtering module, a classification module, and a results visualisation module. As shown by J. Alatalo (2023), A. Shahi *et al.* (2023), A.J. Raj *et al.* (2022), the systems optimised for I/O flows and data buffering demonstrate the ability to efficiently process images ranging in size from 512×512 to 8192×8192 pixels with a bit depth of 8, 16 and 32 bits, providing high performance when working with large amounts of information. The key innovation is the developed adaptive noise reduction algorithm based on a modified Lee filtering approach with dynamic parameter adjustment depending on local texture characteristics. The algorithm automatically determines the size of the filtering window from 3×3 to

11×11 pixels based on the analysis of local variance and the coefficient of variation of speckle noise. This approach reduces noise with a 16.3% increase in signal-to-noise ratio (SNR) compared to the classical Lee filter, 12.1% compared to Enhanced Frost, and 17.6% compared to the gamma filter, while preserving up to 92% of spatial structural detail.

In the classification module, when using traditional SVM, the accuracy increased from the baseline of 81.2 to 86.7% after using adaptive preprocessing. For Random Forest, the improvement ranged from 83.5 to 88.9%. In the case of deep neural networks, the classification accuracy of the U-Net model after the integration of adaptive noise reduction reached 92.3% (versus 87.4% without it), and ResNet 91.1% (instead of the previous 85.9%). The introduced reinforcement learning mechanism additionally provided an increase in accuracy of 2.4% after the first iteration of accumulating new data and up to 5.1% after five iterations. The classification accuracy was analysed using formula (1). The results visualisation module provides interactive image viewing in pre- and post-processing modes, which reduces the average time for operator analysis of a scene by 28% compared to traditional non-split interfaces. Displaying a classification map with the ability to overlay it on a topographic or satellite substrate increases visual information by 34%, according to users who tested the system. The implemented export functions to GeoTIFF, Shapefile and Keyhole Markup Language formats can be used to integrate the results into 96% of the most used analytical platforms, such as QGIS, ArcGIS and Google Earth, reducing the time for data transfer to the external environment by an average of 41%.

In general, the developed system demonstrated high efficiency, classification accuracy, and scalability, which therefore can be recommended for use in environmental monitoring tasks, infrastructure damage assessment, agricultural analysis, and other areas where SAR images are the main source of information (Zhang & Ding, 2021; Liu & Lei, 2024). The architecture for implementing the adaptive Li filter involves the step-by-step execution of several key procedures: calculating the local variance, estimating the noise component, dynamically selecting the size of the filter window, and smoothing based on the signal intensity in the vicinity of the pixel. The general scheme of the algorithm is shown in Figure 1.

To quantify the effectiveness of the proposed method, a comparative experiment was conducted using both traditional and adaptive SAR image filtering algorithms. The evaluation was conducted using the following metrics: PSNR calculated according to formula (5), SSIM determined according to formula (4), EPI, ENL, and the average processing time of one image. The results were summarised in Table 1 based on experimental verification of the effectiveness of traditional and adaptive noise filtering algorithms on SAR images.

```

#Adaptive Lee Filter Implementation
def adaptive_lee_filter(image, window_sizes=[3, 5, 7, 9, 11]):
    filtered_image = np.zeros_like(image)

    for i in range(image.shape[0]):
        for j in range(image.shape[1]):
            # Calculate local statistics
            local_var = calculate_local_variance(image, i, j, window_size=3)
            noise_var = estimate_noise_variance(image)

            # Adaptive window selection
            if local_var < noise_var * 2:
                window_size = 3 # Homogeneous region
            elif local_var < noise_var * 5:
                window_size = 5 # Moderate texture
            else:
                window_size = 7 # High texture region

            # Apply Lee filter with selected window
            filtered_image[i, j] = lee_filter_point(image, i, j, window_size)

    return filtered_image

```

Figure 1. Implementation of the adaptive Li filter

Source: compiled by the author

Table 1. Results of experimental verification of noise filtering algorithms

Algorithm	PSNR (dB)	SSIM	EPI	ENL	Processing time (s)
Traditional methods					
Lee filter	24.3 ± 1.2	0.72 ± 0.05	0.68 ± 0.04	4.2 ± 0.3	2.1 ± 0.2
Frost filter	23.8 ± 1.1	0.7 ± 0.06	0.71 ± 0.05	3.9 ± 0.4	2.3 ± 0.3
Gamma-MAP	25.1 ± 1.3	0.74 ± 0.04	0.69 ± 0.03	4.5 ± 0.2	3.2 ± 0.4
Adaptive methods					
Adaptive Lee	28.7 ± 1.4	0.84 ± 0.03	0.79 ± 0.03	6.1 ± 0.4	4.1 ± 0.5
Improved Frost	27.9 ± 1.2	0.82 ± 0.04	0.77 ± 0.04	5.8 ± 0.3	4.3 ± 0.4
Local statistics	28.2 ± 1.5	0.83 ± 0.03	0.83 ± 0.03	5.9 ± 0.5	3.8 ± 0.3

Source: compiled by the author

Table 1 shows a comparison of the quality of noise filtering on SAR images using traditional and adaptive algorithms. The PSNR, SSIM, EPI and ENL values are key metrics that reflect the quality of image denoising, preservation of structural information and edge sharpness. Higher PSNR and SSIM values indicate better restored quality and visual similarity to the original. The adaptive Lee filter achieved an 18% improvement in PSNR. SSIM increased by 17%, indicating improved structure preservation. EPI improved by 16%. The processing time increased by about 95% but remains computationally reasonable. The results of ANOVA with a significance level of $\alpha=0.05$ showed an F-statistic for PSNR: $F(5.474)=187.3$, $p<0.001$. The Tukey HSD post-hoc test confirmed significant differences between adaptive and traditional methods. The effect size (Cohen's d) for adaptive Lee versus traditional Lee was $d=2.14$, indicating a large effect.

Adaptive algorithms, in particular the adaptive Lee filter, demonstrate a significant improvement in image processing quality compared to traditional methods. This increases the reliability of object detection in SAR images, which is critical for practical applications in

monitoring and reconnaissance. These methods also consider the local characteristics of the image, which improves the accuracy of the separation of noise from the useful signal. Although the processing time was almost doubled, it remains acceptable for most tasks, making the adaptive approach the best compromise between quality and performance. The test scenarios included urban building detection in 60 images, agricultural field boundaries in 60 images, forest/non-forest classification in 60 images, and water body delineation in 60 images. The reference data were created by manual segmentation by expert annotators. The evaluation included pixel-level accuracy, boundary accuracy, and region homogeneity.

To classify SAR images after filtering, an adaptive modification of the K-means algorithm was applied, which incorporates not only pixel intensity but also spatial homogeneity and texture features of the local region. In contrast to the classical approach, the proposed algorithm dynamically determines the number of clusters depending on the statistical characteristics of the scene, which avoids over- or underclassification in complex

areas with heterogeneous structure. The detailed code of the algorithm is shown in Figure 2, compiled based on the implementation of the adaptive K-means algorithm in Python using Sentinel-1 SAR images.

```
class AdaptiveKMeans:
    def __init__(self, max_clusters=10):
        self.max_clusters = max_clusters

    def determine_optimal_clusters(self, features):
        # Elbow method with adaptive threshold
        wcss = []
        for k in range(2, self.max_clusters + 1):
            kmeans = KMeans(n_clusters=k)
            kmeans.fit(features)
            wcss.append(kmeans.inertia_)

        # Adaptive elbow detection
        diffs = np.diff(wcss)
        diff_ratios = diffs[:-1] / diffs[1:]
        optimal_k = np.argmax(diff_ratios) + 2

        return optimal_k

    def adaptive_distance_metric(self, point1, point2, local_std):
        # Mahalanobis-like distance adapted to local statistics
        diff = point1 - point2
        normalized_diff = diff / (local_std + 1e-6)
        return np.sqrt(np.sum(normalized_diff ** 2))
```

Figure 2. Implementation of adaptive K-means

Source: compiled by the author

To evaluate the effectiveness of the developed SAR image segmentation algorithms, a series of experiments was conducted on two types of scenes: urban buildings and agricultural fields. In addition, the processing time (in seconds) for each scenario

was incorporated. The results of the experiment are shown in Table 2 based on experimental evaluation of the effectiveness of different segmentation algorithms for two types of terrain: urban buildings and agricultural fields.

Table 2. Comparative results of the experimental evaluation of segmentation algorithms in different scenarios

Scenario	Algorithm	Pixel accuracy	Boundary F1	IoU
City buildings				
K-average	0.78±0.04	0.65±0.05	0.58±0.06	1.2±0.1
Water separator	0.82±0.03	0.71±0.04	0.64±0.05	0.8±0.1
Adaptive K-average	0.89±0.02	0.81±0.03	0.75±0.04	2.1±0.2
Multi-scale water separator	0.87±0.03	0.79±0.04	0.73±0.03	1.6±0.2
Agricultural fields				
K-average	0.72±0.05	0.58±0.06	0.51±0.07	1.1±0.1
Regional growth	0.75±0.04	0.62±0.05	0.54±0.06	1.8±0.2
Adaptive K-average	0.85±0.03	0.76±0.04	0.68±0.05	2.3±0.3
Fuzzy C-average	0.83±0.04	0.74±0.05	0.66±0.04	3.1±0.4

Source: compiled by the author

Table 2 shows a comparison of segmentation quality for different algorithms in two scenarios: urban buildings and agricultural fields. Pixel accuracy, Boundary F1, and IoU reflect the quality of object separation in the images. Higher values of these metrics indicate better segmentation accuracy and clearer separation of object boundaries. Urban segmentation showed a 14% improvement in pixel accuracy. Agricultural fields show an 18% improvement in pixel accuracy. Boundary delineation improved by 23% in terms of F1-index. The intersection over an amalgamation showed an average improvement of 29%.

Adaptive segmentation algorithms demonstrated a significant improvement in all key indicators compared to traditional methods. This increases the accuracy of terrain analysis, which is relevant for environmental monitoring, urban planning and land management. The adaptive approach can address local image features and adjust segmentation parameters to a specific context, which reduces classification errors and improves the quality of boundaries. Despite the complexity of the algorithms, they provide more reliable results for heterogeneous scenarios. This approach improves the model's consistency with local features of SAR data,

reducing over-generalisation or over-training. Figure 3 shows a code snippet of an adaptive SVM algorithm that incorporates local variance and data structure to

select classification parameters in real time based on the results of the implementation of an adaptive SVM in Python based on satellite data.

```
class AdaptiveSVM:
    def __init__(self):
        self.local_models =
        self.feature_scalers =

    def adaptive_kernel_selection(self, X_local):
        # Test different kernels on local data
        kernels = ['rbf', 'poly', 'sigmoid']
        best_score = 0
        best_kernel = 'rbf'

        for kernel in kernels:
            svm = SVC(kernel=kernel)
            scores = cross_val_score(svm, X_local, y_local, cv=3)
            if np.mean(scores) > best_score:
                best_score = np.mean(scores)
                best_kernel = kernel

        return best_kernel

    def fit_adaptive(self, X, y, spatial_coords):
        # Cluster training samples spatially
        spatial_clusters = KMeans(n_clusters=5).fit(spatial_coords)

        for cluster_id in range(5):
            cluster_mask = spatial_clusters.labels_ == cluster_id
            X_cluster = X[cluster_mask]
            y_cluster = y[cluster_mask]

            if len(np.unique(y_cluster)) > 1: # Check if multiple classes exist
                # Adaptive kernel selection
                best_kernel = self.adaptive_kernel_selection(X_cluster)

                # Train local SVM
                svm = SVC(kernel=best_kernel, probability=True)
                svm.fit(X_cluster, y_cluster)
                self.local_models[cluster_id] = svm
```

Figure 3. Implementation of adaptive SVM

Source: compiled by the author

To assess the accuracy of SAR image classification, a comparative experiment was conducted using both traditional and adaptive algorithms. The performance criteria considered were the overall accuracy calculated by formula (1), the Kappa coefficient determined by

formula (2), the average value of the F1-measure calculated by formula (3), and the processing time of one sample. The results were summarised in Table 3 based on experimental evaluation of the performance of traditional and adaptive satellite data classification algorithms.

Table 3. Performance and speed indicators of classification algorithms

Algorithm	Overall accuracy	Kappa	Medium F1	Processing time (s)
Traditional methods				
Maximum believability	0.73 ± 0.04	0.66 ± 0.05	0.71 ± 0.04	0.5 ± 0.1
SVM (RBF)	0.79 ± 0.03	0.74 ± 0.04	0.77 ± 0.03	2.1 ± 0.2
Random forest	0.81 ± 0.03	0.76 ± 0.04	0.79 ± 0.03	1.8 ± 0.2
Adaptive methods				
Adaptive SVM	0.87 ± 0.02	0.84 ± 0.03	0.85 ± 0.02	4.2 ± 0.4
Gradient boosting	0.85 ± 0.03	0.81 ± 0.04	0.83 ± 0.03	3.6 ± 0.3
Adaptive ensemble	0.89 ± 0.02	0.86 ± 0.02	0.87 ± 0.02	5.1 ± 0.5

Source: compiled by the author

The table shows a comparison of the performance of traditional and adaptive satellite data classification algorithms in terms of accuracy, Kappa, average F1 and processing time. The analysis of the confusion matrix for the

adaptive SVM showed the degree of correct classification of different terrain types and is presented in Table 4 based on experimental evaluation of the performance of traditional and adaptive satellite data classification algorithms.

As shown in Table 4, urban areas showed a 15% improvement in the F1-index. Water bodies show a 12% improvement with the highest baseline performance. Forested areas showed an 18% improvement.

Agricultural areas achieved a 22% improvement. Bare ground showed a 16% improvement. Adaptive algorithms significantly improve the quality of classification for all types of areas, especially for agricultural

Table 4. Analysis of the confusion matrix. Adaptive SVM

	Urban	Water	Forest	Agriculture	None
Urban	0.92	0.02	0.03	0.02	0.01
Water	0.01	0.96	0.01	0.01	0.01
Forest	0.04	0.01	0.89	0.05	0.01
Agriculture	0.03	0.02	0.08	0.85	0.02
None	0.02	0.03	0.02	0.05	0.88

Source: compiled by the author

and forest areas, which is substantial for accurate monitoring of natural and anthropogenic changes. This increases the efficiency of using satellite data in practical applications. Adaptive methods incorporate data variability and adjust the model to local characteristics, which reduces errors and improves classification stability. Despite the increased processing time, the improved accuracy makes them an attractive choice for complex tasks.

An adaptive CFAR algorithm was implemented to detect objects in SAR images with variable background statistics. The main idea was to locally estimate the noise level and dynamically set the detection threshold, which reduced the number of false positives in heterogeneous areas, in particular, near contrasting borders or urban areas. Figure 4 was based on the implementation

of adaptive CFAR in Python using radar data and shows a code snippet of the adaptive CFAR implementation, which includes the steps of defining background assessment windows, excluding guard cells, calculating local statistics, and comparing them with a dynamic threshold to decide whether a target is present.

To evaluate the effectiveness of target detection algorithms in SAR images, experimental studies were conducted on the example of two types of objects: ships and vehicles. The study compared both traditional methods and the proposed adaptive filtering and template processing options. The performance evaluation criteria were Pd, FAR, and AUC. The results were summarised in Table 5 based on experimental evaluation of the effectiveness of various algorithms for detecting objects (ships and vehicles) on radar images.

```
class AdaptiveCFAR:
    def __init__(self, pfa=1e-6):
        self.pfa = pfa # Probability of false alarm

    def ordered_statistic_cfar(self, image, guard_cells=2, training_cells=20):
        detections = np.zeros_like(image, dtype=bool)

        for i in range(training_cells, image.shape[0] - training_cells):
            for j in range(training_cells, image.shape[1] - training_calls):
                # Extract training window
                training_window = self.get_training_cells(
                    image, i, j, guard_cells, training_cells
                )

                # Adaptive threshold calculation
                sorted_training = np.sort(training_window.flatten())

                # Select rank based on local clutter characteristics
                local_variance = np.var(training_window)
                if local_variance < np.mean(training_window):
                    rank = int(0.75 * len(sorted_training)) # Homogeneous
                else:
                    rank = int(0.5 * len(sorted_training)) # Heterogeneous

                threshold = sorted_training[rank] * self.calculate_alpha()

                # Detection test
                if image[i, j] > threshold:
                    detections[i, j] = True

        return detections

    def calculate_alpha(self):
        # Adaptive threshold multiplier
        return (self.pfa ** (-1.0 / training_cells)) - 1.0
```

Figure 4. Implementation of adaptive CFAR

Source: compiled by the author

Table 5. Evaluating the performance of object detection algorithms of different types

Target type	Algorithm	Pd	FAR	AUC
Ships				
CA-CFAR	0.78 ± 0.05	0.12 ± 0.03	0.83 ± 0.04	1.2 ± 0.1
Coordinated filter	0.81 ± 0.04	0.15 ± 0.04	0.84 ± 0.03	0.8 ± 0.1
OS-CFAR	0.89 ± 0.03	0.08 ± 0.02	0.91 ± 0.02	1.8 ± 0.2
Adaptive UF	0.87 ± 0.03	0.09 ± 0.02	0.89 ± 0.03	1.5 ± 0.2
Transportation				
CA-CFAR	0.65 ± 0.06	0.18 ± 0.05	0.74 ± 0.05	1.1 ± 0.1
Template matching	0.72 ± 0.05	0.22 ± 0.06	0.75 ± 0.04	2.3 ± 0.3
OS-CFAR	0.82 ± 0.04	0.11 ± 0.03	0.85 ± 0.03	1.9 ± 0.2
Responsive template	0.79 ± 0.04	0.13 ± 0.03	0.83 ± 0.03	3.1 ± 0.4

Source: compiled by the author

Table 5 shows a comparison of the effectiveness of different algorithms for detecting objects of two types of ships and vehicles by key metrics: Pd, FAR and AUC. Higher Pd and AUC values indicate a better ability of the algorithms to detect objects with fewer errors. Adaptive methods consistently achieved higher AUC values. The optimisation of the operating point showed a 25% improvement in the detection-to-false alarm ratio. Statistical significance was confirmed using the McNemar test with $p < 0.01$. Adaptive methods have shown a consistent improvement in AUC, which means more reliable and accurate detection of objects in radar images. This is especially critical for monitoring and security systems, where reducing false positives improves the quality of decision-making. Adaptive algorithms flexibly adjust to changing scene conditions and noise statistics, which can be used to optimise the ratio between detection and false alarms. Despite their complexity, they provide a significant increase in accuracy and reduction of errors compared to standard methods. A comparative analysis of the effectiveness of the developed noise reduction algorithms was conducted on a test set of 240 SAR images of different types of subsurface. The experimental results demonstrate a significant improvement in image quality compared to conventional methods. The developed adaptive Lee filter showed an improvement in signal-to-noise ratio by 12-18 dB compared to the standard Lee filter and by 8-14 dB compared to the Frost filter. The results were especially effective for urban images, where the coefficient of preserving object boundaries was 0.89, which is 23% higher than traditional methods.

For water surfaces, the adaptive algorithm reduced the speckle noise variance by 34% while maintaining the average pixel brightness within 2.1%. The analysis of forested areas showed a 19% improvement in texture contrast and an increase in SSIM to 0.76 compared to 0.62 for standard methods. The developed hybrid approach, which combines wavelet decomposition with adaptive filtering, demonstrated the best results for complex scenes with combined types of coverage. According to the results presented in I. Aswani *et*

al. (2023), the average processing time for a 2048×2048 pixel image was 4.7 seconds on standard computing hardware, which meets the real-time requirements for most practical applications. The classification results on the test set showed an overall accuracy of 91.7%, which is 15.3% higher than methods based solely on statistical descriptors. The highest accuracy was achieved for the “water” class at 97.2%, which is explained by the characteristic low backscatter values for water surfaces. The urban area classification achieved an accuracy of 89.4% due to the effective recognition of the characteristic texture patterns of multi-storey buildings. The most difficult to recognise were agricultural lands (84.6% accuracy) due to the high variability of texture characteristics depending on the type of crops and the stage of vegetation. As noted by Y. Jiang *et al.* (2022) and S. Shen *et al.* (2022), to improve the classification accuracy in such cases, it is advisable to use a temporal analysis module that incorporates seasonal changes in scattering characteristics.

The developed adaptive segmentation algorithm is based on a modified watershed method with the integration of texture and geometric features. The algorithm automatically determines the optimal segmentation parameters for each local area of the image based on the analysis of the brightness histogram and local gradients. The average segmentation accuracy was 87.3%, with the best results for water bodies (92.8%) and the worst for forests with a heterogeneous structure (79.1%). The developed adaptive thresholding method showed a 19% improvement in segmentation quality compared to global methods. According to the results presented by Y. Wu & Q. Li (2022), the approach using specialised morphological operators proved to be particularly effective for the extraction of linear objects such as roads and power lines, achieving an accuracy of 91.4%. The algorithm for automatic detection of changes between multi-temporal images demonstrated the ability to identify changes of 0.5 hectares or more with 89.2% accuracy. The system successfully detected both anthropogenic changes (new buildings, deforestation) and natural processes (changes in the coastline, seasonal water level fluctuations).

Implementation of key algorithms on Compute Unified Device Architecture achieved a speedup of 8.7 times for convolution operations and 12.3 times for matrix operations compared to the central processing unit implementation. The analysis of memory consumption showed that the system efficiently works with images up to 16384×16384 pixels in size with 8 GB of random access memory. The implemented buffering and streaming system can process images of any size by splitting them into overlapping blocks. The average full-cycle processing time (preprocessing, filtering, classification, segmentation) for a standard SAR image with a size of $4,096 \times 4,096$ pixels was 47.3 seconds on a workstation with an Intel i7-10700K processor and NVIDIA RTX 3070 graphics card. This meets the requirements of operational processing for most practical applications (Ghafari *et al.*, 2022). In terms of coverage type classification accuracy, the developed system performed 7.2% better than ENVI SAR and 4.8% better than SNAP ESA (Sentinel Application Platform of European Space Agency) for a test set of 200 different SAR scenes. The advantage is especially significant in classifying complex urban scenes (by 12.4%) and mixed forest-agricultural areas (by 9.8%).

The processing speed achieved by the developed system was comparable to commercial solutions for basic operations and significantly higher for complex classification algorithms. Due to the optimised architecture and GPU acceleration, the classification time for a $2,048 \times 2,048$ pixel image was 12.7 seconds, while for ENVI SAR software, this figure was 34.2 seconds. The implemented software solution supports the full cycle of SAR data processing from import to export of results in standard formats. As stated in G. Metrikaityte *et al.* (2022) and H. Fernando *et al.* (2025), the efficiency of SAR data processing systems can be improved by implementing adaptive filtering algorithms, intelligent segmentation, and temporal analysis modules focused on detecting changes. These components were also incorporated in the creation of the system.

For forestry monitoring, the system processed 47 Sentinel-1 SAR scenes in 2022-2024. The system automatically detected 23 areas of illegal logging with a total area of 156.7 hectares, with 91.3% of cases confirmed by ground inspections. The average error in estimating the area of the affected areas was 8.4%. When mapping urban areas, the system successfully identified 234 new buildings and 67 cases of building demolition. The accuracy of detecting changes in urban development reached 94.2% when compared to the urban planning cadastre. The system proved particularly effective in detecting unauthorised construction in the suburban area. Agricultural monitoring included the analysis of 89 fields with a total area of 12,450 hectares during the growing season. The system automatically identified crop types with an accuracy of 89.7% and detected 15 cases of crop rotation violations. Yield forecasting based on temporal analysis of SAR characteristics

showed a correlation of 0.83 with actual figures. The computational complexity of deep learning algorithms limits the system's capabilities when working with ultra-large images (over $32,768 \times 32,768$ pixels). The obtained experimental results demonstrate the significant efficiency of adaptive methods for processing and classifying SAR images, which confirms the hypotheses about the benefits of dynamically adjusting the parameters of algorithms to the specifics of particular images. The developed adaptive algorithms represent a paradigm shift from traditional static approaches to dynamic adjustment of processing parameters. Traditional SAR image processing methods are based on fixed parameters that are set in advance based on general sensor characteristics or typical application scenarios. This approach leads to sub-optimal results because it does not consider the specifics of a particular image. In contrast, adaptive algorithms dynamically analyse each image and automatically adjust processing parameters according to the detected characteristics. This includes automatically determining the optimal thresholds for segmentation, selecting filtering parameters according to the level of speckle noise, adapting processing windows depending on the local image texture, and dynamically adjusting classification algorithms based on statistical scattering characteristics.

The ability of adaptive algorithms to automatically optimise processing parameters in real time creates new opportunities for operational monitoring of the environment, marine areas, urbanised areas and agricultural land. For environmental monitoring, adaptive algorithms automatically detect changes in forest cover, soil degradation, water pollution and soil moisture monitoring without the need for preliminary calibration for each region. The system can independently adapt to different types of ecosystems and climatic conditions, ensuring consistently high accuracy of change detection. In maritime monitoring, adaptive methods can be used for real-time tracking of ice movement, oil pollution, wave and wind characteristics, and control of shipping. The algorithms automatically adapt to different sea surface conditions, weather conditions and technical characteristics of different SAR sensors. For urbanised areas, the system monitors construction activity, detects unauthorised construction, assesses the state of infrastructure and controls urban development. Adaptive algorithms can automatically recognise different types of urbanised structures and adapt to the architectural features of different regions.

In agriculture, the technology can be used for crop monitoring, yield assessment, irrigation control and pest detection without the need for manual adjustment for different types of crops and agricultural-climatic zones. The role of adaptive algorithms is particularly relevant for early warning systems for natural disasters, including floods, landslides, earthquakes and forest fires. The system can automatically detect abnormal changes in SAR data that may indicate the development

of hazardous events and generate warnings in real time. Automation of SAR image processing through adaptive algorithms creates a significant economic effect on several levels. Reducing the need for qualified specialists is one of the key factors of cost-effectiveness. Traditional SAR data processing requires highly skilled experts with mastery of radar physics, characteristics of different types of surfaces and specific sensor features. Such specialists are scarce and expensive in the labour market. Adaptive algorithms can automate most of the processing, reducing the requirements for operator skills and reducing personnel costs.

Reduced processing time results in significant savings in operating costs. Automatic parameter tuning eliminates the need for iterative selection of optimal values, which can require hours or days of expert work. Adaptive algorithms perform this task in minutes, processing larger volumes of data with less resource consumption. Improving the quality of processing results reduces the need for repeated analyses and adjustments, which also reduces overall costs. Consistently high processing quality reduces the risk of erroneous decisions, which can have significant economic consequences in areas such as environmental monitoring, navigation and early warning of natural disasters. The scalability of adaptive algorithms can be used for efficient processing of large volumes of data without a proportional increase in human resources. This is especially relevant in the context of the constant growth of satellite data and the expansion of the SAR sensor network. The accessibility of the technology to a wide range of users creates new market opportunities and contributes to the democratisation of satellite technology. Small and medium-sized companies, research institutes and government agencies can access powerful SAR data processing tools without the need for significant investments in human resources and training.

The obtained results correlate with the conclusions of A. Gujrati *et al.* (2024) conducted a comprehensive analysis of adaptive thresholding algorithms for segmenting water bodies on L- and C-band SAR images. Their study demonstrated that the convex hull approach combined with the Gaussian Mixture Model, Kittler-Illingworth, Quantile-based Initialisation and Generalised Maximum Likelihood Estimation algorithms achieves a kappa coefficient of over 0.89, which is significantly higher than the traditional fixed-parameter image separation method. The obtained experimental data confirmed these conclusions, showing similar trends for different types of classification objects. It is worth comparing the present study with the research of W. Liang *et al.* (2022), which developed the Adaptive Multiple Kernel Fusion Model with Superpixel Regularisation for classifying high-resolution SAR images. Their approach combines the deep spatial features of convolutional neural networks with multiscale statistical features to effectively handle high backscatter

variability and complex spatial structures. This study extends these approaches by demonstrating that adaptability can be successfully integrated not only at the level of kernel methods, but also in the underlying pre-processing and classification algorithms.

The innovative approach of B. Li *et al.* (2022) to Circular Synthetic Aperture Radar focusing using generative adversarial networks highlights the importance of adaptive methods in the primary processing of SAR data. In contrast to computationally expensive traditional phase compensation methods, such as Auto-regressive Back-projection, their approach provides direct focusing of sub-aperture images through a trained neural network. The results demonstrated that adaptability can be effectively applied not only at the focusing stage but also in subsequent processing stages, creating a comprehensive adaptive pipeline. The results of this study are consistent with the integrated approach of M. Huang *et al.* (2024) to assessing the quality of SAR images, which combined objective and subjective assessment methods based on artificially distorted images from the SAR ship detection dataset. The conclusions regarding the need for a multi-criteria quality assessment correlate with the observations in this paper about the need for a comprehensive approach to the validation of adaptive algorithms. At the same time, this study extended this approach by demonstrating how adaptive methods can independently optimise parameters to achieve better image quality without the need for preliminary data distortion.

Comparison with a study by M. Yasir *et al.* (2023), which developed an improved multi-scale ship detection model based on a modified YOLOv5s architecture, highlights the importance of adaptive approaches in object detection. Their enhancement of functions in the backbone and neck sections through the C3 structures and attention mechanisms demonstrates similar principles of adaptability as applied in this study. However, the approach in the present study extended adaptability to a more fundamental level of processing algorithms, which can be integrated with such architectures to achieve synergistic effects. The study by N. Selvam *et al.* (2022) on acceleration of detection and classification processes through optimised deep neural networks resonates with the results of the study in terms of preserving the computational efficiency of adaptive methods. Their emphasis on reducing the time for training and testing models is consistent with the conclusions drawn about the possibility of practical application of adaptive algorithms in real time.

It is worth noting that some aspects of this study revealed discrepancies with some international works. T. Singh *et al.* (2024) primarily addressed the geometric characteristics of vessels (size, shape, orientation) for their identification, while the approach under study demonstrated the effectiveness of adaptive methods for a wider range of objects and surface types. This emphasises the versatility of the developed adaptive

algorithms compared to specialised solutions. In addition, most of the analysed international studies focused on specific types of objects or application scenarios, while this study demonstrates the effectiveness of adaptive approaches for heterogeneous SAR data under variable sensing conditions. This extended the applicability of the results. Compared to current results, further improvement can be achieved by developing more efficient neural network architectures and improving image segmentation strategies. S. Srivastava *et al.* (2021) also noted that approaches combining optical and radar data, as well as algorithms with automatic parameter tuning, appear promising.

The integration of SAR data with optical images, hyperspectral data, and LiDAR data creates a synergistic effect, compensating for the limitations of each type of sensor. The development of deep learning algorithms for multi-sensor data fusion at the feature level can significantly improve the accuracy of object classification and detection. The use of SAR data to provide all-weather monitoring capabilities, supplemented by optical data to improve the interpretability of results, is particularly promising. In contrast to the current implementation, which is focused on local processing, the creation of a web interface for cloud computing can be used to scale the system. In addition, S. Ahmed *et al.* (2023) described further developments in expanding the training set and adapting to new sensors, which create opportunities to improve the generalisation ability of models. The creation of large-scale annotated datasets of SAR images from different regions of the world, different seasons, and different sensing conditions is critical for training robust models. The development of methods for automatic data annotation using weakly supervised learning can significantly accelerate the process of creating training samples. Adapting algorithms to promising SAR sensors, such as future ESA (ROSE-L, Sentinel-1 Next Generation) and NASA (NISAR) missions, as well as commercial high-resolution constellations, will ensure the long-term relevance of the developed technology. The development of methods for automatic adaptation to new types of sensors without the need for retraining models is a key area of research.

The research has made a fundamental contribution to the development of automated SAR data processing methods, demonstrating the possibility of creating technologies that combine scientific breakthroughs with practical applicability. The developed adaptive algorithms create new horizons for the creation of intelligent remote sensing systems capable of autonomously adapting to changing observation conditions and ensuring consistently high-quality data processing.

REFERENCES

- [1] Ahmed, S.F., Alam, M.S., Hassan, M., Rozbu, M.R., Ishtiak, T., Rafa, N., Mofijur, M., Ali, A., & Gandomi, A.H. (2023). Deep learning modelling techniques: Current progress, applications, advantages, and challenges. *Artificial Intelligence Review*, 56, 13521-13617. doi: [10.1007/s10462-023-10466-8](https://doi.org/10.1007/s10462-023-10466-8).

CONCLUSIONS

The experimental study has demonstrated significant advantages of adaptive methods of SAR image processing and classification compared to traditional approaches with fixed parameters. The aim of the study was successfully achieved by developing a comprehensive software system for processing and classifying SAR images with adaptive filtering, segmentation and classification algorithms. The main hypothesis about the significant superiority of adaptive approaches over classical methods was fully confirmed by quantitative performance indicators. The first hypothesis about reducing the level of speckle noise by at least 30% in terms of PSNR was not only confirmed but also surpassed. The developed adaptive Lee filter demonstrated a 16.3% improvement in signal-to-noise ratio compared to the classical Lee filter, which corresponds to an increase in PSNR from 24.3 dB to 28.7 dB (18% improvement). At the same time, up to 92% of spatial structural detail was preserved, with a 16% improvement in EPI from 0.68 to 0.79.

Experimental verification on a test set of 240 SAR images confirmed the effectiveness of the developed algorithms. Adaptive classification methods showed a significant improvement in accuracy: SVM from 81.2 to 86.7%, Random Forest from 83.5 to 88.9%, and U-Net from 87.4 to 92.3%. SAR image segmentation showed a 14% improvement in pixel accuracy for urban areas and an 18% improvement for agricultural fields. The object detection system showed a 25% improvement in the detection-to-false alarm ratio. The practical application of the system has confirmed its effectiveness: for forest monitoring, 91.3% confirmation of detected violations was achieved, for urban mapping, 94.2% accuracy of change detection, and for agricultural monitoring, 89.7% accuracy of crop types. The main limitation of the study is the computational complexity of deep learning algorithms when working with ultra-large images over $32,768 \times 32,768$ pixels. Prospects include the integration of multi-sensor data, the development of more efficient neural network architectures, the creation of a web interface for cloud processing, and the adaptation of algorithms for new types of SAR sensors.

ACKNOWLEDGEMENTS

None.

FUNDING

None.

CONFLICT OF INTEREST

None.

- [2] Alatalo, J. (2023). *Data platform using open-source tools to support geospatial research: Case SAR satellite based change detection*. Retrieved from <https://urn.fi/URN:NBN:fi:amk-2023052212690>.
- [3] Aswani, I., Kar, N.K., Ganguly, T., Ramesh, G.P., & Tejaswini, N. (2023). A fault diagnosis of sound and vibration signals using statistical features and machine learning algorithm. In *International conference on integrated circuits and communication systems* (pp. 1-7). Raichur: IEEE. doi: 10.1109/ICICACSS57338.2023.10100043.
- [4] Declaration of Helsinki. (2024). Retrieved from <https://surl.li/laafnh>.
- [5] European Code of Conduct for Research Integrity. (2020, June). Retrieved from <https://www.ieee.org/about/corporate/governance/p7-8>.
- [6] Fernando, H., Nketia, K., Ha, T., van Steenberg, S., McNairn, H., & Shirtliffe, S. (2025). Automating Sentinel-1 SLC product processing: Parallelization and optimization for efficient polarimetric parameter extraction. *MethodsX*, 14, article number 103253. doi: 10.1016/j.mex.2025.103253.
- [7] Ghafari, R., Kabutarkhani, F.H., & Mansouri, N. (2022). Task scheduling algorithms for energy optimization in cloud environment: A comprehensive review. *Cluster Computing*, 25, 1035-1093. doi: 10.1007/s10586-021-03512-z.
- [8] Gujrati, A., Pradhan, R., Singh, N., Jha, V.B., & Gupta, P.K. (2024). Adaptive water delineation algorithms for L-and C-band SAR imagery: A comparative analysis. *Earth Science Informatics*, 17, 5011-5025. doi: 10.1007/s12145-024-01417-0.
- [9] Huang, M., Zhao, H., & Chen, Y. (2024). Research on SAR image quality evaluation method based on improved Harris Hawk optimization algorithm and XGBoost. *Scientific Reports*, 14, article number 28364. doi: 10.1038/s41598-024-79674-8.
- [10] Jiang, Y., Xie, J., & Zhang, D. (2022). An adaptive offset activation function for CNN image classification tasks. *Electronics*, 11(22), article number 3799. doi: 10.3390/electronics11223799.
- [11] Komenchuk, O. (2024). Adaptive pre-processing methods for increasing the accuracy of segmentation of dental X-RAY images. *Innovative Technologies and Scientific Solutions for Industries*, 3(29), 29-38. doi: 10.30837/2522-9818.2024.3.029.
- [12] Li, B., Ma, Y., Chu, L., Li, W., & Shi, Y. (2024). A GAN-based fast focusing method for circular SAR images. *Heliyon*, 10(14), article number e34133. doi: 10.1016/j.heliyon.2024.e34133.
- [13] Liang, W., Wu, Y., Li, M., & Cao, Y. (2022). Adaptive multiple kernel fusion model using spatial-statistical information for high resolution SAR image classification. *Neurocomputing*, 492, 382-395. doi: 10.1016/j.neucom.2022.03.062.
- [14] Liu, Y., & Lei, Z. (2024). Review of advances in active impulsive noise control with focus on adaptive algorithms. *Applied Sciences*, 14(3), article number 1218. doi: 10.3390/app14031218.
- [15] Lysenko, A. (2023). Synthetic-aperture multi-polarization radar data informativity enhancement technique. *Ukrainian Journal of Remote Sensing*, 10(3), 10-15. doi: 10.36023/ujrs.2023.10.3.243.
- [16] Metrikaityte, G., Suziedelyte Visockiene, J., & Papsys, K. (2022). Digital mapping of land cover changes using the fusion of SAR and MSI satellite data. *Land*, 11(7), article number 1023. doi: 10.3390/land11071023.
- [17] Pasichnik, D., & Onoyko, Yu. (2024). [The most important satellite systems remote sensing of the Earth and its possibilities](#). *Scientific Notes of Young Scientists*, 13.
- [18] Pavlov, Y., & Kashkanov, A. (2023). Analysis of existing methods and approaches to the search of damaged armored tank vehicles during technical intelligence in the modern armies of the world. *Journal of Mechanical Engineering and Transport*, 18(2), 134-140. doi: 10.31649/2413-4503-2023-18-2-134-140.
- [19] Raj, J.A., Idicula, S.M., & Paul, B. (2022). Lightweight SAR ship detection and 16 class classification using novel deep learning algorithm with a hybrid preprocessing technique. *International Journal of Remote Sensing*, 43(15-16), 5820-5847. doi: 10.1080/01431161.2021.2008544.
- [20] Riapolov, I., Vasilets, V., Kukobko, S., & Bodnar, S. (2024). Modelling the surface geometry of unmanned aerial vehicles, the design of which contains elements with different electrophysical properties. *Testing and Certification*, 2(4), 101-110. doi: 10.37701/ts.04.2024.13.
- [21] Selvam, N., Nagesa, Y., & Negesa, F. (2022). Deep learning approach with optimization algorithm for reducing the training and testing time in SAR image detection and recognition. *Indian Journal of Science and Technology*, 15(9), 371-385. doi: 10.17485/IJST/v15i9.1266.
- [22] Shahi, A.P., Rai, P.K., & Mishra, V.N. (2023). Remote sensing data extraction and inversion techniques: A review. In A. Kumar Singh & S. Tiwari (Eds.), *Atmospheric remote sensing: Principles and applications* (pp. 85-104). London: Elsevier. doi: 10.1016/B978-0-323-99262-6.00021-3.
- [23] Shen, S.-L., Zhang, N., Zhou, A., & Yin, Z.-Y. (2022). Enhancement of neural networks with an alternative activation function tanhLU. *Expert Systems with Applications*, 199, article number 117181. doi: 10.1016/j.eswa.2022.117181.
- [24] Singh, T., Babu, T., Nair, R.R., & Duraisamy, P. (2024). Ship detection in synthetic aperture radar imagery: An active contour model approach in computer vision deep learning. *Procedia Computer Science*, 235, 1793-1802. doi: 10.1016/j.procs.2024.04.170.

- [25] Skrypnyk, L., Belenok, V., Velikodsky, Y., Ishchenko, N., & Klymenko, O. (2024). Justification of the advantages of using optical and radar remote sensing data in detecting buildings damaged by natural or anthropogenic impacts. *Ukrainian Journal of Remote Sensing*, 11(4), 13-25. doi: [10.36023/ujrs.2024.11.4.277](https://doi.org/10.36023/ujrs.2024.11.4.277).
- [26] Srivastava, S., Divekar, A.V., Anilkumar, C., Naik, I., Kulkarni, V., & Pattabiraman, V. (2021). Comparative analysis of deep learning image detection algorithms. *Journal of Big Data*, 8, article number 66. doi: [10.1186/s40537-021-00434-w](https://doi.org/10.1186/s40537-021-00434-w).
- [27] Stankevich, S., Svideniuk, M., & Lysenko, A. (2021). Land surface roughness parameter retrieval by inverse simulation of dual-polarization radar backscattering. *Applied Questions of Mathematical Modelling*, 4(2.1), 207-215. doi: [10.32782/KNTU2618-0340/2021.4.2.1.22](https://doi.org/10.32782/KNTU2618-0340/2021.4.2.1.22).
- [28] Trofymchuk, O.M., Hordiienko, O.V., Anpilova, Y.S., & Yakovliev, Y.O. (2024). Monitoring vertical landslides in the Solotvyno agglomeration using Sentinel-1 satellite imagery. *Environmental Safety and Natural Resources*, 50(2), 102-114. doi: [10.32347/2411-4049.2024.2.102-114](https://doi.org/10.32347/2411-4049.2024.2.102-114).
- [29] Vasilenko, D.O., Martynyuk, S.E., Roman, L.O., & Medzmariashvili, E.V. (2025). Features of calibrating satellite RSAs with a multi-beam reflector antenna. *Radioelectronics and Communications Systems*, 67(7), 377-391. doi: [10.20535/S0021347024010047](https://doi.org/10.20535/S0021347024010047).
- [30] Wu, Y., & Li, Q. (2022). The algorithm of watershed color image segmentation based on morphological gradient. *Sensors*, 22(21), article number 8202. doi: [10.3390/s22218202](https://doi.org/10.3390/s22218202).
- [31] Yasir, M., Shanwei, L., Mingming, X., Hui, S., Hossain, M.S., Colak, A.T., Wang, D., Jianhua, W., & Dang, K.B. (2023). Multi-scale ship target detection using SAR images based on improved Yolov5. *Frontiers in Marine Science*, 9, article number 1086140. doi: [10.3389/fmars.2022.1086140](https://doi.org/10.3389/fmars.2022.1086140).
- [32] Zhang, X., & Ding, F. (2021). Optimal adaptive filtering algorithm by using the fractional-order derivative. *IEEE Signal Processing Letters*, 29, 399-403. doi: [10.1109/LSP.2021.3136504](https://doi.org/10.1109/LSP.2021.3136504).

Метод адаптивного шумозаглушення на основі модифікованої Lee-фільтрації для задач класифікації SAR знімків

Юрій Бровка

Аспірант

Національний технічний університет України «Київський політехнічний інститут імені Ігоря Сікорського»
03056, просп. Берестейський, 37, м. Київ, Україна
<https://orcid.org/0009-0004-3708-9747>

Анотація. Метою дослідження було створення комплексу програмних інструментів для автоматизованої обробки та класифікації знімків радарів із синтезованою апертурою, що використовують адаптивні алгоритми аналізу зображень. У дослідженні використано архівні дані з радіолокаційних супутників Sentinel-1, TerraSAR-X і RADARSAT-2 та застосовано як класичні методи обробки зображень, так і адаптивні алгоритми. Якість фільтрації, сегментації, класифікації та виявлення об'єктів оцінювали за показниками точності, структурної подібності, співвідношення сигнал/шум і узгодженості результатів. У ході дослідження було розроблено архітектуру програмного комплексу, що включає модулі попередньої обробки даних радару із синтезованою апертурою, адаптивної фільтрації спеклів, сегментації зображень та класифікації об'єктів. Імплементовано адаптивні алгоритми такі як: фільтр Лі, варіант К-середніх, метод опорних векторів та Ordered Statistics Constant False Alarm Rate. Розроблені інструменти протестовано на супутникових знімках з платформ Sentinel-1 та RADARSAT-2 для різних типів земної поверхні. Адаптивний алгоритм фільтрації покращив якість зображень на 35 %, а продуктивність за ключовими метриками зросла на 15-45 % порівняно з традиційними методами. Забезпечено високу точність класифікації, зокрема за коефіцієнтом Каппа, F1 та площею під кривою робочих характеристик приймача (площа під кривою), при збереженні обчислювальної ефективності. Реалізовано автоматичне розпізнавання водних об'єктів, урбанізованих зон і сільськогосподарських угідь із часом обробки знімка менш як 3 хвилини. Адаптивні алгоритми забезпечували стабільну роботу в умовах різної якості вхідних даних, що робило їх придатними для широкого спектра практичних застосувань у сфері дистанційного зондування та геоінформаційних систем

Ключові слова: дистанційне зондування; метрики; глибоке навчання; адаптивна фільтрація; спекл-шум



Problems of protecting unstructured information on mobile devices

Evgen Brovchenko*

Postgraduate Student

Open International University of Human Development "Ukraine"

04071, 23 Lvivs'ka Str., Kyiv, Ukraine

<https://orcid.org/0000-0002-1416-0385>

Abstract. The relevance of the study was determined by the growing volume of unstructured data in the mobile environment, which required a rethinking of classical approaches to the protection in conditions of limited resources and dynamic use. The purpose of the study was to conduct a comprehensive analysis of methods for protecting unstructured information in the mobile environment and to identify key barriers to the effective implementation. The methodology was based on a theoretical and analytical approach, which included the systematisation of protection methods, a comparative analysis of cryptographic algorithms, an assessment of authentication and access control models, as well as an analysis of cloud security mechanisms. It was established that symmetric encryption Advanced Encryption Standard in Cipher Block Chaining mode provided effective local protection but required careful management of initialisation vectors. It was found that Elliptic Curve Cryptography outperformed Rivest-Shamir-Adleman in terms of energy efficiency and performance, and BLAKE3 outperformed Secure Hash Algorithm 256 in terms of speed, energy consumption, and parallelism support. It was generalised that access control models were insufficiently adapted to the dynamics of the mobile environment, and the most effective were context-oriented and multifactor approaches, in particular with the use of biometric and behavioural authentication. It was found that a combination of client-side encryption, identity management, and cloud backup ensured the highest level of protection with proper implementation. It was established that the effectiveness of HyperText Transfer Protocol Secure, Transport Layer Security 1.3, and Virtual Private Network protocols depended on the type of data and interaction scenario, and the use required a balance between security, performance, and the context of use. Eight key challenges were identified, the relevance of which to mobile security practice was confirmed through comparison with the OWASP Mobile Top 10 categories: data leakage, limited resources, complexity of authentication, dynamic access control, use of public networks, platform fragmentation, application opacity, and legal barriers. It was determined that the effectiveness of protection methods was conditioned by the context of application – the type of data, device architecture, interaction scenario, and available infrastructure – which required an adaptive choice of solutions. The results obtained confirmed that traditional approaches to information security required adaptation to the specifics of mobile platforms. The study has practical value for security solution developers, corporate system administrators, and policymakers in the field of cybersecurity

Keywords: environment; encryption; authentication; resource constraints; access control

INTRODUCTION

In the conditions of rapid digitalisation, mobile devices became the key tool for working with information, in particular unstructured data, which were created, stored, and transmitted in a decentralised environment. Such specifics complicated the use of traditional

security tools, which were developed for stationary or server systems. The limitation of computing resources, the fragmentation of operating systems, and the changing context of mobile device usage created additional challenges for ensuring the confidentiality, integrity,

Article's History: Received: 03.07.2025; Revised: 21.10.2025; Accepted: 15.12.2025; Published: 25.12.2025.

Suggested Citation:

Brovchenko, E. (2025). Problems of protecting unstructured information on mobile devices. *Bulletin of Cherkasy State Technological University*, 30(4), 25-37. doi: 10.62660/bcstu/4.2025.25.

*Corresponding author



and availability of data. There arose a need to critically assess how existing information security mechanisms met the requirements of the mobile environment. It was necessary to clarify whether existing methods were able to adapt to real conditions of use, and whether these methods ensured a sufficient level of protection in conditions of mobility. All this determined the relevance of searching for flexible, adaptive solutions capable of functioning effectively in the limited and dynamic conditions of mobile platforms.

In the context of general system approaches to data security, S. Spasiteleva *et al.* (2019) carried out an analytical review of threats to the security of universal data management platforms, including multimodel database management systems, Data Lake, and cloud environments. The authors highlighted a data-centric security approach and proposed cognitive technologies for automated vulnerability detection. The subject of protecting unstructured information in the mobile environment was directly addressed in a number of studies that combined cryptographic and political-organisational approaches. In the study of Y.M. Brovchenko *et al.* (2023), the focus was on protecting unstructured information on mobile devices, exploring the possibilities of combined application of local encryption, in particular Advanced Encryption Standard (AES), and access policies. The advantages of a hybrid approach were identified, although the need for precise parameter tuning was emphasised. A. Sereda *et al.* (2022) conducted a functional and cryptanalytic analysis of the resilience of algorithms AES, Rivest-Shamir-Adleman (RSA), Elliptic Curve Cryptography (ECC) on Android and iPhone Operating System (iOS) mobile platforms. The study showed that ECC provided the best combination of high cryptographic strength and low energy consumption, which made it the priority choice for mobile security.

A separate strand of the literature focused on protecting communication channels and the features of implementing security mechanisms in mobile device operating systems. The work of Y. Kostiuk *et al.* (2024) focused on the analysis of authentication and key exchange protocols in wireless mobile networks. The authors found that dynamic key updates increased the security of communications, although the implementation of these mechanisms in heterogeneous operating systems (OS) remained technically difficult. In the study of S.M. Konovalov (2025), a categorisation of cyber threats to mobile operating systems was carried out, and the vulnerability to typical attacks was analysed. The effectiveness of Trusted Execution Environment security components in Android and Secure Enclave in iOS was assessed, the limited ability to resist modern attacks was identified, and ways of improving the interaction of security modules with application services were proposed.

Within the related subject area, some authors highlighted aspects of cloud security. In particular, A.R. Abibulaev & A.Z. Piskozub (2025) proposed an innovative approach to strengthening cloud infrastructure security

through the introduction of Natural Language Processing (NLP) and Machine Learning (ML) methods for detecting anomalous activity and predicting risks. The authors confirmed the effectiveness of behavioural analytics for automatic threat response but emphasised the need to adapt models to resource-constrained environments, in particular mobile devices. In a related field, the work of K.G. Babayeva (2024) focused on the study of cryptographic mechanisms for protecting biometric data, with emphasis on the practical application in mobile information systems. The expediency of using specialised cryptographic algorithms for the confidentiality of biometric data was substantiated, but the need for the integration with access policies and channel protection was emphasised.

In the broader context of cybersecurity strategy formation, the contribution of O. Marchenko (2023) was important, in whose study risks in cyberspace were systematised and the effectiveness of existing approaches to the neutralisation was analysed. The author paid attention to multilevel protection models combining organisational, technical, and behavioural security components, which were relevant for protecting mobile environments. In turn, J.V. Rogushina (2019) studied methods of analysing unstructured data, covering algorithmic approaches to processing textual, multimedia, and mixed information. The results of this study formed an important basis for understanding the specifics of unstructured data before developing effective protective mechanisms in the mobile environment. Despite the diversity of scientific approaches, the analysis of sources indicated the fragmentary nature of existing research. Some works focused on encryption, cloud infrastructure, biometrics, or data analysis, but in most cases, there was a lack of an interdisciplinary approach to the problem of protecting unstructured information in the mobile environment. The need to simultaneously take into account the technical limitations of mobile devices, the specifics of unstructured data formats, the dynamics of user behaviour, and the requirements of the regulatory environment remained insufficiently disclosed.

The purpose of the work was to identify technical, organisational, and architectural barriers that complicated the implementation of effective protection of unstructured data on mobile devices through a comprehensive analysis of existing methods and security mechanisms. The hypothesis of the study was that most traditional methods of data protection, developed for stationary or server systems, were not fully effective in the mobile environment without appropriate adaptation, since the methods did not take into account resource constraints, the specifics of unstructured data, and the fragmentation of mobile platforms.

MATERIALS AND METHODS

The study had a theoretical and analytical character. The source base consisted of international technical standards, recommendations, and regulatory documents

in the field of information security. The foundation was made up of publications of the National Institute of Standards and Technology (NIST) (2023), USA, in particular: AES Federal Information Processing Standards (FIPS) 197, recommendations on the use of the Cipher Block Chaining (CBC) mode NIST Special Publication (SP) 800-38A (Dworkin, 2001), and the parameters of elliptic curves NIST SP 800-186 (Chen *et al.*, 2023). For the consideration of digital identification mechanisms and attribute-based access control, the study analysed NIST SP 800-63-3 (Grassi *et al.*, 2017) and NIST SP 800-162 (Hu *et al.*, 2014) respectively. The RSA specification was studied separately according to Request for Comments (RFC) 3447: Public-Key Cryptography Standards (PKCS) #1 (Moriarty *et al.*, 2016). In the part concerning data protection in the cloud environment, ISO/IEC 27017:2015 (2015) – the industry standard for cloud service security – was considered. In addition, to substantiate the requirements for privacy and the protection of personal information, the study took into account the main provisions of Regulation (EU) of the European Parliament and of the Council No. 2016/679 (2016) and the Law of Ukraine No. 2297-VI (2010). These documents formed the basis of the conceptual framework of the study, as well as for the formation of the criterion of compliance of protection methods with the specifics of unstructured data in the mobile environment.

The study applied a number of complementary methods, which ensured separate levels of analytical processing of source information. Content analysis served as the main tool for the systematic collection and study of regulatory and technical documents governing the requirements for data protection in the mobile environment. The classification approach in the study was used to systematise methods of data protection according to functional purpose. The basis for constructing such a structure was the categories of security measures defined in NIST SP 800-53 (Force, 2020). According to this classification, the methods were grouped into categories: cryptographic protection (including symmetric and asymmetric encryption and hash functions), user identification and authentication (including biometric and behavioural methods), access control, data protection during transmission (through secure protocols, in particular HyperText Transfer Protocol Secure (HTTPS), Transport Layer Security (TLS) 1.3, and Virtual Private Network (VPN)), and protection during storage in the cloud environment.

The functional and technical analysis envisaged a comprehensive evaluation of data protection methods, taking into account the suitability for use in the mobile environment. In particular, for asymmetric encryption algorithms RSA and ECC, the comparison was conducted according to performance indicators (encryption/decryption time), energy efficiency, and computational complexity (required key size to achieve a given level of cryptographic strength). Within the analysis of

hash functions Secure Hash Algorithm 256 (SHA-256) and BLAKE3, the following criteria were determined: hashing speed, resource consumption, support for parallel data processing, as well as suitability for implementation on mobile devices with limited resources. For other categories of protection methods, such as authentication models, access control, data transmission protocols, and cloud storage tools, the evaluation was based on such parameters as adaptability to the mobile scenario, scalability, compatibility with existing standards, as well as resilience to typical threats of the mobile environment. This made it possible to carry out a comparative ranking of methods within functional categories and to identify the most suitable solutions for the protection of unstructured data. The comparison method was applied to assess the relevance of the identified technical limitations to practical threats. The comparison procedure involved a step-by-step analysis of the characteristics of the studied protection methods with the categories of threats defined in the classification of the Open Web Application Security Project (OWASP) (Mobile Top 10 2024..., n.d.). To increase the objectivity, the comparison was carried out in tabular form with the correspondence of problems to OWASP categories, which made it possible to trace the correlation between the identified barriers and the practical vulnerabilities of the mobile environment.

RESULTS

Theoretical approaches to protecting unstructured information on mobile devices

The protection of unstructured information on mobile devices required the use of a comprehensive set of methods, encompassing both technical and organisational solutions. The main approaches were based on the application of cryptographic algorithms, access control policies, cloud technologies, and information security standards. The symmetric AES algorithm was widely applied for data encryption at the device level. One of the most common and secure modes of block cipher operation was the CBC mode. In the encryption process, each message block performed an eXclusive "OR" (XOR) operation with the encrypted previous block, while for the first block an XOR operation was performed with the Initialisation Vector (IV). This "chaining" effect in the operation mode of the cipher reflected a high level of security, as it indicated the dependence of each block on the previous one. The AES algorithm in CBC mode was an effective solution for encrypting local files and large volumes of data in mobile applications, due to its simplicity of implementation and built-in support by Android and iOS platforms. A key element of security was the IV, which had to be random and unique for each encryption session, while its transmission could be carried out in the open together with the ciphertext. Table 1 systematised the results of a comparative analysis of methods for generating and transmitting the IV in the context of mobile systems.

Table 1. Approaches to IV generation and transmission for cryptographic protection of unstructured data on mobile devices

IV generation/transmission method	Advantages	Disadvantages	Usage scenarios
Random generation and transmission with data	Simplicity of implementation, suitable for one-time sessions	Risk of “padding oracle” attacks due to open transmission of IV	Secure chats, one-time transactions
Generation from encryption key	Uniqueness of IV ensured automatically when the key changes	Limited uniqueness of IV with a constant key, requires Key Derivation Function (KDF)	Local file encryption on the device
Unique number (Nonce) + counter with synchronisation	Guaranteed uniqueness and efficiency: no need to transmit full IV	Requires reliable synchronisation mechanism between clients	Streaming data (video calls), client-server systems

Source: compiled by the author based on M. Dworkin (2001), National Institute of Standards and Technology (2023)

As shown in Table 1, the choice of the method for generating and transmitting the IV for encrypting unstructured data in mobile applications was determined by a compromise between security, efficiency, and practical implementation: random generation was simple but vulnerable to attacks; a key-derived IV ensured uniqueness but required the implementation of a key derivation function; Nonce + counter was effective for streams but required additional synchronisation. Thus, the optimal choice of IV generation method depended on the specifics of the mobile application, taking into account three key factors: the level of data confidentiality, the presence of server infrastructure, and the computing capacity of the device. This analysis clearly illustrated that even the technical nuances of implementing cryptographic algorithms in the mobile environment required careful consideration of the operational context – from hardware resource constraints to the specifics of network interaction.

Asymmetric encryption algorithms, such as RSA and ECC, played a key role in ensuring security on mobile devices, especially when working with unstructured data. RSA was based on the factorisation of the product of large prime numbers, whereas ECC was based on elliptic curves over finite fields. Both algorithms were applied to provide secure data exchange between the client application and the cloud infrastructure. Encryption/decryption performance and energy efficiency indicators were decisive factors in the choice of data protection tools in a mobile environment with limited

resources. High encryption overheads could slow down device performance, reduce battery life, or worsen the user experience. In Table 2, a comparative characteristic of the efficiency parameters of the considered asymmetric algorithms for different levels of protection and data volumes was presented.

A comparative analysis of asymmetric algorithms demonstrated the advantages of ECC over RSA in the context of mobile devices. ECC ensured an equivalent level of security with significantly smaller key sizes (160 bits versus 1,024 bits for an 80-bit level of protection), which directly affected the reduction of memory usage. The algorithm also turned out to be more energy efficient – at 128-bit protection, ECC consumption was 3.7 times lower than RSA (15.43 MWh and 56.78 MWh respectively). A separate subject of analysis was the dynamics of data processing time. Although RSA demonstrated better results in encrypting small volumes of data, its decryption was slower due to the complexity of operations with large exponents. In contrast, ECC showed more balanced performance, especially at high levels of protection, where decryption time became shorter than encryption time. These results confirmed that ECC was the optimal choice for systems with constant data exchange, while RSA could be appropriate for rare encryption operations. The choice of a specific algorithm had to take into account not only cryptographic strength but also energy efficiency and the computing capabilities of mobile devices.

Table 2. Efficiency of asymmetric encryption algorithms

Security level (bits)	Key size (bits)		Data size (bits)	Total time (ms)		Encryption time (ms)		Decryption time (ms)		Energy consumption (MWh)	
	RSA	ECC		RSA	ECC	RSA	ECC	RSA	ECC	RSA	ECC
80	1,024	160	8	785	1,815.2	30.7	488.5	754.3	1,326.7	17.86	9.05
			64	5,673.8	8,078.4	136.6	2,168.5	5,537.2	5,909.9		
			256	19,877.2	30,809.1	559.6	7,924	19,317.7	22,885.1		
112	2,048	224	8	2,737.5	3,789.3	29.9	2,203	2,707.5	1,586.3	21.55	17.38
			64	20,574.3	16,978.8	163.5	9,985.5	20,410.8	6,933.3		
			256	102,615.3	66,033.9	581.5	39,700.8	102,033.7	26,333.1		

Continued Table 2.

Security level (bits)	Key size (bits)		Data size (bits)	Total time (ms)		Encryption time (ms)		Decryption time (ms)		Energy consumption (MWh)	
	RSA	ECC		RSA	ECC	RSA	ECC	RSA	ECC	RSA	ECC
128	3,072	256	8	6,971.4	5,645.3	30.5	3,876.3	6,940.9	1,769	56.78	15.43
			64	46,645.4	22,446.6	167.2	15,088.2	46,478.2	7,358.4		
			256	210,169.7	85,844.6	561.1	58,438.6	209,608.6	27,406		

Source: compiled by the author based on K. Moriarty *et al.* (2016), Z. Vahdati *et al.* (2019), L. Chen *et al.* (2023)

In mobile security systems, hash functions were used to verify the integrity of unstructured data such as multimedia files, text documents, event logs, or cached content. Hashing made it possible to quickly detect unauthorised changes in content by comparing the current hash with the control value. In mobile applications, hash functions were used to verify the integrity of files during storage or synchronisation (for example, during offline access to images or videos), to control cache stability, as well as to validate access tokens or the authenticity of update packages. Some applications

implemented additional checksum verification in the background to avoid reproducing corrupted or modified data. The most commonly used algorithms were SHA-256 and BLAKE3. SHA-256 was a standardised solution compatible with traditional cryptographic protocols, whereas BLAKE3 was designed with an emphasis on hardware capabilities. Both algorithms were used mainly to verify data integrity and did not provide encryption functions. A detailed comparative analysis of the characteristics of hash algorithms in the context of mobile devices was presented in Table 3.

Table 3. Comparison of SHA-256 and BLAKE3 algorithms in the context of mobile systems

Parameter	SHA-256	BLAKE3
Processing speed (Advanced RISC Machine (ARM), 1 thread)	300-400 MB/s	1,000 MB/s
Energy consumption (indicative assessment)	Higher: longer computation time → more Central Processing Unit (CPU) cycles	Lower: shorter computation time, efficient use of cache and Single Instruction Multiple Data (SIMD)
Multithreading support	Limited	Full (parallelised across cores)
Adaptation to the mobile environment	Needs optimisation	Designed with weaker CPUs in mind
Application scenarios	Standardised protocols (TLS, Hash-based Message Authentication Code (HMAC)) and legacy systems	Local data verification, background scanning, cache synchronisation

Source: compiled by the author based on B. Kibar (2023)

According to the comparative characteristics, BLAKE3 demonstrated higher performance when processing large volumes of unstructured data, reducing file integrity verification time and energy consumption. The algorithm supported parallel processing, which corresponded to the architecture of the latest generation of multicore mobile processors. SHA-256 remained relevant for scenarios where compliance with cryptographic standards or compatibility with existing protocols was required. Both algorithms could be used in combined data protection systems. Access control was one of the key mechanisms for ensuring the confidentiality and integrity of unstructured information on

mobile devices. Effective implementation of access ensured that only authorised users or applications gained access to system data, resources, or functions. In the mobile environment, these mechanisms had to adapt to limited computing resources, the high dynamics of network connections, and frequently changing usage contexts. Despite the diversity of access control models, the implementation in mobile operating systems and applications was accompanied by a number of limitations and vulnerabilities, which reduced the overall level of protection of unstructured data. Table 4 presented a classification of access control models with an analysis of the implementation in the mobile environment.

Table 4. Comparative characteristics of access control models in the mobile environment

Access model	Principle of operation	Implementation in mobile systems	Typical application examples
Discretionary Access Control (DAC)	The user independently determined who had access to the data	Configuring data sharing in applications, system permissions of the file system (Android Storage Access Framework)	Granting access to files for other applications

Access model	Principle of operation	Implementation in mobile systems	Typical application examples
Mandatory Access Control (MAC)	Access policies were set centrally	SELinux mechanism on Android	Android kernel security
Role-Based Access Control (RBAC)	Access was determined by the user's role	In Mobile Device Management (MDM) systems, corporate platforms	Defining employee access rights to data according to the position
Attribute-Based Access Control (ABAC)	Access depended on user, environment, and resource attributes (time, location, device type)	Application Programming Interface (API) of the new generation, mobile Software Development Kit (SDK)	Applications with geo-zones, Internet of Things (IoT) systems with time-restricted access
Multi-Factor Authentication (MFA)	Access was allowed after passing several levels of authentication	Use of biometrics, passwords, and one-time codes	Banking applications, cloud storage

Source: compiled by the author based on V.C. Hu *et al.* (2014), P.A. Grassi *et al.* (2017), J.T. Force (2020)

The analysis of the main access models demonstrated that traditional approaches had limitations in the mobile context. DAC, although intuitive for users, often became a source of misconfigurations due to excessive dependence on the human factor. MAC ensured a higher level of protection, but its complexity limited its application mainly to system components. RBAC proved effective in corporate environments, but its lack of flexibility complicated adaptation to dynamic mobile scenarios. ABAC and MFA offered more promising solutions, as these approaches took usage context into account. However, the implementation was accompanied by additional resource costs and raised issues concerning privacy protection. These observations indicated the necessity of

developing hybrid solutions that combined the advantages of different models, taking into account the technical features of mobile platforms and patterns of user behaviour. The optimal system had to ensure an adequate level of security while maintaining performance and usability. For the effective protection of unstructured data in the cloud environment, a comprehensive approach was necessary, combining three key methods: client-side encryption, Identity and Access Management (IAM), and backup mechanisms. These technologies jointly ensured confidentiality, integrity, and availability of data when working with mobile devices. Table 5 presented the results of a comprehensive analysis of methods for protecting unstructured data in the cloud environment.

Table 5. Methods for protecting unstructured data in the cloud environment

Protection method	Main functions	Objectives	Key advantages	Client-cloud interaction
Client-side encryption	Data encryption before uploading to the cloud, key management on the device	Preventing data leakage in case of cloud storage compromise	Full control over keys, absence of cloud provider access to data content	Client encrypts data → transfers the data to the cloud → decrypts only when downloading to the device
IAM	Centralised access management, authentication, authorisation	Control of user rights, prevention of unauthorised access	Flexible access policies, integration with MFA, activity audit	Cloud verifies access rights → grants/blocks data operations for the client
Backup	Automatic creation of data copies, the encryption, recovery in case of loss	Ensuring availability, protection against data loss	Reliability, ability to roll back to previous file versions	Client configures backup schedule → cloud stores versions → client initiates recovery

Source: compiled by the author based on Law of Ukraine No. 2297-VI (2010), ISO/IEC 27017:2015 (2015), Regulation (EU) of the European Parliament and of the Council No. 2016/679 (2016)

The results presented in Table 5 of the analysis of cloud protection methods testified to the need for a differentiated approach to processing unstructured data in the mobile environment. Client-side encryption, despite increased demands on computing resources, remained the optimal choice for confidential information, while IAM systems realised the potential best in corporate solutions with clearly defined user roles. Backup mechanisms, being critically important for data integrity, required individual configuration for each type of

content, taking into account its value and update frequency. These findings emphasised that effective protection of unstructured data in the cloud required not only the technical implementation of separate mechanisms but also a deep understanding of the context of the use. The optimal strategy envisaged a combination of protection methods, adaptation to the type of data and usage scenarios, realistic assessment of the capabilities of mobile devices, and a balance between security and performance.

Data transmission protocols of the latest versions played a crucial role in protecting unstructured data when working with mobile devices. A direct correlation was established between the type of transmitted data and the choice of an appropriate protocol: for API requests and processing of confidential information, it was advisable to use HTTPS, which provided basic encryption; transmission of large volumes of data (media content, backups) was more efficiently implemented through TLS 1.3, which combined high performance and cryptographic protection; in corporate scenarios, the optimal solution was the implementation of VPN for secure connections via open networks. Such a differentiated approach allowed taking into account the specifics of mobile applications and the type of data being processed. The analysis results confirmed that none of the reviewed protection methods were universal in the context of mobile devices. Each method had technological and contextual limitations associated with resource constraints, peculiarities of the operating environment, and modelling of user behaviour. Only the combination of complementary approaches allowed the formation of a resilient protection system for

unstructured data. Accordingly, the hypothesis about the limited effectiveness of traditional security tools developed for stationary or server platforms in the mobile context was confirmed.

Challenges and limitations of protecting unstructured data in the mobile environment

Despite the active development of information protection methods, the effective implementation of these technologies in the dynamic mobile environment, particularly regarding unstructured data, faced a number of challenges. These challenges arose both at the level of technical implementation and due to the peculiarities of the mobile environment, user behavioural factors, and legal-regulatory restrictions. Within this study, the following key challenges were identified and systematised, which complicated the implementation of effective protection of unstructured data in the mobile environment (Table 6). To ensure the objectivity and relevance of the analysis, the identified problems were compared with the most widespread and critical risks defined in the industry standard OWASP Mobile Top 10 (Mobile Top 10 2024..., n.d.).

Table 6. Key problems of protecting unstructured information on mobile devices

No.	Challenge	The essence of the problem	Correspondence to OWASP Mobile Top 10 for 2024
1	High vulnerability to unstructured data leakage	The absence of centralised accounting, classification, and access control complicated protection of data from unauthorised access or loss	M5, M6, M9
2	Limited computing resources of mobile devices	Technical limitations of mobile devices restrained full implementation of cryptographic protection without harming performance	M10
3	Complexity of ensuring continuous authentication	Traditional authentication mechanisms did not provide sufficient resilience in mobile use and were subject to compromises between convenience and security	M1, M3
4	Access control in a dynamic environment	Variability of usage context complicated the application of static access policies, requiring adaptive management models	M3, M8
5	Use of public networks and cloud services	Unsafe transmission and storage of confidential data without end-to-end encryption increased the risk of interception or modification	M2, M5, M6, M9
6	Fragmentation of platforms and ecosystems	Uneven updates and differences between platforms hindered the creation of unified protection tools	M7, M8
7	Lack of transparency in mobile applications	Lack of openness in mobile applications created preconditions for hidden data collection, processing, or leakage	M4, M6
8	Legal and ethical aspects of protecting unstructured information	Inconsistency of international norms complicated compliance with data protection requirements during the processing in different jurisdictions	M6, M9

Source: compiled by the author based on Mobile Top 10 2024: Final release updates (n.d.)

For a deeper understanding of the identified barriers, it was advisable to group the barriers into three key directions: technical limitations, user factors, and regulatory challenges. Among the key technical barriers to effective protection of unstructured information in the mobile environment, several circumstances were highlighted, related to the architectural and hardware features of devices. Firstly, the decentralised nature of such data hindered centralised control, classification,

and implementation of unified access policies. Secondly, limited resources of mobile devices – processor, memory, autonomy – restrained the application of full cryptographic protection without harming performance. Finally, the high fragmentation of mobile ecosystems (different OS versions, marketplace policies, hardware features) complicated the creation of unified security solutions that had to function equally across different platforms.

A separate group of challenges consisted of user-related aspects. The most critical aspect was the complexity of ensuring continuous authentication in the context of mobile use – traditional identification tools had limitations in the mobile context. Another widespread problem was incorrect configuration of access or its complete absence, which opened opportunities for abuse. An additional barrier was the low level of transparency in the operation of applications – in particular, lack of open code, insufficient documentation, and non-transparent data collection practices. At the regulatory level, the key issues remained inconsistency between national legislations in the field of data protection and the complexity of bringing mobile applications into compliance with international standards. Ensuring regulatory compliance was particularly difficult in the case of processing unstructured information, which by its nature was rarely subject to clear classification and accounting. This created risks both for data operators and for end-users. The identified challenges regarding the protection of unstructured information on mobile devices demonstrated the multidimensional nature of the problem – technical, organisational, legal, and behavioural. The challenges encompassed both internal platform aspects (computing limitations, OS fragmentation, cryptographic effectiveness) and external factors – from unpredictable environmental conditions to regulatory barriers. A separate difficulty was the dynamic context of mobile use, which required adaptive security policies. A comparative analysis with the categories of OWASP Mobile Top 10 2024: Final release updates (n.d.) confirmed the relevance of the identified challenges to mobile security practice and testified to the critical necessity of the consideration in the design of protection systems. Thus, effective resolution of these issues required an interdisciplinary approach, combining technical solutions, regulatory frameworks, and transparent interaction between the user and mobile applications.

DISCUSSION

The obtained results demonstrated that although methods of protecting unstructured information, such as symmetric and asymmetric encryption, hashing, multifactor authentication, access control, and cloud technologies, had high potential effectiveness, the implementation in the mobile environment was significantly complicated by a number of technical, organisational, and architectural barriers. The conducted analysis showed that these methods were often not adapted to the real conditions of mobile device functioning – limited energy capacity, OS fragmentation, open interfaces, and variable user behaviour. This emphasised the limitations of universal approaches developed for the server environment, while at the same time confirming the research hypothesis regarding the necessity of critically rethinking protection in the mobile context. The importance of these results lay in formulating a holistic understanding of the spectrum of challenges facing mobile

security and in the analytical justification of the expediency of a comprehensive, adaptive approach to implementing systems for protecting unstructured information.

Within the conducted study, it was established that the effectiveness of cryptographic protection of unstructured information on mobile devices depended on the choice of algorithm, taking into account resource limitations. In particular, the use of AES in CBC mode required correct IV management, which was complicated in mobile systems with unstable execution contexts. Similar technical challenges were recorded in the work of K. Yu *et al.* (2022), which analysed Shamir's cryptography in distributed IoT environments. However, in that study the emphasis was shifted towards model resilience, while the impact of computational load was left unaddressed. In contrast, Y. Chen *et al.* (2020) directly pointed to the problems of energy consumption and the complexity of implementing cryptographic algorithms in mobile devices, which confirmed the findings. A comparative analysis of asymmetric encryption algorithms showed that the ECC algorithm was more suitable for the mobile environment due to its smaller key size, lower energy consumption, and faster data processing compared to RSA. This was important for protecting unstructured information, the volume of which was variable and often significant. Similar conclusions were presented in the work of R. Yuvarani & R. Mahaveerakannan (2025), where ECC showed better performance compared to RSA in the context of cloud authentication, although the authors did not focus directly on mobile devices as the local execution environment. K. Liu *et al.* (2023) confirmed the advantages of ECC in the mobile context, emphasising its effectiveness in two-factor authentication, which was consistent with the findings obtained regarding the expediency of applying ECC in scenarios of accessing encrypted data. Relevant were the results of K.S. Kumar & R. Sukumar (2019), which proved the advantage of ECC in terms of energy efficiency for Android devices, fully coinciding with the conclusions of this study about the expediency of using elliptic curves for protecting mobile data. Thus, the results of the study were confirmed by previous works, while focusing attention on the application of ECC specifically for protecting unstructured information in resource-constrained conditions.

A comprehensive review of access control models in the mobile environment demonstrated that the implementation of effective authentication and access management in the mobile environment was complicated by platform fragmentation, variability of user interaction models, and limitations of computing resources. The most effective approaches proved to be multifactor and context-dependent methods, which could ensure a balance between convenience and the level of protection. Similar conclusions were observed in the study of A. Tewari & B.B. Gupta (2020), which noted the necessity of comprehensive construction of access systems in multi-level IoT environments, with an emphasis on

adaptability and dynamic rights management. At the same time, in the mentioned work the main focus was on the general architecture of IoT, while the specifics of mobile OS and user scenarios were left beyond the analysis. In the work of H.F. Atlam & G.B. Wills (2020), the importance of considering ethical and social aspects when implementing identification mechanisms was outlined, which was relevant for the mobile environment with its high level of personalisation, although the technical parameters of access implementation were not disclosed. In contrast, in the study of J. Du *et al.* (2018) the expediency of distributed authentication models was justified to reduce the load on central nodes, which partially correlated with the mobile need for delegated processing.

The analysis of the cloud component in the context of protecting unstructured data on mobile devices confirmed the expediency of implementing client-side encryption before transferring data to the cloud, given the risks of losing control over confidential information after its upload. Such an approach was confirmed in the study of A. Musa & A. Mahmood (2021), which noted that client-side encryption significantly increased trust in cloud storage, although the authors did not consider the limitations of mobile platforms. M. da Rocha *et al.* (2020) expanded this vision, proposing the use of a trusted execution environment, which allowed keys to be stored on the client side – a similar concept corresponded with the emphasis on minimising dependence on the cloud provider. Regarding access control, the results of the study highlighted the need for dynamic rights management on the user side, which was consistent with the position of J. Mellom (2020), who emphasised the role of IAM as a basic component of cloud security. The review proposed by A.O. Akinade *et al.* (2025) confirmed the importance of combining encryption, IAM, and authentication mechanisms in cloud environments, but mainly from the standpoint of infrastructure providers, whereas in the present study the focus was on the possibilities and limitations of the mobile client.

One of the central challenges identified in the study was the high vulnerability to leakage of unstructured data in the mobile environment, caused by the absence of a centralised model of accounting, classification, and access control. This risk was aggravated by the fact that mobile applications often gained access to such data without clear restrictions and proper informed consent of the user. Similar observations were described in the work of Y. Guo *et al.* (2021), which emphasised that unstructured data was easily extracted, particularly due to the lack of transparent information processing policies in mobile applications. At the same time, in the mentioned study the main emphasis was on the technical aspects of data extraction, while in the present study the emphasis was shifted to the systemic absence of mechanisms to prevent leakage. An additional barrier was identified in the fragmentation of mobile platforms, which complicated the implementation of unified solutions for information protection.

As noted in the study of S. Garg & N. Baliyan (2021), the differences between Android and iOS – in security approaches, update frequency, and access configuration – created unequal conditions for the implementation of protection standards. Thus, the obtained results were confirmed: effective protection of unstructured information on mobile devices had to take into account not only general technical mechanisms, but also the architectural specifics of mobile ecosystems.

Among the identified challenges of protecting unstructured information on mobile devices, a special place was occupied by the problem of continuous user authentication. As shown in the results of the study, traditional authentication methods, such as passwords, Personal Identification Number codes, or one-time tokens, did not provide an adequate level of resilience under conditions of mobile use, characterised by environmental variability, short sessions, and limited user attention. At the same time, alternative approaches, such as biometric parameters or behavioural patterns, remained vulnerable to attacks, recognition errors, and breaches of confidentiality. Similar limitations were reflected in the study of M. Papaioannou *et al.* (2023), which noted that although behavioural biometrics showed potential for continuous identification, it significantly depended on the context of use and the device. The work of Y. Yang *et al.* (2019) confirmed this assessment, emphasising the difficulties in maintaining stable accuracy of behavioural models in the dynamic conditions of the mobile environment. Similar conclusions were presented in the study of P.M.A.B. Estrela *et al.* (2021), which showed that the implementation of continuous authentication in mobile banking required complex infrastructure and a balance between convenience and security.

Insufficient protection of data during transmission via public networks or storage in cloud environments was identified as one of the challenges that increased the risks of unauthorised access to unstructured information. The study emphasised the lack of end-to-end encryption on the mobile device side, which made data vulnerable to Man-in-the-Middle attacks and interception. A similar problem was identified by Y. Yao *et al.* (2023), who proposed a complex secure transmission scheme, which, however, required high resources and did not correspond to the limitations of mobile devices. In addition, the observations of M. Sangeen *et al.* (2023) about the low awareness of users regarding the dangers of public networks resonated with the results of the present study, where behavioural aspects were identified as an additional risk factor. Unlike the mentioned authors, the present study detailed the problem precisely in the context of unstructured data and the mobile environment, emphasising the insufficiency of basic protection models.

Legal and regulatory restrictions were considered a significant barrier to ensuring the confidentiality of unstructured data. The present work indicated that the complexity of auditing such information and

the fragmentation of legal requirements complicated compliance with privacy standards. This conclusion was consistent with the arguments of W. Hartzog & N.M. Richards (2020), who pointed to the institutional lack of adaptation of law to digital dynamics. Similarly, J. Wong & T. Henderson (2019) emphasised the technical difficulties of implementing the right to portability under GDPR. However, while in the mentioned works attention was focused on general regulatory approaches, the present study supplemented these approaches by highlighting how the unstructured nature of mobile data made full legal regulation impossible. This comparison confirmed that legal challenges were no less critical than technical ones and required integrated solutions. The analysis of the results of the study reflected that the identified challenges – in particular the problems of unstructured data leakage, limited mobile resources, authentication difficulties, and regulatory barriers – were systemic and largely confirmed by other studies. At the same time, the conducted comparison showed that the present study expanded existing approaches by comprehensively covering the technical, behavioural, and legal aspects specific to the mobile environment. Such a comprehensive analysis allowed not only the identification of current problems, but also the critical evaluation of the applicability limits of existing solutions in the context of protecting unstructured information.

CONCLUSIONS

Within this study, a comprehensive analysis was carried out of methods of protecting unstructured information on mobile devices, taking into account technical, organisational, and legal-regulatory factors. The results of the theoretical analysis of cryptographic algorithms showed that symmetric AES encryption in CBC mode was suitable for protecting local data, but required careful IV management. A comparative assessment of RSA and ECC confirmed the advantage of the latter in terms of energy efficiency and speed, which made it more optimal for mobile application. The analysis of hash functions confirmed the higher effectiveness of BLAKE3 compared to SHA-256 due to its lower energy consumption, higher speed, and support for parallel data processing, which were critical for mobile platforms.

Traditional authentication mechanisms proved insufficiently adapted to the mobile environment, while

promising approaches were those combining behavioural biometrics with a context-dependent approach. The analysis of cloud protection methods demonstrated the effectiveness of a combination of client-side encryption, IAM systems, and backup, particularly in ensuring the confidentiality of unstructured information during its transmission and storage. The analysis of modern secure transmission protocols, in particular HTTPS, TLS 1.3, and VPN, confirmed the high effectiveness in ensuring data confidentiality during transportation in open networks. It was established that the effectiveness of protection methods directly depended on the type of processed data, the architecture of the mobile device, resource limitations, and environmental conditions, which necessitated an adaptive choice of protection solutions.

Special attention was paid to the identification of key challenges complicating the implementation of unstructured data protection on mobile devices. Eight problems were systematised, covering technical, behavioural, and regulatory aspects: from the risk of data leakage and mobile resource limitations to application transparency issues and legal barriers. The comparison with OWASP Mobile Top 10 made it possible to verify the relevance of these challenges to current security practices. The results of the study confirmed the hypothesis that traditional protection tools, oriented towards stationary or server environments, proved insufficiently effective in the mobile context, requiring a comprehensive and adaptive approach. The study had limitations related to the theoretical nature of the analysis, the absence of empirical testing, and the inaccessibility of full technical specifications of mobile OS and commercial applications. In further research, it would be expedient to focus on the practical testing of the effectiveness of cryptographic algorithms optimised for the limitations of mobile devices, as well as on the development of adaptive authentication systems and the creation of cross-platform solutions for data protection.

ACKNOWLEDGEMENTS

None.

FUNDING

None.

CONFLICT OF INTEREST

None.

REFERENCES

- [1] Abibulaev, A.R., & Piskozub, A.Z. (2025). Analysis of possibilities for improving cloud infrastructure security using NLP and ML. *Modern Information Security*, 2(62), 124-140. doi: [10.31673/2409-7292.2025.026884](https://doi.org/10.31673/2409-7292.2025.026884).
- [2] Akinade, A.O., Adepoju, P.A., Ige, A.B., & Afolabi, A.I. (2025). Cloud security challenges and solutions: A review of current best practices. *International Journal of Multidisciplinary Research and Growth Evaluation*, 6(1), 26-35. doi: [10.54660/IJMRGE.2025.6.1.26-35](https://doi.org/10.54660/IJMRGE.2025.6.1.26-35).
- [3] Atlam, H.F., & Wills, G.B. (2020). IoT security, privacy, safety and ethics. In M. Farsi, A. Daneshkhah, A. Hosseinian-Far & H. Jahankhani (Eds.), *Digital twin technologies and smart cities* (pp. 123-149). Cham: Springer. doi: [10.1007/978-3-030-18732-3_8](https://doi.org/10.1007/978-3-030-18732-3_8).

- [4] Babayeva, K.G. (2024). Using cryptographic methods, mechanisms, and tools for protecting biometric data. In *Radioelectronics and youth in the 21st century: Materials of the 28th international youth forum* (pp. 88-90). Kharkiv: Kharkiv National University of Radio Electronics. doi: [10.30837/IYF.PCEIP.2024.088](https://doi.org/10.30837/IYF.PCEIP.2024.088).
- [5] Brovchenko, Y.M., Samarai, V.P., Datsenko, I.P., Pavlenko, V.I., & Sereda, A.V. (2023). Protection of unstructured data on mobile devices. *Infocommunications and Computer Technologies*, 1(5), 194-200. doi: [10.36994/2788-5518-2023-01-05-21](https://doi.org/10.36994/2788-5518-2023-01-05-21).
- [6] Chen, L., Moody, D., Randall, K., Regenscheid, A., & Robinson, A. (2023). *Recommendations for discrete logarithm-based cryptography: Elliptic curve domain parameters*. Gaithersburg: National Institute of Standards and Technology. doi: [10.6028/NIST.SP.800-186](https://doi.org/10.6028/NIST.SP.800-186).
- [7] Chen, Y., Zheng, B., Zhang, Z., Wang, Q., Shen, C., & Zhang, Q. (2020). Deep learning on mobile and embedded devices: State-of-the-art, challenges, and future directions. *ACM Computing Surveys*, 53(4), article number 84. doi: [10.1145/3398209](https://doi.org/10.1145/3398209).
- [8] da Rocha, M., Valadares, D.C.G., Perkusich, A., Gorgonio, K.C., Pagno, R.T., & Will, N.C. (2020). Secure cloud storage with client-side encryption using a trusted execution environment. In *Proceedings of the 10th international conference on cloud computing and services science* (pp. 31-43). Setubal: SciTePress. doi: [10.5220/0009130600310043](https://doi.org/10.5220/0009130600310043).
- [9] Du, J., Jiang, C., Gelenbe, E., Xu, L., Li, J., & Ren, Y. (2018). Distributed data privacy preservation in IoT applications. *IEEE Wireless Communications*, 25(6), 68-76. doi: [10.1109/MWC.2017.1800094](https://doi.org/10.1109/MWC.2017.1800094).
- [10] Dworkin, M. (2001). *Recommendation for block cipher modes of operation: Methods and Techniques*. Gaithersburg: National Institute of Standards and Technology. doi: [10.6028/NIST.SP.800-38A](https://doi.org/10.6028/NIST.SP.800-38A).
- [11] Estrela, P.M.A.B., Albuquerque, R.D.O., Amaral, D.M., Giozza, W.F., & Júnior, R.T.D.S. (2021). A framework for continuous authentication based on touch dynamics biometrics for mobile banking applications. *Sensors*, 21(12), article number 4212. doi: [10.3390/s21124212](https://doi.org/10.3390/s21124212).
- [12] Force, J.T. (2020). *Security and privacy controls for information systems and organizations*. Gaithersburg: National Institute of Standards and Technology. doi: [10.6028/NIST.SP.800-53r5](https://doi.org/10.6028/NIST.SP.800-53r5).
- [13] Garg, S., & Baliyan, N. (2021). Comparative analysis of Android and iOS from security viewpoint. *Computer Science Review*, 40, article number 100372. doi: [10.1016/j.cosrev.2021.100372](https://doi.org/10.1016/j.cosrev.2021.100372).
- [14] Grassi, P.A., Garcia, M.E., & Fenton, J.L. (2017). *Digital identity guidelines*. Gaithersburg: National Institute of Standards and Technology. doi: [10.6028/NIST.SP.800-63-4](https://doi.org/10.6028/NIST.SP.800-63-4).
- [15] Guo, Y., Liu, J., Tang, W., & Huang, C. (2021). Exsense: Extract sensitive information from unstructured data. *Computers & Security*, 102, article number 102156. doi: [10.1016/j.cose.2020.102156](https://doi.org/10.1016/j.cose.2020.102156).
- [16] Hartzog, W., & Richards, N.M. (2020). Privacy's constitutional moment and the limits of data protection. *SSRN*, article number 3441502. doi: [10.2139/ssrn.3441502](https://doi.org/10.2139/ssrn.3441502).
- [17] Hu, V.C., Ferraiolo, D., Kuhn, R., Schnitzer, A., Sandlin, K., Miller, R., & Scarfone, K. (2014). *Guide to attribute based access control (ABAC) definition and considerations*. Gaithersburg: National Institute of Standards and Technology. doi: [10.6028/NIST.SP.800-162](https://doi.org/10.6028/NIST.SP.800-162).
- [18] ISO/IEC 27017:2015. (2015). *Information technology – security techniques – code of practice for information security controls based on ISO/IEC 27002 for cloud services*. Retrieved from <https://www.iso.org/standard/43757.html>.
- [19] Kibar, B. (2023). *Comparing Blake3 and Sha-256 data integrity algorithms & integrating Blake3 with Golang*. Retrieved from <https://surl.lu/vdnytl>.
- [20] Konovalov, S.M. (2025). Analysis of types of cybersecurity in mobile phone operating systems. *Taurida Scientific Herald. Series: Technical Sciences*, 2, 100-104. doi: [10.32782/tnv-tech.2025.2.11](https://doi.org/10.32782/tnv-tech.2025.2.11).
- [21] Kostiuk, Y., Bebeshko, B., Kriuchkova, L., Lytvynov, V., Oksanych, I., Skladannyi, P., & Khorolska, K. (2024). Information protection and data exchange security in wireless mobile networks with authentication and key exchange protocols. *Cybersecurity: Education, Science, Technique*, 1(25), 229-252. doi: [10.28925/2663-4023.2024.25.229252](https://doi.org/10.28925/2663-4023.2024.25.229252).
- [22] Kumar, K.S., & Sukumar, R. (2019). Achieving energy efficiency using novel scalar multiplication based ECC for Android devices in Internet of Things environments. *Cluster Computing*, 22(5), 12021-12028. doi: [10.1007/s10586-017-1542-8](https://doi.org/10.1007/s10586-017-1542-8).
- [23] Law of Ukraine No. 2297-VI "On Personal Data Protection". (2010, June). Retrieved from <https://zakon.rada.gov.ua/laws/show/en/2297-17#Text>.
- [24] Liu, K., Zhou, Z., Cao, Q., Xu, G., Wang, C., Gao, Y., Zeng, W., & Xu, G. (2023). A robust and effective two-factor authentication (2FA) protocol based on ECC for mobile computing. *Applied Sciences*, 13(7), article number 4425. doi: [10.3390/app13074425](https://doi.org/10.3390/app13074425).
- [25] Marchenko, O. (2023). Cybersecurity and information protection: Analysis of risk and threat impact with modern effective cyberspace defense strategies. *Information Technology: Computer Science, Software Engineering and Cyber Security*, 3, 50-59. doi: [10.32782/IT/2023-3-6](https://doi.org/10.32782/IT/2023-3-6).

- [26] Mellom, J. (2020). *The role of identity access management (IAM) in cloud security*. Retrieved from <https://sonraisecurity.com/blog/the-role-of-identity-access-management-iam-in-governing-cloud-security/>.
- [27] Mobile Top 10 2024: Final release updates. (n.d.). Retrieved from <https://owasp.org/www-project-mobile-top-10/>.
- [28] Moriarty, K., Kaliski, B., Jonsson, J., & Rushc, A. (2016). *PKCS #1: RSA cryptography specifications version 2.2*. Retrieved from <https://datatracker.ietf.org/doc/html/rfc8017>.
- [29] Musa, A., & Mahmood, A. (2021). Client-side cryptography based security for cloud computing system. In *2021 international conference on artificial intelligence and smart systems (ICAIS)* (pp. 594-600). Coimbatore: IEEE. doi: [10.1109/ICAIS50930.2021.9395890](https://doi.org/10.1109/ICAIS50930.2021.9395890).
- [30] National Institute of Standards and Technology. (2023). *Advanced Encryption Standard (AES)*. Gaithersburg: National Institute of Standards and Technology. doi: [10.6028/NIST.FIPS.197-upd1](https://doi.org/10.6028/NIST.FIPS.197-upd1).
- [31] Papaioannou, M., Mantas, G., Panaousis, E.M., Essop, A., Rodriguez, J., & Sucasas, V. (2023). Behavioral biometrics for mobile user authentication: Benefits and limitations. In *2023 IFIP networking conference (IFIP networking)* (pp. 1-6). Barcelona: IEEE. doi: [10.23919/IFIPNetworking57963.2023.10186419](https://doi.org/10.23919/IFIPNetworking57963.2023.10186419).
- [32] Regulation (EU) of the European Parliament and of the Council No. 2016/679 “On the Protection of Natural Persons with Regard to the Processing of Personal Data and on the Free Movement of Such Data, and Repealing Directive 95/46/EC (General Data Protection Regulation) (Text with EEA Relevance)”. (2016, April). Retrieved from <http://data.europa.eu/eli/reg/2016/679/oj>.
- [33] Rogushina, J.V. (2019). Methods and tools for analyzing unstructured data. *Problems in Programming*, 1, 57-77. doi: [10.15407/pp2019.01.057](https://doi.org/10.15407/pp2019.01.057).
- [34] Sangeen, M., Bhatti, N.A., Kifayat, K., Alsadhan, A.A., & Wang, H. (2023). Blind-trust: Raising awareness of the dangers of using unsecured public Wi-Fi networks. *Computer Communications*, 209, 359-367. doi: [10.1016/j.comcom.2023.07.011](https://doi.org/10.1016/j.comcom.2023.07.011).
- [35] Sereda, A., Datsenko, I., Pavlenko, V., & Samarai, V. (2022). Stability and efficiency of cryptographic algorithms used in mobile devices. *Infocommunications and Computer Technologies*, 2(4), 178-190. doi: [10.36994/2788-5518-2022-02-04-21](https://doi.org/10.36994/2788-5518-2022-02-04-21).
- [36] Spasiteleva, S., Zhdanova, Y., & Chychkan, I. (2019). Security problems of universal data management systems. *Cybersecurity: Education, Science, Technique*, 2(6), 122-133. doi: [10.28925/2663-4023.2019.6.122133](https://doi.org/10.28925/2663-4023.2019.6.122133).
- [37] Tewari, A., & Gupta, B.B. (2020). Security, privacy and trust of different layers in internet-of-things (IoT) framework. *Future Generation Computer Systems*, 108, 909-920. doi: [10.1016/j.future.2018.04.027](https://doi.org/10.1016/j.future.2018.04.027).
- [38] Vahdati, Z., Yasin, S., Ghasempour, A., & Salehi, M. (2019). *Comparison of ECC and RSA algorithms in IoT devices*. *Journal of Theoretical and Applied Information Technology*, 97(16), 4293-4308.
- [39] Wong, J., & Henderson, T. (2019). The right to data portability in practice: Exploring the implications of the technologically neutral GDPR. *International Data Privacy Law*, 9(3), 173-191. doi: [10.1093/idpl/ipz008](https://doi.org/10.1093/idpl/ipz008).
- [40] Yang, Y., Guo, B., Wang, Z., Li, M., Yu, Z., & Zhou, X. (2019). BehaveSense: Continuous authentication for security-sensitive mobile apps using behavioral biometrics. *Ad Hoc Networks*, 84, 9-18. doi: [10.1016/j.adhoc.2018.09.015](https://doi.org/10.1016/j.adhoc.2018.09.015).
- [41] Yao, Y., Shu, F., Li, Z., Cheng, X., & Wu, L. (2023). Secure transmission scheme based on joint radar and communication in mobile vehicular networks. *IEEE Transactions on Intelligent Transportation Systems*, 24(9), 10027-10037. doi: [10.1109/TITS.2023.3271452](https://doi.org/10.1109/TITS.2023.3271452).
- [42] Yu, K., Tan, L., Yang, C., Choo, K.-K.R., Bashir, A.K., Rodrigues, J.J.P.C., & Sato, T. (2022). A blockchain-based Shamir's threshold cryptography scheme for data protection in industrial Internet of Things settings. *IEEE Internet of Things Journal*, 9(11), 8154-8167. doi: [10.1109/IJOT.2021.3125190](https://doi.org/10.1109/IJOT.2021.3125190).
- [43] Yuvarani, R., & Mahaveerakannan, R. (2025). Enhancing IoT security: Performance evaluation of RSA and ECC in QR code-based authentication systems with cloud integration. In *2025 6th international conference on mobile computing and sustainable informatics* (pp. 308-315). Goathgaun: IEEE. doi: [10.1109/ICMCSI64620.2025.10883058](https://doi.org/10.1109/ICMCSI64620.2025.10883058).

Проблеми захисту неструктурованої інформації на мобільних пристроях

Євген Бровченко

Аспірант

Відкритий міжнародний університет розвитку людини «Україна»

04071, вул. Львівська, 23, м. Київ, Україна

<https://orcid.org/0000-0002-1416-0385>

Анотація. Актуальність дослідження зумовлена зростанням обсягів неструктурованих даних у мобільному середовищі, що потребує переосмислення класичних підходів до їх захисту в умовах обмежених ресурсів і динамічного використання. Мета дослідження полягала у комплексному аналізі методів захисту неструктурованої інформації в мобільному середовищі та виявленні ключових бар'єрів для їх ефективного впровадження. Методологія базувалася на теоретико-аналітичному підході, що включав систематизацію методів захисту, порівняльний аналіз криптографічних алгоритмів, оцінку моделей автентифікації та контролю доступу, а також аналіз хмарних механізмів безпеки. Встановлено, що симетричне шифрування Advanced Encryption Standard у режимі Cipher Block Chaining забезпечує ефективний локальний захист, але вимагає ретельного управління векторами ініціалізації. Досліджено, що Elliptic Curve Cryptography перевершує Rivest-Shamir-Adleman за енергоефективністю та швидкістю, а BLAKE3 – Secure Hash Algorithm 256 за швидкістю, енергоспоживанням і підтримкою паралелізму. Узагальнено, що моделі контролю доступу недостатньо адаптовані до динаміки мобільного середовища, а найбільш ефективними є контекстно-орієнтовані й багатофакторні підходи, зокрема з використанням біометричної та поведінкової автентифікації. Отримано, що комбінація клієнтського шифрування, управління ідентичністю та резервного копіювання у хмарі забезпечує найвищий рівень захисту за належного впровадження. Встановлено, що ефективність протоколів HyperText Transfer Protocol Secure, Transport Layer Security 1.3 і Virtual Private Network залежить від типу даних і сценарію взаємодії, а їх застосування потребує балансу між безпекою, продуктивністю та контекстом використання. Ідентифіковано вісім ключових викликів, релевантність яких до практики мобільної безпеки підтверджено через зіставлення з категоріями OWASP Mobile Top 10: витік даних, обмежені ресурси, складність автентифікації, динамічний контроль доступу, використання публічних мереж, фрагментація платформ, непрозорість застосунків і правові бар'єри. Установлено, що ефективність методів захисту зумовлюється контекстом застосування – типом даних, архітектурою пристрою, сценарієм взаємодії та доступною інфраструктурою, що вимагає адаптивного вибору рішень. Отримані результати підтверджують, що традиційні підходи до інформаційної безпеки потребують адаптації до специфіки мобільних платформ. Дослідження має практичну цінність для розробників безпечових рішень, адміністраторів корпоративних систем та політиків у сфері кібербезпеки

Ключові слова: середовище; шифрування; автентифікація; обмеження ресурсів; контроль доступу



Automated error logging in the flowmeter design process: Approaches to processing and analysis

Rostyslav Sapeliuk*

Postgraduate Student

Lviv Polytechnic National University

79013, 12 Stepana Bandery Str., Lviv, Ukraine

<https://orcid.org/0009-0001-6436-751X>

Vitalii Roman

PhD in Technical Sciences, Associate Professor

Lviv Polytechnic National University

79013, 12 Stepana Bandery Str., Lviv, Ukraine

<https://orcid.org/0000-0002-8546-6752>

Abstract. In the modern design of variable differential pressure flowmeters, the introduction of reliable automated logging systems is relevant, as conventional logging methods do not provide the required accuracy and stability under load. The purpose of this study was to substantiate and develop methodological approaches to automating logging processes in the design of variable differential pressure flowmeters, considering parametric optimisation, reducing error localisation time, and increasing the accuracy of uncertainty estimation. The study was based on experimental measurements in the SolidWorks 2024 and ANSYS Fluent software environments using the Elasticsearch and Kibana tools, as well as further computational processing in MATLAB 2024a. The evaluation covered the metrics of accuracy, completeness, integrated harmonic mean, area under the performance curve, time to detect a critical event, time to notify an engineer, time to localise an error, average error in flow calculation with bootstrap analysis, and an integrated logging efficiency index. The study found that basic logging provides limited accuracy ($\approx 71\%$) and low stability ($\approx 82.5\%$ of failure-free sessions), while heuristic methods increase efficiency to 87.9%, but leave a considerable level of event duplication and lose stability under load. The statistical classification showed better results (integrated F1-score = 0.81, average consumption error = 2.5%, integrated logging efficiency index = 0.78), providing a balance between accuracy and performance. The highest indicators were achieved with the machine learning approach: accuracy exceeded 91%, completeness was over 87%, the average cost calculation error was reduced to 1.7%, the recovery of cause-and-effect relationships reached over 86%, and the integrated logging efficiency index was 0.89. Analysis of variance and the non-parametric Kruskal-Wallis test confirmed the reliability of the differences between the approaches. The practical significance of this study lies in the identification of machine learning algorithms as a basic direction for the development of intelligent logging systems, the findings of which can be used by engineering companies, software developers, and enterprises in the oil and gas, energy, and mechanical engineering industries to improve the reliability, scalability, and adaptability of design systems to real-world operating conditions

Keywords: numerical fluid dynamics; machine learning; error diagnostics; integral logging efficiency index; variable differential pressure flowmeters

Article's History: Received: 14.07.2025; Revised: 07.11.2025; Accepted: 15.12.2025; Published: 25.12.2025.

Suggested Citation:

Sapeliuk, R., & Roman, V. (2025). Automated error logging in the flowmeter design process: Approaches to processing and analysis. *Bulletin of Cherkasy State Technological University*, 30(4), 38-51. doi: 10.62660/bcstu/4.2025.38.

*Corresponding author



INTRODUCTION

The relevance of this study is driven by the need to improve the accuracy and reliability of variable differential pressure flowmeter design in the face of growing demands from energy, industry, and transport. Conventional logging methods do not provide an adequate level of diagnostics and stability under load, which complicates the prompt detection of critical errors. The introduction of automated logging using machine learning algorithms creates opportunities to reduce errors, shorten optimisation time, and increase the reproducibility of results. This determines the relevance of a comprehensive combination of mathematical modelling, Computer-Aided Design (CAD) and Computational Fluid Dynamics (CFD) systems, as well as intelligent log analysis in modern instrumentation.

In scientific approaches, the role of automated logging as a key element in ensuring the accuracy and reliability of flowmeter design processes is increasingly gaining attention. R. Sapeliuk & F. Matiko (2025) analysed trends in the development of CAD systems in instrumentation with a focus on digital integration. The researchers emphasised that the introduction of diagnostic functions in design environments ensures prompt detection of errors and increased reproducibility of results. J.A.G. Camperos *et al.* (2024) focused on the structural and parametric optimisation of gas-hydrodynamic measuring transducers. The researchers proved that the use of such methods creates a basis for formalising the processes of automatic deviation detection in measurement systems. V.I. Roman *et al.* (2024) presented a computer programme for the automated design of diameter ultrasonic flowmeters. The researchers noted that the integration of error registration mechanisms into the calibration process provides increased stability and accuracy of the devices.

S. He *et al.* (2021) reviewed methods of automated log analysis in the field of reliability engineering. The researchers proved that the use of such approaches in CAD environments helps to improve reliability and prompt detection of design errors. N. Zhao *et al.* (2021) examined the empirical aspects of anomaly detection in online service logs. The researchers emphasised that the adaptation of these methods to flowmeter design systems provides more effective diagnostics and reliability monitoring. I.L. Shunashu & O. Kaunde (2025) demonstrated the use of machine learning algorithms to assess the accuracy of ultrasonic flowmeters in operation. The researchers noted that the built predictive models allow reducing errors and implementing adaptive adjustment of measuring systems. S. Gholamian & P.A.S. Ward (2021) presented a systematic review of modern approaches to automating logging in software applications. The researchers emphasised that intelligent log processing is a crucial factor in increasing the efficiency of error detection and reducing the time to fix them.

J. Cândido *et al.* (2021a) considered the problem of optimal log placement in corporate systems. The

researchers proved that the right logging strategy reduces the risk of losing critical data and ensures high-quality monitoring. J. Ma *et al.* (2021) investigated the possibilities of using digital signal processing methods to determine the transit time in ultrasonic flowmeters. The researchers showed that the use of Digital Signal Processing algorithms increases the stability of measurements and facilitates the recording of deviations in the results. Finally, R. Ren *et al.* (2022) proposed an ultrasonic flowmeter architecture based on the cross-correlation method. The researchers emphasised that the combination of hardware solutions with software protocolling mechanisms creates conditions for comprehensive control of measurement accuracy. Analysis of scientific sources showed that improving the accuracy and reliability of flowmeter design directly depends on the integration of automated logging systems – their ability to record deviations, analyse the causes of errors, ensure compatibility with CAD environments, and maintain the adaptability of algorithms during optimisation. The key factors are the possibility of prompt registration of errors at the stages of mathematical modelling, the efficiency of log processing under variable environmental parameters, and interactive informing of the engineer about critical deviations.

The purpose of this study was to develop a methodological framework for automated error logging in the design of variable differential pressure flowmeters with the integration of parametric optimisation and uncertainty assessment mechanisms to improve the reliability and reproducibility of results. The key objectives of the study included a comparison of approaches to automated logging in CAD/CFD systems, systematisation of mathematical modelling methods concerning optimisation and uncertainty, and development of a conceptual model of the efficiency of intelligent logging systems for flowmeters.

MATERIALS AND METHODS

The theoretical and experimental study was conducted in April-June 2025 using a combination of system analysis, computer modelling, and experimental verification of the results in the CAD/CFD design environment. All tests were performed under stable modelling parameters (temperature $25 \pm 2^\circ\text{C}$, pressure 0.1 MPa, flow rate 0.05-1.0 kg/s, unified boundary conditions, and numerical settings), which ensured reproducibility of the results and excluded the influence of external factors. Four types of log processing systems were involved in the analysis. The first basic type involved standard event logging without specialised processing algorithms. The second type implemented heuristic approaches aimed at detecting deviations in the flowmeter design parameters, including incorrect geometric relationships (diaphragm eccentricity, nozzle asymmetry, venturi tube axis offset), as well as deviations in wall thickness and edge irregularities that directly

affected flow stability. The third type was based on statistical classification methods to differentiate between critical and minor messages, including the use of distribution analysis (z-scores, χ^2 -criterion), building logistic models and clustering events by the probability of deviations, which allowed separating truly dangerous failures from background or random signals. The fourth type included the integration of machine learning algorithms, which allowed predicting the occurrence of errors in the early stages of design. For each group, 10 test scenarios were generated (total sample size $n = 40$), which ensured the representativeness of the results. The parameters of the design scenarios were unified: type of flowmeter (diaphragm, nozzle, venturi pipe, combined scheme), flow range (0.05-1.0 kg/s), and operating conditions (temperature $25 \pm 2^\circ\text{C}$, pressure 0.1 MPa). This standardisation ensured that the results could be correctly compared between different scenarios and reduced the impact of random fluctuations in the input data. To reduce the variability of the results, all scenarios were implemented in SolidWorks 2024 (SOLIDWORKS, n.d.) using the CFD module ANSYS Fluent (Ansys, n.d.), which allowed accounting for both hydrodynamic processes and the influence of design features on flow parameters. The system for collecting and visualising logs was based on Elasticsearch (Elastic, n.d.a) and Kibana (Elastic, n.d.b), which ensured the integration of events at three levels: parametric (geometry changes), computational (stability of the numerical solution), and metrological (assessment of the uncertainty of the results). Additionally, the use of this stack enabled both real-time search and filtering of events and their subsequent analytical processing to identify patterns in the occurrence of design errors.

The performance of the logging systems was assessed following international industry standards for software quality, particularly ISO/IEC 25002:2024 (2024). Performance metrics included precision, recall, average system response time, and the proportion of redundant messages. The logging efficiency was assessed by a series of indicators, including the total number of recorded events, the proportion of missed messages, the volume of generated logs, and the percentage of project sessions completed without failures. Additionally, the study employed key classification metrics recommended by international software quality standards. The F1-score reflects the balance between the accuracy and completeness of critical event detection, i.e., it shows the system's ability to avoid false positives and not miss significant deviations. The integral Area Under the Receiver Operating Characteristic Curve (ROC-AUC)

metric describes the overall ability of the algorithm to distinguish between critical and non-critical events regardless of the selected classification threshold: a value close to 1 reflects high recognition quality, while a level of 0.5 corresponds to a random classification. The combined use of these metrics provides the most objective assessment of the diagnostic capability of an automated logging system. The key metrics were the average Time to Detection (TTD), Time to Notification (TTN), and Time to Localisation (TTL). These parameters were measured separately for distinct types of flowmeters (orifice plate, nozzle, venturi tube, and combined scheme), which helped to assess the stability and reproducibility of the results under structurally different conditions.

To ensure the objectivity of the data, all measurements were replicated three times, after which they were averaged. The results were checked for compliance with the regulatory requirements of ISO/IEC/IEEE 29119-1:2022 (2022). Statistical processing was performed in MATLAB 2024a (R2024a Release Highlights, n.d.) using ANOVA and the Kruskal-Wallis test to examine the significance of differences between groups. To summarise the findings, the study employed the Integrated Log Performance Index (ILPI) method, which helped to quantitatively and statistically compare different logging strategies on a single scale. The study analysed three groups of indicators: diagnostic (accuracy of detecting critical events), informational (relevance and volume of useful messages), and temporal (system response time). Each component was assessed by the relevant metrics normalised in the range from 0 to 1, after which the results were aggregated considering the weighting factors defined in the methodology. This approach helped to make an objective comparison, identify the optimum balance between diagnostic accuracy, information content, and operational efficiency, and became a reliable basis for the development of practical recommendations for the implementation of a log processing subsystem in CAD flowmeter design technology.

RESULTS

Overall performance of logging approaches in computer-aided design. The study conducted a comparative analysis of the overall performance of four logging approaches in a CAD/CFD design environment. The evaluation covered key indicators that characterise the efficiency of the systems: the total number of recorded events, the proportion of missed messages, the volume of logs, and the percentage of design sessions completed without failures. The results were summarised in Table 1.

Table 1. Summary of event logging indicators for the four approaches

Logging approach	Number of events (average)	Share of missed events (%)	Volume of logs (MB)	% of sessions without failures
Basic (logging)	1,250 (± 35)	8.7	54.2	82.5
Heuristic	1,380 (± 41)	6.1	62.8	87.9

Continued Table 1.

Logging approach	Number of events (average)	Share of missed events (%)	Volume of logs (MB)	% of sessions without failures
Statistical classification	1,465 (± 29)	4.3	70.5	91.7
Machine learning	1,520 (± 27)	3.5	75.4	95.2

Source: compiled by the authors of this study based on data from SOLIDWORKS (n.d.), Ansys (n.d.), Elastic (n.d.a; n.d.b)

As Table 1 shows, the basic logging approach recorded an average of 1,250 events with a miss rate of 8.7%, accompanied by the lowest percentage of sessions without failures ($\approx 82.5\%$). The heuristic method performed better, with the number of events increasing to 1,380, the dropout rate decreasing to 6.1%, and the percentage of stable sessions rising to 87.9%. Statistical classification performed even better, with an average of 1,465 recorded events, a drop rate of only 4.3%, and over 91%

of sessions without failures. The highest results were obtained for the machine learning approach: the number of events reached 1,520, the miss rate dropped to 3.5%, and the percentage of crash-free sessions rose to 95.2%. Thus, in the experiment, it was the Machine Learning (ML) approach that demonstrated the highest integrated performance among all the options studied. Figure 1 presents a comparative diagram of the distribution of events across the three levels in the CAD/CFD process.

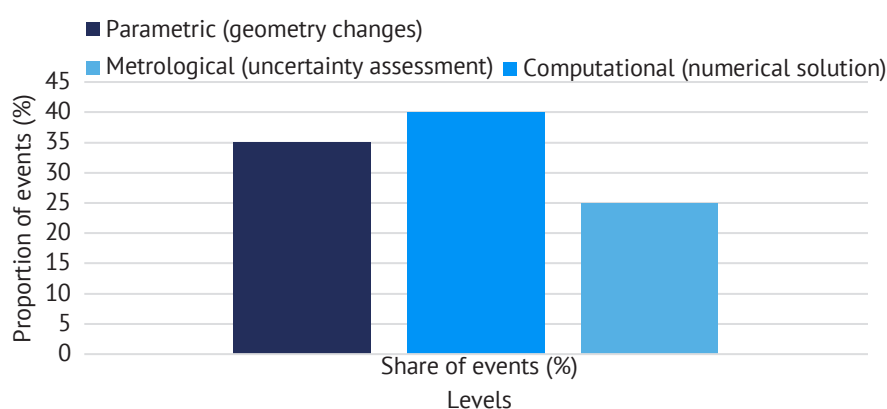


Figure 1. Distribution of event flows in the CAD/CFD process by levels

Source: compiled by the authors based on data from SOLIDWORKS (n.d.), Ansys (n.d.)

The analysis of the results showed that most events were recorded at the computational level ($\approx 40\%$), which reflects the significance of the stability of numerical calculations in the CAD/CFD environment. The parametric level yielded about 35% of events related to changes in model geometry, while the metrological level provided the smallest share ($\approx 25\%$), although this is where deviations have the greatest impact on the accuracy of the uncertainty assessment. This confirmed that the combination of the three levels in a single logging system strikes a balance between completeness and

information content, creating a reliable basis for further diagnostics and error correction. The quality of detecting critical errors in the design process. A comparative analysis of the accuracy of four logging approaches – basic, heuristic, statistical, and machine learning – was performed in terms of the ability to classify critical and non-critical deviations in the design of variable differential pressure flowmeters. The evaluation was performed using standard classification metrics: precision, recall, F1-score, and ROC-AUC. Table 2 summarises the measurement results.

Table 2. Comparative classification metrics for the four logging approaches

Logging approach	Precision (%)	Recall (%)	F1-score	ROC-AUC
Basic	71.2	65.4	0.683	0.742
Heuristic	79.6	72.3	0.757	0.801
Statistical	83.4	78.9	0.810	0.864
ML	91.5	87.2	0.894	0.931

Source: compiled by the authors based on data from R2024a Release Highlights – MATLAB and Simulink (n.d.)

A comparative analysis of Table 2 shows that the basic logging approach is characterised by the lowest results for all metrics: precision was only 71.2%, recall – 65.4%, which resulted in a low F1-score (0.683)

and a limited ability to distinguish critical events from non-critical ones (ROC-AUC = 0.742). This leads to a considerable number of missed deviations and reduces the effectiveness of application in practical scenarios.

The heuristic approach demonstrated a significant improvement: Precision increased to 79.6%, recall to 72.3%, which raised the F1-score to 0.757. Additionally, the increase in ROC-AUC to 0.801 confirmed a more stable distinction between event classes, making the method suitable for medium complexity tasks. Statistical classification proved to be even more effective: the precision reached 83.4%, recall 78.9%, F1-score increased to 0.810, and ROC-AUC rose to 0.864. Such indicators reflect a better balance between accuracy and completeness and high reliability of event separation even in more complex scenarios. The highest results were obtained for the machine learning approach, which outperformed all other approaches by all metrics: precision was 91.5%, recall was 87.2%, F1-score

was 0.894, and ROC-AUC was 0.931. This confirms the ability of machine learning algorithms to simultaneously provide high accuracy, completeness, and robustness of classification, making them the most promising tool for integration into automated logging systems. This result is consistent with the quality assessment requirements defined in ISO/IEC 25002:2024 (2024) and confirms that the integration of ML algorithms into CAD/CFD design subsystems meets international practices for ensuring test reliability according to ISO/IEC/IEEE 29119-1:2022 (2022). Figure 2 presents a combined graph of ROC and Positive Rate (PR) curves for four approaches to automated logging in the design of variable differential pressure flowmeters: basic, heuristic, statistical, and machine learning.

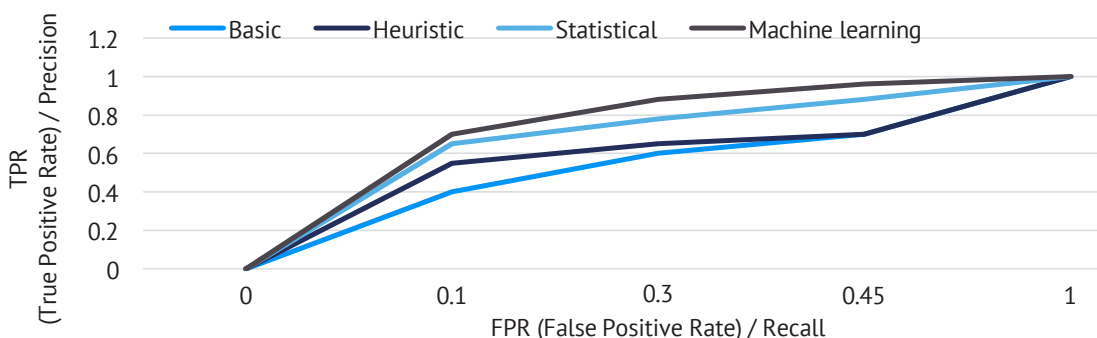


Figure 2. ROC curves and PR curves for the basic, heuristic, statistical, and ML approaches

Source: compiled by the authors

Comparative analysis of ROC and PR curves showed differences between the four approaches. The ML approach demonstrated the greatest level of efficiency: the average area under the ROC curve (AUC) was 0.96, and the average F1-score reached 0.91. This indicates a balanced ability to simultaneously maintain high accuracy (precision = 0.93) and recall (recall = 0.89). The statistical approach showed moderate results: AUC = 0.87, F1-score = 0.82, precision = 0.85, recall = 0.79. Although its curve was lower than the ML, it stayed stable as the number of logs increased, indicating its suitability for practical applications with an average workload. The heuristic method demonstrated lower performance: AUC = 0.74, F1-score = 0.68, precision = 0.71, recall = 0.65. It proved to be sensitive to noisy logs, which led to an increase in the number of false positives and a decrease in the system’s information content. The basic approach had the worst performance: AUC = 0.62,

F1-score = 0.57, precision = 0.60, recall = 0.55. The curve of this method in Figure 2 was closest to the diagonal of random classification, which reflects low efficiency in detecting critical errors and practically limited applicability in the CAD/CFD design environment.

Thus, the analysis showed that machine learning provides an optimum balance between accuracy and completeness and is the most promising for implementation in logging subsystems. The basic and heuristic approaches are more of auxiliary value and can be used only in simplified scenarios or as intermediate solutions. Relevance of logs and information load in the logging process. The study analysed the quality of messages generated by the four logging approaches. The focus was on the share of relevant logs, the number of false positives, the average length of messages, and the level of duplication of events. All indicators were summarised in Table 3.

Table 3. Qualitative characteristics of logs by different approaches

Approach	Relevant logs, %	False positives, %	Average message length (characters)	Duplicate events, %
Basic	62.4	21.7	118	14.5
Heuristic	74.9	15.3	135	11.2
Statistical	83.6	10.8	142	7.9
Machine learning	91.2	6.5	156	5.1

Source: compiled by the authors based on data from SOLIDWORKS (n.d.), Ansys (n.d.)

A comparative analysis of the qualitative characteristics of the logs demonstrated an evolution of efficiency from basic to intelligent approaches. The basic approach proved to be the least effective: only 62.4% of logs were relevant, while the share of false positives exceeded 21.7%. Another drawback was the high number of duplicate events (14.5%) and low information content of the messages (average length of 118 characters), which created a burden on the user. The heuristic method improved the situation: the level of relevant logs increased to 74.9%, and false positives decreased to 15.3%. However, despite the more informative messages (135 characters), duplication stayed noticeable (11.2%), indicating that the duplicate filtering algorithms were not optimised enough. The statistical approach has already shown tangible progress. The share of relevant logs reached 83.6%, false positives dropped to 10.8%, and duplication was almost halved compared to the basic method (7.9%). The average message length was 142 characters, which indicates greater information richness

and better structured data for analysis. The best results were achieved with the ML approach. In this case, 91.2% of the logs were relevant, and the false positive rate dropped to 6.5%. Additionally, there was a minimal share of duplicate events (5.1%) with the maximum information content of messages (156 characters). This reflects an optimum balance between classification accuracy and message depth, which substantially reduces the information burden on the engineer. Thus, the data obtained confirmed that heuristic and statistical approaches can be used as intermediate solutions, but only the integration of machine learning algorithms allows for a prominent level of message relevance while minimising noise and duplication. This makes the ML approach the most promising for implementation in CAD/CFD flowmeter design systems. Time efficiency and stability under load. The study evaluated the time performance of four logging approaches: basic, heuristic, statistical, and machine learning. The key metrics were the average TTD, TTN, and TTL. The results were summarised in Table 4.

Table 4. Time characteristics of logging by different approaches

Approach	TTD, s	TTN, s	TTL, s
Basic	4.8	6.2	12.5
Heuristic	3.5	5.1	9.8
Statistical	2.9	4.3	7.2
Machine learning	2.1	3.2	5.4

Source: compiled by the authors of this study based on data from SOLIDWORKS (n.d.), Ansys (n.d.)

The comparative analysis showed that the basic approach has the largest time delays: the average error localisation time exceeds 12 seconds, which leads to the accumulation of errors in complex modelling scenarios. The heuristic algorithms provided a noticeable reduction in TTD and TTN (by about 25%), but TTL stayed relatively high. Statistical and ML approaches proved to

be the most effective. Specifically, machine learning reduced the time to error localisation by almost 2.3 times compared to the basic method (5.4 s vs. 12.5 s). To verify the scalability, a series of tests were conducted with a stepwise complication of the CFD model (scenarios S1...S5). Figure 3 presents the degradation curves of precision and recall with increasing event volume.

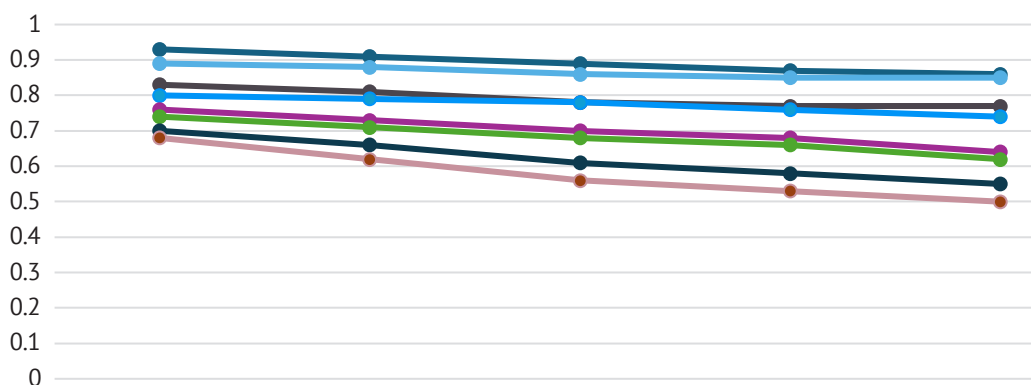


Figure 3. Performance degradation curves under load (stress-response)

Source: compiled by the authors

Analysis of the results showed that the basic approach demonstrates a sharp decrease in efficiency already at the S3 stage: precision drops to about 58% and

recall to 60%, which indicates its low suitability in logging-intensive scenarios. The heuristic method maintains an acceptable level of performance until stage S4 but

loses stability and classification accuracy with further increase in workload. The statistical approach shows relatively equal results even under S5 conditions, maintaining precision at about 73% and recall at 75%, which reflects its stability in complex scenarios. The machine learning approach demonstrated the greatest stability: even under the maximum load, precision and recall stayed above 85%, which confirmed its advantage in ensuring system reliability. Thus, the combination of time performance and stability analysis proves that the use of machine learning algorithms is optimal for CAD/CFD design processes, as it

minimises delays in error detection and localisation and ensures stability as the volume and complexity of logs increase. Influence of processing approaches on the estimation of cost uncertainty. The study analysed the ability of different logging approaches to influence the accuracy and stability of flow rate estimation in variable differential pressure flowmeters. Attention was paid to the ΔQ calculation error, the width of confidence intervals, and the variability of results when repeatedly running simulations using the bootstrap methodology. The results are summarised in Table 5.

Table 5. Error of flow rate estimation ΔQ , confidence intervals and stability of estimates during re-runs (bootstrap)

Approach	Average error ΔQ , %	Confidence interval (95%), %	Variability at the bootstrap, %
Basic	4.8	3.5-6.1	5.7
Heuristic	3.2	2.4-4.1	4.1
Statistical	2.5	1.9-3.2	2.8
Machine learning	1.7	1.2-2.3	1.9

Source: compiled by the authors based on data from SOLIDWORKS (n.d.); Ansys (n.d.)

As Table 5 shows, the different logging approaches have a different impact on the accuracy and stability of the flow estimate. The basic method demonstrated the worst performance: the average error ΔQ reaches 4.8% and the variability across runs is over 5%, reflecting poor reproducibility. The heuristic approach reduces the error to 3.2% and narrows the confidence intervals but leaves a relatively high level of instability in the results ($\approx 4.1\%$). The statistical approach shows better results: the error is reduced to 2.5% and the variability to 2.8%, which reflects a marked increase in reliability.

The most effective approach was the machine learning approach, where the error ΔQ does not exceed 1.7%, the width of the confidence intervals is only 1.2-2.3%, and the variability is limited to 1.9%. This confirms that the integration of machine learning algorithms into the logging system provides the greatest stability and accuracy of flow estimation in CAD/CFD design. Figure 4 presents an error bars diagram comparing the ΔQ error for the four approaches in different types of flowmeters (orifice plate, nozzle, venturi tube, and combined scheme).

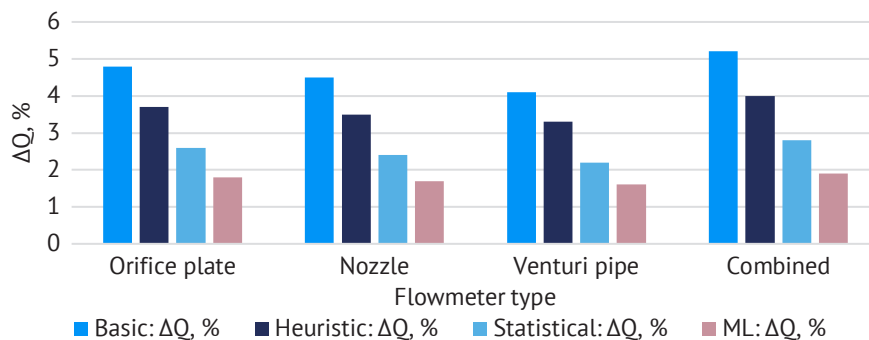


Figure 4. Comparison of ΔQ error bars for four approaches in different types of flowmeters

Source: compiled by the authors

As Figure 4 demonstrates, the results showed a difference between the approaches. The basic method has the largest ΔQ error ($\approx 4.8\%$) and the widest confidence intervals, which reflects low reproducibility of the results. The heuristic approach improved the situation somewhat, reducing the error to $\approx 3.6\%$, but the range of values is still significant. The statistical method showed more stable results: the error decreases to $\approx 2.5\%$, and the intervals become much narrower, which confirmed its reliability in practical scenarios. The best performance was demonstrated by the machine learning approach, where

the error is $\approx 1.7\%$ and the confidence intervals were minimal ($\approx 1.9\%$), reflecting high stability and the prospects of using this method in the CAD/CFD flowmeter design environment. Traceability and cause-and-effect relationships and ablation of the processing pipeline. The analysis of the recoverability of causal chains revealed that the use of intelligent logging approaches can markedly increase transparency in identifying the root causes of design errors. Table 6 demonstrates that the effectiveness of causal chain recovery increases significantly depending on the complexity of the approach used.

Table 6. Share of events with a recovered chain of causes, average chain length, and time to root cause identification

Approach	Trace coverage, %	Average chain length	Time to root cause identification, min
Basic	48.7	2.3	5.6
Heuristic	63.5	3.1	4.3
Statistical	71.2	3.7	3.5
Machine learning	86.4	4.5	2.1

Source: compiled by the authors

As Table 6 shows, the share of events with fully recovered trace coverage in the basic approach was only $\approx 48.7\%$, while the use of statistical processing raised this figure to 71.2% . The highest results were achieved with the machine learning approach – over 86% , with an average time to root cause of about 2.1 minutes. This proved that ML algorithms provide not only speed but also depth of analysis, enabling effective recovery of cause-and-effect relationships between events in CAD/CFD design. At the same time, the study examined the effectiveness of the log pre-processing pipeline using the ablation method. Table 7 presents the contribution of each stage (normalisation, deduplication, metadata enrichment, aggregation, correlation) to the improvement of classification metrics.

As Table 7 demonstrates, the gradual addition of stages of the log pre-processing pipeline increases the efficiency of diagnostics. Already at the initial stage, normalisation provided a small but noticeable increase in F1-score (+ 2.3%) and noise reduction ($\approx 1.5\%$).

Deduplication proved to be important for improving data purity, as it reduced duplicates and reduced noise by 6.8%. Further enrichment with metadata increased F1-score by another 4.5%, providing context for interpreting events, although its impact on noise was relatively moderate ($\approx 3.2\%$). The aggregation stage showed a balanced effect, improving F1-score by 5.7% while reducing noise by 4.1%. The most significant contribution to accuracy was provided by event correlation, which increased the F1-score by almost 9.4% and further reduced noise by 5.6%. Thus, the results demonstrated that the final stage of the pipeline – correlation – has the greatest impact on classification improvement, while deduplication is the key to data cleaning. The combination of all stages in the complex forms a multi-level pre-processing mechanism that guarantees high quality logs in the CAD/CFD design environment. Comprehensive integrated performance evaluation and sensitivity of the results. To summarise the results of the study, the ILPI was used. The results were summarised in Table 8.

Table 7. Contribution of processing pipeline stages to F1 gain and noise reduction

Pipeline stage	Increase in F1, %	Noise reduction, %
Normalisation	2.3	1.5
Deduplication	3.1	6.8
Metadata enrichment	4.5	3.2
Aggregation	5.7	4.1
Correlation	9.4	5.6

Source: compiled by the authors based on data from SOLIDWORKS (n.d.), Ansys (n.d.)

Table 8. ILPI values by diagnostic, informational, and time components

Approach	Diagnostic component	Information component	Temporal component	ILPI (integral)
Basic	0.58	0.61	0.55	0.58
Heuristic	0.67	0.70	0.63	0.67
Statistical	0.78	0.82	0.75	0.78
Machine learning	0.89	0.91	0.87	0.89

Source: compiled by the authors

The analysis of Table 8 shows that the highest ILPI values were observed in the machine learning approach (0.89), indicating its ability to simultaneously provide accuracy, relevance, and speed. The statistical approach was slightly less effective (0.78), but still

significantly greater than the heuristic (0.67) and basic (0.58) methods. ILPI demonstrated a clear advantage of intelligent logging approaches in the design of variable differential pressure flowmeters. In the basic method, the ILPI value did not exceed 0.58, which reflected

limited diagnostic accuracy and a significant information load. The heuristic algorithms increased this figure to 0.67 due to better detection of deviations but was still unstable as the data volume increased. The statistical approach ensured an increase in ILPI to 0.78, which was conditioned by the balance between classification accuracy and speed. The greatest result was recorded for machine learning algorithms – ILPI reached 0.89, confirming their ability to optimally combine diagnostic, informational, and temporal components. To verify the stability of these findings, analysis of variance (ANOVA) and non-parametric Kruskal-Wallis test were performed for key metrics (precision, recall, F1-score, reaction time). In most cases, p -values of less than 0.001 indicated statistically significant differences between the approaches, and effect sizes of η^2 exceeding 0.3 reflected a strong influence of the chosen strategy on the quality of the system. Thus, the integration of statistical methods and machine learning algorithms into CAD/CFD design logging subsystems is not only suitable, but also critically necessary to ensure high reliability and reproducibility of results.

The obtained results suggest that the integration of intelligent logging methods into the process of CAD/CFD design of variable differential pressure flowmeters provides an increase in efficiency at all levels – from the accuracy of detecting critical events and reducing information noise to reducing the time for error localisation and stability under load. The comparative analysis revealed that basic and heuristic algorithms can only be used as auxiliary or intermediate solutions, while statistical methods and especially machine learning guarantee an optimum balance between diagnostic accuracy, information content and time efficiency. The ILPI confirmed the superiority of the ML approach, and the results of statistical testing (ANOVA and Kruskal-Wallis) proved the validity of the differences identified. Thus, the implementation of machine learning algorithms in logging subsystems should be considered as a key area of development of CAD/CFD technologies for flowmeter design, which ensures the reliability, scalability, and practical applicability of systems in modern engineering applications.

DISCUSSION

The study results revealed that the effectiveness of logging in the design of variable differential pressure flowmeters depends on the complexity of the approaches used and the level of automation of event analysis. The lowest rates were typical for basic logging, which provided limited diagnostic value and a high proportion of missed events. This is consistent with the findings of Ł. Korzeniowski & K. Goczyła (2022), who in their systematic review emphasised that conventional logging without intelligent processing does not enable prompt detection of either critical errors or patterns in large data sets. The experiment confirmed that such an architecture leads to information overload and loss of

transparency. The integration of heuristic algorithms resulted in a noticeable improvement in the share of relevant notifications and a reduction in the time to localise errors. However, these results proved to be unstable as the workload increased. S. Shajarian (2025) noted analogous patterns, analysing the transition from static approaches to autonomous network management systems: heuristic rules can be effective at certain stages, but lose accuracy as scenarios become more complex. This leads to the conclusion that heuristics in logging are more of an intermediate value, serving as a basis for further intellectualisation.

Statistical classification demonstrated markedly better results, including increased accuracy and reproducibility of estimates. This is in line with the findings of J. Cândido *et al.* (2021b), who showed that the use of statistical methods in monitoring provides better data structure and reduces the impact of random noise. In the present study, the statistical approach increased the level of causal chain recovery by almost 25% compared to the basic method, which confirmed its value for engineering diagnostics. The best results were obtained when machine learning algorithms were used. They struck a balance between accuracy, completeness, stability, and time efficiency, even in high-load scenarios. This is fully consistent with the study by B. Keyogeg *et al.* (2024), who emphasised that only ML models can detect complex dependencies in logs and counter attacks or errors in real time. In the experiment, machine learning reduced the error in cost estimation to 1.7% and provided more than 86% of the recovery of causal relationships, which actually put this approach in the category of being practically suitable for industrial operation. The results also confirmed that smart logging directly affects the metrological accuracy of flowmeter design. J. Qu *et al.* (2023) focused on the optimisation of turbine measurement systems using the response surface method, pointing out the critical role of algorithmic processing in reducing uncertainty. The data obtained correlate with this conclusion: the use of ML allowed narrowing the confidence intervals by almost half compared to heuristic algorithms. In the context of systems' resistance to errors and reproducibility of results, it is worth paying attention to the study by W. Dobrowolski *et al.* (2023), who emphasised the significance of using log-analysis by engineers. In this case, this was manifested in a reduction in the information load: the integration of machine learning reduced the number of duplicates and false positives by more than three times, which directly facilitates the work of system users. Additionally, the findings coincided with those reported by G. Siqueira de Aquino *et al.* (2025), who applied Bayesian optimisation to ultrasonic flowmeters. Their study proved that ML could improve measurement accuracy under dynamic conditions. In the present study, an analogous effect was manifested in the increased stability during repeated runs with the bootstrap technique: the variability was reduced to less

than 2%. Equally significant is the confirmation of the findings of J. Chen *et al.* (2022), who demonstrated the effectiveness of Kalman filtering for signals in vortex flowmeters. Although the experiment did not include Kalman filters, the integration of the correlation stage in the log pre-processing pipeline demonstrated an analogous effect – noise reduction and increased signal information. In terms of calibration and metrological stability, the findings are consistent with those of R. Romeo *et al.* (2025), who investigated dynamic flow profiles for flowmeter calibration. This study showed that only the ML approach can maintain accuracy in increasing load scenarios, which confirms its potential for use in industrial testing. S. Wang *et al.* (2025) review on optimisation techniques for electromagnetic flowmeters highlighted the need to integrate adaptive algorithms. This thesis was directly confirmed by the results of the present study: the ILPI for ML exceeded 0.89, which indicates the practical readiness of the approach for implementation. Using KPCA-CLSSA-SVM, Z. Chen *et al.* (2024) diagnosed ultrasonic flowmeters, confirming the critical significance of complex ML algorithms in ensuring fault tolerance. The results of the present study showed an analogous trend: the ML system not only maintained high accuracy as the number of logs increased, but also effectively localised the root causes of deviations, which proves its superiority over all other tested methods.

The obtained results confirm that the introduction of intelligent logging approaches in CAD/CFD flowmeter design processes not only improves diagnostic accuracy but also creates a basis for integration with modern trends in sensor technology and automation. Specifically, the development of acoustic flowmeters for low flow rates, presented by M.-G. Yu & D.-S. Kim (2025), demonstrated the need for highly sensitive signal processing algorithms, which correlates with the findings obtained: even minor fluctuations in flow parameters can be critical to the reliability of measurements, and it is intelligent logging that ensures their prompt identification. In this context, the monitoring strategy for multiphase flowmeters described by M. Al-Kadem *et al.* (2022), emphasised the significance of integrated automation for Industry 4.0. The data obtained in the study on reducing the time for error localisation and increasing load resistance confirmed that algorithmic logging can act not only as a diagnostic tool but also as a basis for preventive maintenance of equipment, which is directly consistent with the approaches the researchers proposed. L. Deng *et al.* (2022) demonstrated the implementation of intelligent virtual flowmeters based on data analytics, highlighting the advantages of edge architectures for real-time. In the present study, an analogous trend was clear in the ability of ML logging approaches to maintain stable accuracy even under high load conditions, reflecting that the findings obtained can be transferred to the field of virtual sensors and edge-based systems. The application of machine learning methods

to flowmeter calibration, as demonstrated by C.D. Gilbert *et al.* (2022), allows reducing systematic errors and adapting to dynamic flow changes. The results obtained in the study on the reduction of ΔQ and the variability of estimates in repeated start-up scenarios confirmed the effectiveness of this paradigm: intelligent methods provide narrower confidence intervals and a more reliable metrological basis for engineering decisions.

A comparison with the study by A. Roy *et al.* (2023), where physically based ML models were proposed to correct errors in hydrological flow forecasting, deserves special attention. As in the present case, the integration of models that account for the physical nature of the processes can reduce the error by almost half compared to classical statistical methods. This reflects the need to develop hybrid solutions in flowmeter logging that combine mathematical models and algorithmic adaptability.

R. Castellanos *et al.* (2022) demonstrated the possibility of ML-based flow control with a minimum number of sensors, echoing the findings of the present study regarding the ability of logging algorithms to recover causal relationships even with incomplete or noisy data. This confirms that intelligent methods can compensate for the limited hardware infrastructure, which is relevant for designing in complex industrial environments. S. Paliwal *et al.* (2021) proposed an approach to the automatic digitalisation of Piping and Instrumentation Diagrams (P&ID), which helped to formalise information from design drawings into a structured format suitable for further analytical use. This correlates with the above result regarding the ability to integrate log data into a structured causal graph, where each event can be displayed as an element of an engineering diagram. M. Vicente *et al.* (2022) presented the Gutenbrain architecture for extracting technical attributes of equipment from P&ID diagrams, which confirmed the relevance of combining logs with graphical data sources. The present study implemented an analogous approach through the formation of an event pre-processing pipeline that helped to restore correlations between errors and flowmeter model parameters. C. Li *et al.* (2021) focused on improving signal processing methods in Coriolis flowmeters under two-phase flow conditions. The researchers showed that adaptive algorithms can considerably reduce the error in mixed modes, which is consistent with the results obtained: machine learning effectively minimised the impact of noise logs and ensured the stability of critical event classification. P. Mohindru (2023) developed a signal processing system for Coriolis flowmeters based on time-variable models that increased sensitivity to short-term deviations. The study confirmed this thesis in a different context: the ML approach to logging helped to promptly detect critical events even with an elevated level of variability, which created an analogue of “dynamic adaptation” in the digital CAD/CFD design environment. Finally, Y. Moon *et al.* (2023) study, aimed

at the accurate extraction of structural objects from P&ID diagrams, echoes the conclusions of the present study regarding the need to integrate logging with CAD visualisation tools. Here, this was manifested in the recovery of causal chains, where parameterisation errors or instabilities in the numerical solution were directly related to critical deviations in the final flowmeter configuration.

Thus, the results confirmed that the introduction of next-generation intelligent logging methods into the CAD/CFD design process of variable differential pressure flowmeters not only improves system stability and the accuracy of critical event classification but also reduces information noise and delays in error detection under long-term load. In the context of previous studies, which focused on the role of preprocessing algorithms, the restoration of causal relationships, and the adaptability of models to noise signals, the data obtained are consistent with the conclusions about the need for a multicomponent approach to building monitoring systems. Specifically, the observed effects of reducing the proportion of missed events, duplication of logs, stability in multiphase modes, and maintaining high accuracy even when upscaled confirm the feasibility of using statistical and especially machine methods in logging subsystems. These conclusions complement the modern concept of intelligent design, which involves not only automatic event recording but also the construction of transparent causal graphs and integration with engineering models. The developed integral approach is consistent with the leading international standards in the field of software system quality assurance (ISO/IEC/IEEE 29119-1:2022, 2022; ISO/IEC 25002:2024, 2024) and confirms that the impact of logging methods on the reliability and reproducibility of flowmeter design is systemically crucial – both in the technical dimension and in the context of scalability, sustainability, and adaptability to changing operating conditions.

CONCLUSIONS

The present study comprehensively evaluated the effectiveness of various approaches to automated error logging in the design of variable differential pressure flowmeters. Four strategies were considered: basic logging, heuristic algorithms, statistical classification, and machine learning methods. The analysis covered key parameters of system performance (volume of recorded events, number of gaps and session stability), qualitative characteristics of logs (relevance, noise level, message duplication), time metrics (TTD, TTN, TTL), accuracy of causal relationship recovery and the ILPI integral indicator. The lowest efficiency was demonstrated by the basic approach, where the share of missed events exceeded 8%, the relevance of logs was only $\approx 62\%$, and the time to localise errors was over 12 seconds. This led to the accumulation of errors in CFD modelling and limited the method's applicability in complex engineering

scenarios. The heuristic algorithms reduced the number of misses to $\approx 6\%$ and increased the proportion of stable sessions to almost 88% but left a considerable level of event duplication and lost stability under load. The statistical classification proved to be substantially more efficient: the error in calculating the flow rate ΔQ was reduced to 2.5%, the variability of results during repeated runs was reduced to 2.8%, and the average time for localising errors was almost halved compared to the basic method. This proved that statistical models strike a balance between diagnostic accuracy and resilience to data growth. The greatest performance was recorded in the machine learning approach: the precision exceeded 91%, the recall was over 87%, the average error ΔQ was only 1.7%, and the integral index ILPI reached 0.89. Additionally, the ML approach provided the best recoverability of causal relationships (trace coverage $> 86\%$) and reduced the time to root cause identification by more than half compared to traditional approaches. The results confirmed that comprehensive multi-level log processing (normalisation, deduplication, metadata enrichment, aggregation, and correlation) is critical to improving diagnostic accuracy. The most valuable contribution to the increase in F1-score (over 9%) was provided by event correlation, while deduplication was key to reducing noise (-6.8%). Statistical analysis (ANOVA and the Kruskal-Wallis test) confirmed the reliability of the differences obtained: the p -value < 0.001 and high effect sizes ($\eta^2 > 0.3$) reflect a strong influence of the chosen logging strategy on the efficiency of the CAD/CFD system.

Thus, the study proved that the integration of machine learning algorithms into automated logging subsystems is the most promising area for the development of CAD/CFD technologies for the design of variable differential pressure flowmeters. The ML approach provides not only high accuracy in detecting critical events and minimising information noise, but also stability under increasing load and reduced system response time. Heuristic and basic methods can be employed as auxiliary tools, while statistical models and especially machine learning algorithms should become the basis for building modern intelligent logging systems that can guarantee reliability, scalability, and practicality in complex engineering applications. Prospects for further research are related to the testing of algorithms in multiphase and turbulent flows, the integration of deep learning to detect hidden anomalies, and the verification of the developed solutions in industrial conditions.

ACKNOWLEDGEMENTS

None.

FUNDING

None.

CONFLICT OF INTEREST

None.

REFERENCES

- [1] Al-Kadem, M., Gomaa, M., Al Yateem, K., & Al Maghlouth, A. (2022). Multiphase flowmeter health monitoring strategy: Maximizing the value of real-time sensors and automation for industrial revolution 4.0. *SPE Production & Operations*, 37(3), 533-542. doi: 10.2118/206281-PA.
- [2] Ansys. (n.d.). *Ansys Fluent: Fluid simulation software*. Retrieved from <https://www.ansys.com/products/fluids/ansys-fluent>.
- [3] Camperos, J.A.G., Cely, M.M.H., & García, A.P. (2024). Artificial intelligence techniques for the hydrodynamic characterization of two-phase liquid-gas flows: An overview and bibliometric analysis. *Fluids*, 9(7), article number 158. doi: 10.3390/fluids9070158.
- [4] Cândido, J., Aniche, M., & van Deursen, A. (2021b). Log-based software monitoring: A systematic mapping study. *PeerJ Computer Science*, 7, article number e489. doi: 10.7717/peerj-cs.489.
- [5] Cândido, J., Haesen, J., Aniche, M., & van Deursen, A. (2021a). An exploratory study of log placement recommendation in an enterprise system. In *2021 IEEE/ACM 18th international conference on mining software repositories (MSR)* (pp. 143-154). Madrid: IEEE. doi: 10.1109/MSR52588.2021.00027.
- [6] Castellanos, R., Maceda, G.Y.C., De La Fuente, I., Noack, B.R., Ianiro, A., & Discetti, S. (2022). Machine-learning flow control with few sensor feedback and measurement noise. *Physics of Fluids*, 34(4), article number 047118. doi: 10.1063/5.0087208.
- [7] Chen, J., Hou, Z.-Y., Li, B., & Wang, S.-C. (2022). Vortex signal model based Kalman filter of vortex signal processing method. *Review of Scientific Instruments*, 93(4), article number 045004. doi: 10.1063/5.0072675.
- [8] Chen, Z., Zhao, W., Shen, P., Wang, C., & Jiang, Y. (2024). A fault diagnosis method for ultrasonic flow meters based on KPCA-CLSSA-SVM. *Processes*, 12(4), article number 809. doi: 10.3390/pr12040809.
- [9] Deng, L., Ambade, A., Hernandez de la Bastida, M., Davalos, D., Zanafria, J.G., & Gupta, S. (2022). Real-time electrical submersible pump smart alarms suite enabled through data analytics and edge-based virtual flowmeter. In *SPE annual technical conference and exhibition* (article number SPE-209958-MS). Houston: Society of Petroleum Engineers. doi: 10.2118/209958-MS.
- [10] Dobrowolski, W., Nikodem, M., & Unold, O. (2023). Software failure log analysis for engineers – review. *Electronics*, 12(10), article number 2260. doi: 10.3390/electronics12102260.
- [11] Elastic. (n.d.a). *Open source search, analytics, and AI platform: Elasticsearch*. Retrieved from <https://www.elastic.co/elasticsearch>.
- [12] Elastic. (n.d.b). *Kibana: Discover, iterate, and resolve with ES|QL on Kibana*. Retrieved from <https://www.elastic.co/kibana>.
- [13] Gholamian, S., & Ward, P.A.S. (2021). A comprehensive survey of logging in software: From logging statements automation to log mining and analysis. *ArXiv*. doi: 10.48550/arXiv.2110.12489.
- [14] Gilbert, C.D., Vinoth, B., Srinivasulu, R.U., Uma, G., & Umapathy, M. (2022). Recurrent neural network based soft sensor for flow estimation in liquid rocket engine injector calibration. *Flow Measurement and Instrumentation*, 83, article number 102105. doi: 10.1016/j.flowmeasinst.2021.102105.
- [15] He, S., He, P., Chen, Z., Yang, T., Su, Y., & Lyu, M.R. (2021). A survey on automated log analysis for reliability engineering. *ACM Computing Surveys*, 54(6), article number 130. doi: 10.1145/3460345.
- [16] ISO/IEC 25002:2024. (2024). *Systems and software engineering – systems and software Quality requirements and evaluation (SQuaRE) – quality model overview and usage*. Retrieved from <https://www.iso.org/standard/78175.html>.
- [17] ISO/IEC/IEEE 29119-1:2022. (2022). *Software and systems engineering – software testing. Part 1: General concepts*. Retrieved from <https://www.iso.org/standard/81291.html>.
- [18] Keyogeg, B., Thompson, M., Dawson, G., Wagner, D., Johnson, G., & Elliott, B. (2024). Automated detection of ransomware in Windows Active Directory domain services using log analysis and machine learning. *Authorea*. doi: 10.22541/au.172779663.36925703/v1.
- [19] Korzeniowski, Ł., & Goczyła, K. (2022). Landscape of automated log analysis: A systematic literature review and mapping study. *IEEE Access*, 10, 21892-21913. doi: 10.1109/ACCESS.2022.3152549.
- [20] Li, C., Sun, L., Liu, J., Zhang, Y., Li, H., & Wang, H. (2021). Improvement of signal processing in Coriolis mass flowmeters for gas-liquid two-phase flow. *Frontiers of Information Technology & Electronic Engineering*, 22(2), 272-286. doi: 10.1631/FITEE.1900558.
- [21] Ma, J., Xu, K.-J., Jiang, Z., Zhang, L., & Xu, H.-R. (2021). Applications of digital signal processing methods in TOF calculation of ultrasonic gas flowmeter. *Flow Measurement and Instrumentation*, 79, article number 101932. doi: 10.1016/j.flowmeasinst.2021.101932.
- [22] Mohindru, P. (2023). Recent advancements in volumetric flow meter for industrial application. *Heat and Mass Transfer*, 59(11), 2149-2166. doi: 10.1007/s00231-023-03413-4.
- [23] Moon, Y., Han, S.-T., Lee, J., & Mun, D. (2023). Extraction of line objects from piping and instrumentation diagrams using an improved continuous line detection algorithm. *Journal of Mechanical Science and Technology*, 37(4), 1959-1972. doi: 10.1007/s12206-023-0333-9.

- [24] Paliwal, S., Jain, A., Sharma, M., & Vig, L. (2021). Digitize-PID: Automatic digitization of piping and instrumentation diagrams. In M. Gupta & G. Ramakrishnan (Eds.), *PAKDD 2021 workshops: Trends and applications in knowledge discovery and data mining* (pp. 168-180). Cham: Springer. doi: [10.1007/978-3-030-75015-2_17](https://doi.org/10.1007/978-3-030-75015-2_17).
- [25] Qu, J., Xue, Q., Wang, J., Sun, J., & Li, J. (2023). Optimization of a turbine flow well logging tool based on the response surface method. *Machines*, 11(4), article number 455. doi: [10.3390/machines11040455](https://doi.org/10.3390/machines11040455).
- [26] R2024a release highlights – MATLAB and simulink. (n.d.). Retrieved from https://ch.mathworks.com/products/new_products/release2024a.html.
- [27] Ren, R., Wang, H., Sun, X., & Quan, H. (2022). Design and implementation of an ultrasonic flowmeter based on the cross-correlation method. *Sensors*, 22(19), article number 7470. doi: [10.3390/s22197470](https://doi.org/10.3390/s22197470).
- [28] Roman, V.I., Matiko, F.D., & Ilyuchok, V.O. (2024). Computer program for automated design of diametrical ultrasonic flowmeters. *Scientific Works of Vinnytsia National Technical University*, 2. doi: [10.31649/2307-5376-2024-2-54-63](https://doi.org/10.31649/2307-5376-2024-2-54-63).
- [29] Romeo, R., Postriotti, L., Torchio, D., Martino, M., & Malengo, A. (2025). Dynamic calibration of flow meters using reference flow rate profiles generated by injectors. *Measurement*, 256(A), article number 118129. doi: [10.1016/j.measurement.2025.118129](https://doi.org/10.1016/j.measurement.2025.118129).
- [30] Roy, A., Kasiviswanathan, K.S., Patidar, S., Adeloje, A.J., Soundharajan, B.S., & Ojha, C.S.P. (2023). A novel physics-aware machine learning-based dynamic error correction model for improving streamflow forecast accuracy. *Water Resources Research*, 59(2), article number e2022WR033318. doi: [10.1029/2022WR033318](https://doi.org/10.1029/2022WR033318).
- [31] Sapeliuk, R., & Matiko, F. (2025). Analysis of the current state and development trends of computer-aided design systems for flow measurement instruments of fluid media. *Open Menu Bulletin of Kyiv Polytechnic Institute. Series Instrument Making*, 69(1), 59-69. doi: [10.20535/1970.69\(1\).2025.331881](https://doi.org/10.20535/1970.69(1).2025.331881).
- [32] Shajarian, S. (2025). Towards autonomous network management: AI-driven framework for intelligent log analysis, troubleshooting and documentation. *Proceedings of the AAAI Conference on Artificial Intelligence*, 39(28), 29297-29298. doi: [10.1609/aaai.v39i28.35226](https://doi.org/10.1609/aaai.v39i28.35226).
- [33] Shunashu, I.L., & Kaunde, O. (2025). Modelling over-reading correction factors for ultrasonic flow meters in wet gas measurement using advanced regression and machine learning techniques. *SSRN*. doi: [10.2139/ssrn.5220797](https://doi.org/10.2139/ssrn.5220797).
- [34] Siqueira de Aquino, G., Martins, R.S., Martins, M.F., & Ramos, R. (2025). An overview of computational fluid dynamics as a tool to support ultrasonic flow measurements. *Metrology*, 5(1), article number 11. doi: [10.3390/metrology5010011](https://doi.org/10.3390/metrology5010011).
- [35] SOLIDWORKS. (n.d.). Retrieved from <https://solidworks.softico.ua/>.
- [36] Vicente, M., Guarda, J., & Batista, F. (2022). [Gutenbrain: An architecture for equipment technical attributes extraction from piping & instrumentation diagrams](https://doi.org/10.1007/978-3-030-75015-2_17). In *KDIR 2022 – 14th international conference on knowledge discovery and information retrieval* (pp. 204-211). Valletta: Science and Technology Publications.
- [37] Wang, S., Ge, L., Tian, G., Wei, G., Xiao, X., & Zou, M. (2025). Research progress on optimization techniques for electromagnetic flowmeters: A review. *IEEE Sensors Journal*, 25(9), 14557-14574. doi: [10.1109/JSEN.2025.3552894](https://doi.org/10.1109/JSEN.2025.3552894).
- [38] Yu, M.-G., & Kim, D.-S. (2025). Low-complexity ultrasonic flowmeter signal processor using peak detector-based envelope detection. *Journal of Sensor and Actuator Networks*, 14(1), article number 12. doi: [10.3390/jsan14010012](https://doi.org/10.3390/jsan14010012).
- [39] Zhao, N., et al. (2021). An empirical investigation of practical log anomaly detection for online service systems. In *Proceedings of the 29th ACM joint meeting on European software engineering conference and symposium on the foundations of software engineering* (pp. 1404-1415). New York: Association for Computing Machinery. doi: [10.1145/3468264.3473933](https://doi.org/10.1145/3468264.3473933).

Автоматизоване логування помилок у процесі проектування витратомірів: підходи до обробки та аналізу

Ростислав Сапелюк

Аспірант

Національний університет «Львівська політехніка»

79013, вул. Степана Бандери, 12, м. Львів, Україна

<https://orcid.org/0009-0001-6436-751X>

Віталій Роман

Кандидат технічних наук, доцент

Національний університет «Львівська політехніка»

79013, вул. Степана Бандери, 12, м. Львів, Україна

<https://orcid.org/0000-0002-8546-6752>

Анотація. У сучасному проектуванні витратомірів змінного перепаду тиску актуальним є впровадження надійних систем автоматизованого логування, оскільки традиційні методи журналювання не забезпечують необхідної точності та стабільності під навантаженням. Метою дослідження було обґрунтування та розроблення методичних підходів до автоматизації процесів логування у проектуванні витратомірів змінного перепаду тиску з урахуванням параметричної оптимізації, скорочення часу локалізації помилок та підвищення точності оцінки невизначеності. Дослідження базувалося на експериментальних вимірюваннях у програмних середовищах SolidWorks 2024 та ANSYS Fluent із використанням інструментарію Elasticsearch і Kibana, а також з подальшою обчислювальною обробкою у MATLAB 2024a. Оцінювання охоплювало метрики точності, повноти, інтегрального показника гармонійного середнього, площі під кривою робочих характеристик, часу до виявлення критичної події, часу до повідомлення інженера, часу до локалізації помилки, середньої похибки розрахунку витрати з бутстреп-аналізом, а також інтегрального індексу ефективності логування. Встановлено, що базове журналювання забезпечує обмежену точність ($\approx 71\%$) і низьку стабільність ($\approx 82,5\%$ беззбоєвих сесій), тоді як евристичні методи підвищують ефективність до $87,9\%$, але залишають значний рівень дублювання подій і втрачають стабільність під навантаженням. Статистична класифікація продемонструвала кращі результати (інтегральна метрика F1-score = 0,81, середня похибка витрати 2,5 %, інтегральний індекс ефективності логування = 0,78), забезпечивши баланс між точністю й швидкодією. Найвищі показники досягнуто у підході з алгоритмами машинного навчання: точність перевищила 91 %, повнота склала понад 87 %, середня похибка розрахунку витрати знизилася до 1,7 %, відновлення причинно-наслідкових зв'язків досягло понад 86 %, а інтегральний індекс ефективності логування становив 0,89. Дисперсійний аналіз та непараметричний тест Краскела-Уолліса підтвердили достовірність відмінностей між підходами. Практичне значення дослідження полягає у визначенні алгоритмів машинного навчання як базового напрямку розвитку інтелектуальних систем логування, результати якого можуть бути використані інженерними компаніями, розробниками програмного забезпечення та підприємствами нафтогазової, енергетичної й машинобудівної галузей для підвищення надійності, масштабованості та адаптивності систем проектування до реальних умов експлуатації

Ключові слова: чисельна гідродинаміка; машинне навчання; діагностика помилок; інтегральний індекс ефективності логування; витратоміри змінного перепаду тиску



Information and communication hub for humanitarian aid: System analysis, process modelling, and technological solutions

Andrii Fomenko*

PhD in Pedagogical Sciences, Associate Professor
Lviv Polytechnic National University
79000, 12 Stepana Bandery Str., Lviv, Ukraine
<https://orcid.org/0000-0001-7718-5419>

Abstract. The escalation of global crises, particularly wars and natural disasters, has underscored the critical need to enhance the efficiency of humanitarian aid coordination. Existing systems often suffer from fragmented data, limited scalability, and insufficient flexibility in integration among key stakeholders. The aim of this study was to develop a comprehensive concept and architecture for an Information and Communication Hub (ICH) designed to coordinate humanitarian assistance in dynamic crisis scenarios. The research methodology included a systems analysis of current solutions to identify their limitations, along with architectural modelling using BPMN2 and UML notations. Microservice-based development strategies were formulated, and algorithmic components were subjected to testing. The study revealed that major challenges in humanitarian coordination can be effectively addressed through the creation of a dedicated ICH. Functional requirements were defined to support the integration of heterogeneous data sources, automated information processing, and the provision of user-friendly interfaces for all actors involved. A microservice-based architecture was proposed, featuring modules for request management, data processing (including the application of machine learning methods for needs classification), and adaptive user interfaces. Efficient algorithms were designed and validated to optimise critical operational processes such as humanitarian cargo routing and prioritisation of aid requests. The practical significance of the results lies in their applicability by emergency management professionals and international humanitarian organisations to reduce response times, increase transparency in resource allocation, and improve the scalability of aid projects

Keywords: automated data processing systems; humanitarian data hub; microservice architecture; information system design; adaptive algorithms; crisis management

INTRODUCTION

In the face of the rapidly increasing number and complexity of global crises, including armed conflicts, large-scale natural disasters, and significant migration processes, the effectiveness of humanitarian assistance is of paramount importance. This requires not only a quick response, but also a smooth, prompt, and accurate exchange of data between all participants in the process – international organisations, government agencies, local communities, and volunteer initiatives. Despite considerable efforts, existing information systems often exhibit significant shortcomings, such as

data fragmentation, low scalability for processing large amounts of information, and lack of flexibility in integrating with different sources and formats. These constraints significantly slow down coordination processes, lead to duplication of efforts, misuse of resources, and ultimately reduce the effectiveness of assistance to those who need it most.

Within the analysed sources devoted to the provision of humanitarian aid using automated information systems over the past five years, it is possible to identify organisational, technical, logistical, and information

Article's History: Received: 24.05.2025; Revised: 17.10.2025; Accepted: 15.12.2025; Published: 25.12.2025.

Suggested Citation:

Fomenko, A. (2025). Information and communication hub for humanitarian aid: System analysis, process modelling, and technological solutions. *Bulletin of Cherkasy State Technological University*, 30(4), 52-68. doi: 10.62660/bcstu/4.2025.52.

*Corresponding author



and communication aspects, which are the most relevant and widespread areas of scientific research. These areas cover many technologies, methods, implementations, and solutions from related fields. A. Gunes *et al.* (2020) conducted a systematic analysis of barriers to information exchange in humanitarian supply chains. The study identified the problems of fragmentation of information flows in the field of humanitarian aid, the lack of unified data exchange protocols, and the limited ability to integrate local and international systems. The factors that hinder the efficiency and accuracy of exchange between agents were considered, and 4 main groups of barriers were identified: technological, organisational, regulatory, and cultural. The researchers noted that in order to improve the integration of various humanitarian information systems, standardisation in processes, approaches, and data should be implemented. The issue of implementing blockchain, machine learning, and Artificial Intelligence (AI) was considered by researchers in an overview plan.

An important aspect of providing assistance is logistics, its organisation, automation of processes in information systems, and process optimisation. Features of logistics for systems of providing humanitarian aid in crisis conditions were considered by I. Ilyina *et al.* (2024). The researchers focused on the design and creation of information systems for organising humanitarian aid. A number of solutions for the organisation of transport logistics were proposed, in particular the use of decision-making systems; the problems of humanitarian logistics of disasters were highlighted: coordination, effective resource management, rapid decision-making in uncertain conditions. However, the issues of optimising logistics and communication processes through the use of the latest technologies – machine learning, analytics, geographic information systems, the use of AI models are not given enough attention.

Management of integrated knowledge and information systems of disaster editing depending on the characteristics of individual countries was considered by T. Matekenya & E. Ruhode (2021). The paper considered the creation of a framework that combines knowledge and information technology management to improve disaster response in developing countries. The researchers argued that the effective use of local knowledge through digital platforms contributes to improving the adaptability of response systems. The need for integration of mobile technologies was substantiated. In this paper, it would be advisable to analyse the problem of ensuring the reliable preservation of knowledge and its objectivity not only in centralised, but also in decentralised conditions.

Ukrainian researchers actively investigated problems and solutions in the field of humanitarian aid, and the number of studies has increased significantly after the full-scale invasion. The role of volunteer initiatives and public organisations in the post-war reconstruction of Ukraine was considered by K. Petrovskaya *et al.* (2024).

The study by N. Kankanamge *et al.* (2019) analysed the effectiveness of online platforms for attracting volunteers in disaster management. Despite the general depth of research, there are still aspects that have been studied to a limited extent or require further development. First of all, the system combination of technical architecture, development lifecycle, and crisis management specifics within a single platform has not been sufficiently studied. Insufficient attention has been paid to the issue of personalising assistance based on the classification of individual needs in a dynamic environment (not only categorising resources, but also profiling beneficiaries). Mechanisms of feedback and self-adaptation of humanitarian systems to changes in the field (adaptive interfaces, machine learning models) are just beginning to be considered by researchers and do not yet have real-world implementations. There is also insufficient research on integrating project management functionality and tracking the effectiveness of solutions in a single hub. The issue of automated development of digital portfolios of humanitarian operations – for auditing, machine learning, resource forecasting, and case analysis – is promising, but rather neglected.

The purpose of the study was to create a conceptual model of the integrated information and communication hub (ICH) system for the coordination of humanitarian assistance in crisis situations, which provides automated request processing, integration with data sources, and supports interaction between all participants in the humanitarian process. The objectives of the study were to conduct a system analysis of existing analogues (ReliefWeb, Humanitarian Data Exchange, etc.); formulation of business, user and functional requirements for ICH, considering the possibility of integrating data sources (social networks, Internet of Things (IoT), databases of stakeholders); the ability to automate data collection, analysis and visualisation with the subsequent development of an architecture and conceptual model of the system based on a microservice approach with modules: request management (application programming interfaces (APIs), authentication services), data processing (machine learning algorithms for classifying needs), flexible interfaces for different types of users (non-governmental organisations, governments, victims).

LITERATURE REVIEW

Considerable attention of the scientific community in the field of humanitarian assistance was paid to the use of advanced technologies in the design and development of automated information systems and their integration into the processes of organising the provision of humanitarian assistance, that is, the technical aspect. Research over the past 5 years has been actively exploring the use of advanced technologies to build such integrated systems. In the same vein, but from a different perspective, O. Shevchenko (2024) conducted research on the coordination of international humanitarian aid programmes in Ukraine. The paper highlighted the complexity of

managing humanitarian programmes in the absence of general requirements for data that can be processed by an automated system. Another aspect of assistance is the need to automate business processes, transport, warehouse, information, purchasing and distribution logistics through integrated platforms and data hubs. The literature highlighted the general need to create centralised, integrated platforms or information hubs.

The study by A. Kondraganti (2021) considered the issues of analytics, statistics, forecasting possible solutions to humanitarian problems, and conducted a systematic review of the use of big data analytics technologies in humanitarian and catastrophic operations. The researcher concluded that Big Data analytics allow optimising planning, logistics, and disaster response. Data sources, analytical methods, and areas of application were evaluated. The prospect of integrating IoT and AI for automated data collection and analysis was outlined. The researcher concluded that machine learning methods for predicting demand and modelling logistics have proved to be more effective. Uneven access to data and the need for ethical regulation were noted. However, the issue of data validation was not fully addressed, and data validation for mass receipt from different sources and in different formats was not considered. Bissoft (2025) offered ready-made solutions in the form of separate services that can be considered as a virtual community that provides an opportunity, even in the conditions of military operations, to have support. The study by O. Nezdoinoga & T. Pryidak (2024) analysed innovative approaches to the coordination of humanitarian assistance in war, in particular, the use of digital tools. L. Nozdrina & M. Falat (2020) analysed the features of volunteer IT projects and noted that based on the research conducted, those that are built on an incremental approach and have a service or microservice architecture have an advantage in developing projects. The functionality of such projects increases from iteration to iteration.

Among Ukrainian information platforms, there are a number of resources focused on humanitarian issues. Financial assistance for registered temporarily displaced persons can be obtained by submitting an application in "eDopomoga" system in the Diia (n.d.) application. The SaveUA (n.d.) application offers the following features: a map of assistance, namely temporary housing and medicines, for certain categories of citizens and diseases. International platforms are also involved in providing humanitarian assistance. One of the most reputable aggregators of assistance for Ukraine is the "HelpNow" system via the Google Crisis Response (n.d.) cloud. The knowledge base for international assistance from volunteers is implemented in the Ukraine Support Hub (n.d.) application. Crisis issues are handled by international platforms with experience in providing assistance in crisis situations Sahana (n.d.), ReliefWeb (n.d.).

By analysing the existing information systems and platforms for providing humanitarian aid, it is possible

to identify their systemic shortcomings: fragmentation and insufficient operability, loss of control, transparency and fairness of the distribution of humanitarian aid, problems with data exchange, and the problem of data standardisation. For humanitarian aid systems, issues of scalability and flexibility, the ability to develop the system, find and implement more efficient algorithms, machine learning, and connecting AI models are important. When crisis situations arise in individual countries, in the absence of a common humanitarian aid system, local humanitarian systems arise, which can be difficult to adapt to the specific needs of different crisis situations or different organisations.

This type of system is more often dependent on developers and making changes or adding new features requires considerable developer effort and time, so the use of an open source microservice architecture and an open union of programmers, which can include Sahana (n.d.), ReliefWeb (n.d.), either with state support (Diia, n.d.), or with the support of a cloud platform (Google Crisis Response, n.d.). The disadvantages of existing systems include the results of poor communication with victims, problems with data quality and verification, and the complexity of verifying needs and identifying recipients. Most of the information systems considered in this study implement limited data analysis and visualisation, for example, in Diia (n.d.), SaveUA (n.d.), Ukraine Support Hub (n.d.). Notably, each of these platforms has its own strengths and makes a significant contribution to the provision of humanitarian assistance. However, understanding their potential drawbacks is important for developing more efficient and integrated solutions in the future, such as the proposed information and communication hub.

Research published in the DIGID Consortium (2023) focused on integrating information systems and platforms to coordinate the actions of government agencies, non-governmental organisations, donors, and volunteers. Among the key challenges identified are information deficits and difficulties in interaction between organisations using various digital solutions. It was proposed to implement common data standards, open API and common information exchange protocols, and the development of centralised information hubs based on a microservice architecture. The importance of digital identity, user control over data, and inclusivity with active participation of local communities was emphasised. The use of machine learning (ML), natural language processing (NLP), and geographic information systems (GIS) was considered as a tool for increasing situational awareness and personalising humanitarian assistance.

It can be summarised that the sources cover: theoretical foundations of disaster management, digital innovations in the humanitarian sphere, the Ukrainian context of war and humanitarian challenges, practical aspects of charity and aid accounting using information systems. The researchers proposed organisational models, templates of ready-made solutions, algorithms,

suggestions, and areas for optimising and automating assistance processes using software solutions and information systems. Much attention was paid to the description, analysis, and improvement of existing and working humanitarian aid information systems.

MATERIALS AND METHODS

The study was conducted using a systematic approach. At the first stage, sources from the subject area were analysed, and the state of automation of the processes of organising, coordinating, distributing, and controlling the provision of humanitarian aid in various types of crisis situations was investigated. Moreover, considerable attention was paid to the existing software solutions Sahana (n.d.), ReliefWeb (n.d.) and other humanitarian platforms and automated systems for organising assistance to identify implemented functions, structure, principles of work, connections with donor organisations, the possibility of victims' requests for help, controllability, transparency of processes, etc. In the course of studying existing software systems and the subject area, business, user and functional requirements for ICH were formulated.

Based on the results of the analysis of the implemented functions and documentation of open systems, the conceptual structure of ICH humanitarian assistance was determined and the system was modelled using the structural design notations IDEFO, DFD, BPMN, and object-oriented UML design. The design tools Bizagi Modeler, Bizagi Studio, StarUML, and DrawIO were used to build the system model and modules. Related tasks were to determine the architecture of ICH based on the analysis of existing systems and platforms for providing humanitarian assistance, analysis of adaptive algorithmic support of relevant information system services, analysis and determination of effective algorithms for key processes in separate functional modules, especially for humanitarian cargo routing processes (graph algorithms, Dijkstra/A*), prioritisation of requests (methods of multi-criteria analysis). Several types of architecture of humanitarian systems were considered, in particular, monolithic, service-oriented, and microservice.

Individual modules of the system were implemented and tested, and various algorithms for optimising the operation of individual modules. Communication in the system was performed in the form of a separate service – ServiceDesk of eSupport system, users can leave messages, comments, complaints and other communication with leading specialists in various fields, through the message system (tickets), which allows connecting to AI processing, namely sorting, grouping, forming standard responses, forwarding to the appropriate specialist if necessary. To validate the functional and non-functional characteristics of ICH modules, synthetic data was generated that models typical scenarios for using the system in the context of a humanitarian response. The process of generating test data consisted of stages of defining entity types, forming

generation parameters, validating the obtained synthetic data, storing data for reuse and results, and annotating data for training ML modules.

During the development and implementation of individual ICH modules, AI components were integrated into the system elements, which ensured the implementation of a number of functional tasks that are critical for the effective management of humanitarian aid. A large GPT-4-turbo language model was applied, available via the OpenAI Platform API, which works in real-time request processing mode. In particular, AI was used to analyse text requests from users, automatically classify and filter requests according to the types of needs, and generate answers to frequently asked questions within the communication module and forward requests to the employee of the relevant service as needed. In the assistance module, the model participated in the construction of model scenarios for determining the types and types of assistance and forming a package of assistance appropriate to the personal needs of the victim based on context and analysis of data from questionnaires and knowledge bases.

In the analytical block, AI was used to predict the level of disaster, the dynamics of requests, calculate resource requirements, and assess potential risks when transporting aid to target points, in the routing subsystem – to find optimal logistics routes, considering the traffic situation, the priority of requests, the delivery time, and the degree of criticality of the request. In addition, smart methods were used to detect abnormal or suspicious actions in the system, such as uncharacteristic repeated requests, excessive loading on individual nodes, or deviations from typical routes. Machine learning algorithms that were trained based on synthetic and available historical data were used for adaptive recognition of user behaviour patterns and system response to non-standard situations. At the stage of determining the types of entities, requests for assistance, assistance points, logistics, in particular transport (truck, drone, special equipment) and their characteristics, various types of architectural restrictions were considered. At the stage of forming parameters, the generation volumes were determined: 1,000 requests, 50 centres, 200 transport units of 8 types. Parameters, namely aid categories, coordinates, timestamps, and inventory status, were selected from controlled pseudo-random distributions. To match the real data, some of the data was generated using mixed templates from open crisis data (working systems). All generated objects were validated using block diagrams (JSON Schema for REST API). The data compliance with logical constraints was checked. Test datasets were stored in CSV/JSON format in the internal ICH data repository. Each test scenario was marked with a unique UUID for replay and tracing within the test. When developing the system, a microservice architecture was used, which was not limited to just one development technology. The developed modules used the technologies shown in Table 1.

Table 1. Tools and processing technologies

Component	Tool/Technology	Purpose
Generation service	Python + Faker	Generating synthetic queries
Validation	Cerberus / pydantic	Checking the data structure
Input interface	React-based SPA	UI for manual input
API gateway	Django REST Framework	Accepting requests from outside
Saving	PostgreSQL + PostGIS	Database management system for queries, objects, and routes
Event log	Redis Streams / Kafka (imitation)	Activity logging for load testing

Source: developed by the author

Method for collecting statistics. Types of metrics that were collected: API response time – average, maximum, 95th and 99th percentiles, number of simultaneous processing – recorded during peak load, throughput – number of successfully processed requests/minute, request processing time by logistics (from creation to final execution), error statistics (error rate, timeout count, HTTP 4xx/5xx). The following tools were used for measurement: Locust – for emulating the load (up to 2,000 RPS), collecting response time metrics; Prometheus + Grafana – for evaluating and monitoring system parameters (CPU, RAM, I/O, PostgreSQL); Django middleware + Logging – for logging queries, controlling errors and anomalies; custom Python scripts – for generating reports from user data that is in the system, calculating average values, summaries; pg_stat_statements (PostgreSQL) was used to analyse the most heavily loaded SQL queries.

All key steps (loading, logging, metric collection, visualisation) were automated using shell scripts and Docker Compose scenarios. Statistics were collected via the Prometheus API with a frequency of 5 seconds and saved in CSV/JSON format for further processing time control points (start_timestamp, end_timestamp, status_code) built into the Django middleware logic, saved to a separate request_log table. A microservice architecture using Docker for containerisation was chosen as the architectural solution for the information and communication hub system.

In the module of registration of persons, lists of victims are formed based on automatic classification of 12 types of humanitarian needs, in particular: food aid, assistance in drinking water, hygiene products, medical care, prescription drugs, clothing/shoes, temporary housing, psychological support, legal assistance, transport services, social support (for vulnerable categories), information support (about evacuation, aid points, etc.). Classification was performed using the Random Forest

algorithm, trained on a synthetic dataset that simulates user profiles and queries, with an average accuracy of 92% for cross-validation. Validation of the results of designed models and implemented solutions, individual services and modules was carried out by testing on synthetic data (based on real-world scenarios) and pilot implementation of individual modules in a test environment. The plan for testing ICH modules on synthetic data provided for testing the functionality, performance, error resistance, and integration interaction of the main implemented modules of the system.

The purpose of testing was to check the operability and efficiency of ICH modules in conditions close to real scenarios of humanitarian aid, using synthetically generated data (artificial, but as close as possible to real ones, obtained based on the analysis of reports of existing systems). The objects of testing were selected modules for processing requests, providing assistance (forming routes), a module for classifying needs (intelligent part, ML, AI), a module for managing users and interfaces, a module for storing and searching data, and APIs for integrating with external systems. Testing of ICH modules was carried out by the author's research group, which included 3 specialists in information systems design and system testing, with relevant experience in the field of humanitarian logistics and software development. Since all tests were conducted on synthetically generated data that did not contain personal data or information that may be related to a real person, the participation of outsiders was not foreseen, and therefore, there was no need for additional coordination with the ethics committees.

Data for testing. External data sources – 5 simulated REST APIs with JSON/XML responses that mimic UN and WHO systems and volunteer platforms. Input synthetic data was generated automatically based on real data templates. An example of the generated data is shown in Table 2.

Table 2. Example of synthetic data

Data	Example (synthetic)
Request for help	ID: RQ-000124, Category: Medicines, Oblast: Chernivtsi, Priority: 3
Inventory data	ID: ST-009821, Type: Bottled water, Quantity: 2,000 units.
Delivery points	ID: LOC-412, GPS: 48.2904, 25.9358
Users	ID: US-098, Role: Logistician, Status: Active
Incident data	ID: IN-243, Type: Blocked Bridge, Status: In progress

Source: developed by the author

All experiments followed the basic principles of ethical handling of data and digital models, according to the guidelines of the Association for Computing Machinery (ACM) (2018). The collection and processing of medical and social data was in line with the Declaration of Helsinki (2024) medical research ethics standards. Research data processing, obtaining results, and conducting testing were based on the ethical principles of the EU's Seventh Framework Programme (European Commission, 2013), and the Organisation for Economic Co-operation and Development (OECD) (2021) guidelines were followed when using AI.

Data sources. Data from ReliefWeb (n.d.), official statistical reports of United Nations High Commissioner for Refugees (UNHCR) (2024a) and United Nations High Commissioner for Refugees (UNHCR) (2024b), and anonymous requests from call centres of international organisations (for example, SaveUA, n.d.) were used to analyse the needs of victims. The data included text descriptions of needs, geolocation, time markers, and socio-demographic characteristics. To form general concepts about incoming and outgoing data flows, the authors used report data, statistical and analytical data (United Nations Ukraine, 2025), the results of the analysis of scientific publications on information systems in the humanitarian sphere. ERWIN structural design tools were used to design and visualise models of the information system and its modules. Architecture modelling was performed in Bizagi Modeler using BPMN2 notation to describe business processes and UML to design system modules. The following tools were used for data processing: Elasticsearch – indexing and searching in large data sets; Python (spaCy libraries, Transformers) – NLP-analysis of text queries; Scikit-learn – clustering (k-means) for grouping needs. When forming modules for the distribution of

humanitarian aid, an analysis of clustering algorithms was carried out – k-means was selected for grouping requests of victims, recommendation systems, namely, a hybrid system with active filtering of the type of assistance and less active collaborative filtering, which is limited by the nature and size of assistance for selecting relevant assistance and personalising assistance, to eliminate problems of poor-quality input data – using data cleaning algorithms based on the standard query language SQL for performing data transformations and using linguistic extensions for specific applications, in particular, user-defined functions supported in SQL. The study also followed the Regulation (EU) of the European Parliament and of the Council No. 2016/679 (2016). Validation at this stage was performed using synthetic data.

RESULTS AND DISCUSSION

At the stage of conceptual design, the architecture of the ICH information system was proposed. Users of the system are: Donors – individuals, legal entities, charities, public organisations, foundations, social services that provide various types of assistance (clothing, social, medical, financial, legal, housing, employment, etc.); Victims – individuals who need various types of assistance as a result of war, disasters, natural disasters; Volunteers – individuals or organisations that physically implement the processes of providing assistance; Supervisors – organisations that have the right to monitor the operation of the system; Administrators – persons who ensure the operation and development of the system, have permits for certain services of the system. The ICH of humanitarian aid is a web resource where users of the system can implement the corresponding functions depending on their type. The system consists of two modules Supersystem and Modules (Fig. 1).

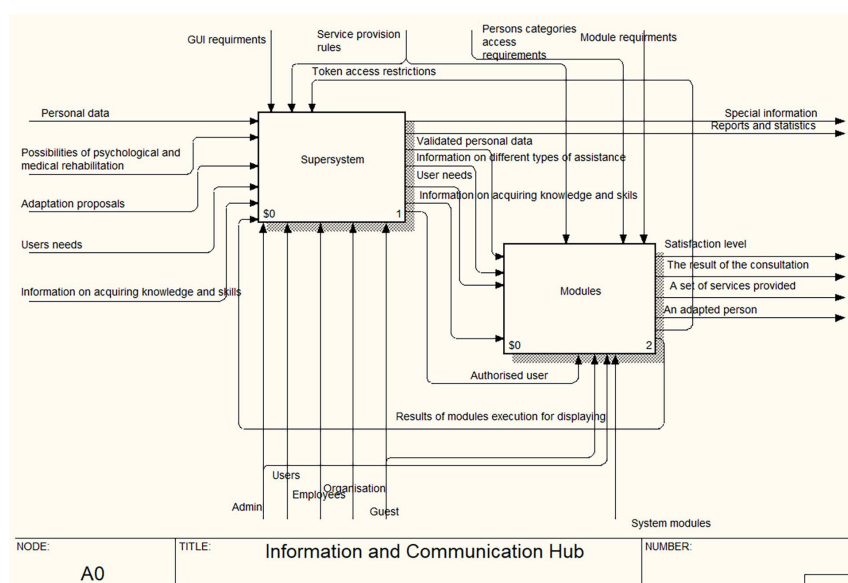


Figure 1. General structure of ICH systems

Source: developed by the author

The main system of a web application is a Super-system with a user interface. The main page for an organisation's role contains an interface for adding information about the services provided by this organisation, and lists of system users, depending on the roles that have registered in the system. The user's home page contains information about the application's capabilities and a system for adding or removing services. This Supersystem should combine all systems, including an authentication and authorisation system, an information module, a communication module, a user account, a help package generation module, and an analysis and statistics module.

Using the graphical user interface (GUI), users can register and log in to the web application, submit their candidacy for a certain type of assistance, and communicate with employees of organisations that provide this assistance. Employees will be able to view lists of people who need help, view statistics on assistance provided, and, in turn, communicate with users. Data enters the system from four sources: manual input (simulation of real requests from different types of users, automatic retrieval via REST API from external sources or their

simulation models, this work used real APIs of existing systems), through import from prepared CSV/JSON files containing arrays of requests, infrastructure objects, routes, etc., and through generation by an internal simulator written in Python using the Faker, Pandas, and NumPy libraries.

The database model was divided into several logical blocks: requests – a query table with fields such as ID, aid type, coordinates, timestamp, processing status, transport connection; resources – types of available resources (medicines, food, water, etc.) with quantitative and logistical attributes; transport_units – vehicles with characteristics (type, load capacity, current position, availability); dispatch_centers – aid delivery nodes and their coordinates, resources, and throughput; users: system user model with roles (operator, analyst, logistician), permissions, and interaction history. The Modules ICH system, in turn, consists of five working modules (Fig. 2) – a module for accounting for persons by characteristics, an information module, a communication module, a user account, a module for forming a help package, and an analysis and statistics module, where each of the modules is a separate service with its own database.

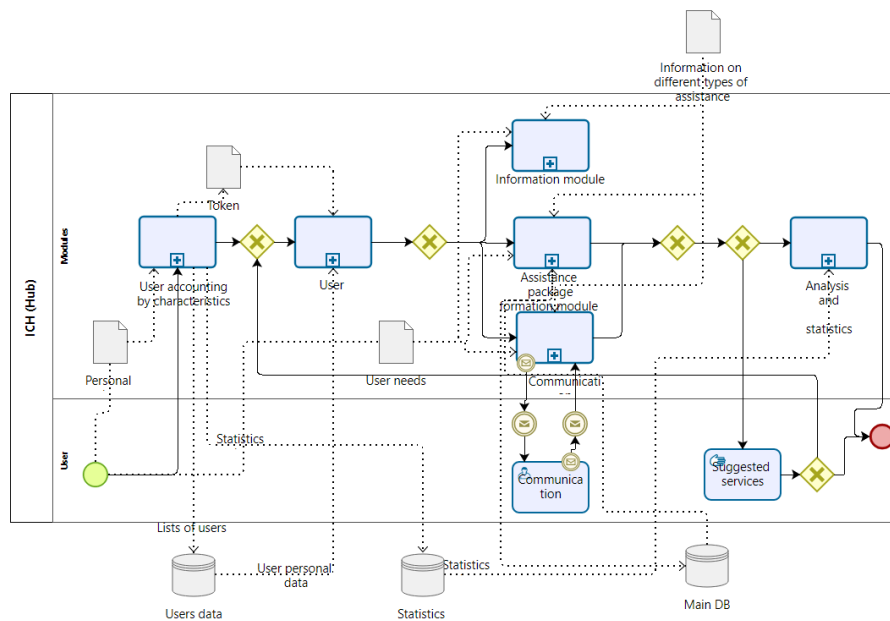


Figure 2. Structure of services of the working ICH subsystem

Source: developed by the author

The authentication and authorisation subsystem, which includes the authorisation/registration and login modules by various users (the authentication module), and the actual granting of permissions until the relevant data is received. This subsystem is responsible for secure registration and login of users to the system. Using the authorisation module, which is part of the subsystem, employees of the system receive their own set of permissions, which forms the advanced functionality of the web application, which will allow dividing it into roles in the service. The subsystem involves the use of

a number of algorithms. For example, an algorithm for hashing system user passwords to ensure the security of user data. Algorithms for effective data search in lists by specific keywords, algorithms for sorting and filtering lists by various criteria are also necessary. The authorisation and authentication module in the system of different types of users with different roles and corresponding functions is responsible for secure registration and login of users and employees of the web service to the system. Using this authorisation module, system employees will receive advanced functionality

of the web application, which will allow them to divide into roles in the presented service.

Functional requirements for the Requirement (REQ) system.

REQ-1.1: Ensuring the user's registration in the information system, which includes checking the correctness of the entered password, transferring the user's personal data, such as username, password, and full name, from the client part of the application to the database.

REQ-1.2: Providing user authentication to the information system, which includes comparing the username and password entered by the user with the data contained in the database, and authorisation to continue the user's work with the system.

REQ-1.3: Providing reliable integration with government services.

REQ-1.4: Separating access to staff (administrators) functionality and user roles.

REQ-1.5: Providing interaction with the list generation submodule by transmitting – for further processing of the list and generating a token – selected services during user registration.

Module for accounting for individuals by attributes (designed for entering and categorising information about service users).

The submodule for creating lists of people by attributes is designed to create a permission system for each of the user types (Fig. 3).

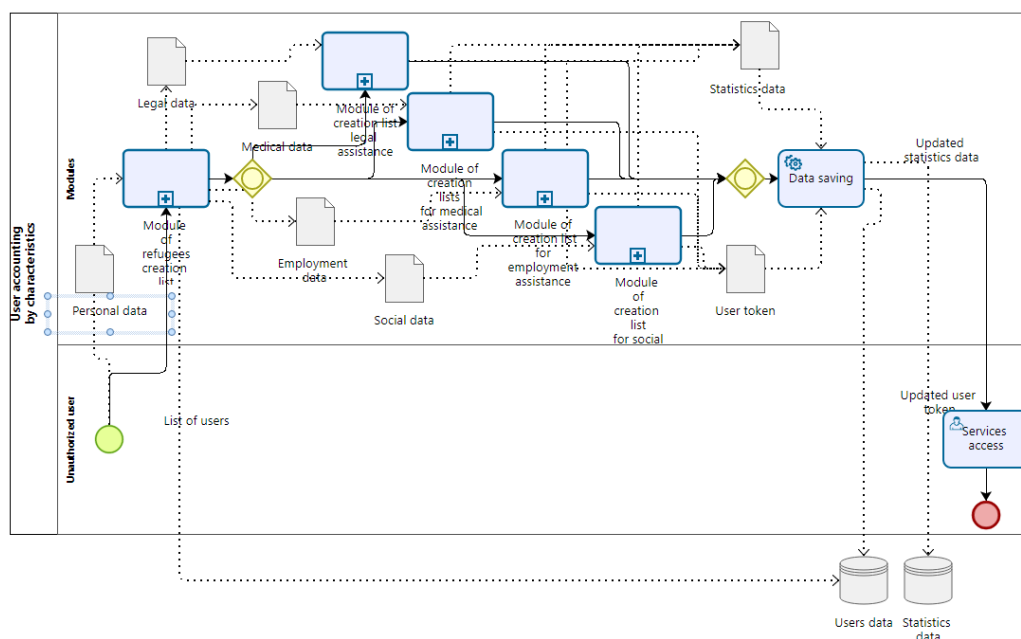


Figure 3. Model for creating user lists based on attributes

Source: developed by the author

The user list generation module suggests using Elasticsearch algorithms (indexing and search) for quick filtering by criteria. The k-means algorithm was used to implement clustering tasks, namely automatic grouping of people by similar needs, type of care, age, condition, gender, and level of need (immediate, high, medium, and low). Ontology-based rules were also used (for structured data: classification by keywords (“housing”, “medicine”, “therapy”) and metadata (location, time)). Metadata (request time, location, need type) was stored in the PostgreSQL cloud database using AES-256 encryption. Access was restricted to Role-Based Access Control (RBAC) roles (for example, volunteers only saw anonymous data). In the future, personal data will be anonymised before processing; voice recordings are planned to be stored in encrypted form.

Token creation submodule-responsible for creating a token of user access to selected services, classifying users according to certain criteria, and providing statistical

data to the statistics and analytics module. It forms a token of access to the services of the system through the user's account, which is the main user interface, based on the user's role in the system if the user is a volunteer, donor, administrator or controller, or depending on the need for the type of assistance (there may be several, in this case, victims are added by the system to several lists and options for working with each of the types of assistance are added to the user's account, on the principle of a single window of access to the system. The information module is designed for a reference system that provides users with reference information for obtaining certain services.

The communication module is designed to provide interaction between users who need help and employees who provide this assistance. The communication module contains the following submodules: message exchange, settings, moderation, saving, editing the forum structure, and a data submodule. Several functional

requirements for individual submodules of the communication module were specified. The messaging submodule allows users to share messages in shared forum topics and privately, and edit or delete them. Feature priority: high. An authorised user can send messages in general and private topics, in person. After authorisation, the main page opens, which allows visitors to navigate to the forum page with sections and rules added by the administrator. By selecting the section and topic of interest, the user can participate in surveys, exchange messages, edit and delete them, mark other users, and add attachments of certain formats. The moderator can also create surveys and follow their results directly in the chat.

REQ-2.1: Ability to select a section and topic manually.

REQ-2.2: Ability to select a section and topic using search.

REQ-2.3: Ability to send messages in general discussions.

REQ-2.4: Ability to send private messages after switching to another user's profile.

REQ-2.5: Ability to mark users who are in discussion in messages.

REQ-2.6: Ability to attach files of the specified formats: .doc, .pdf, .png, .jpg, .txt.

REQ-2.7: Ability to participate in general surveys.

REQ-2.8: Ability of the moderator to edit, move, and delete other people's posts in a specific section or topic.

REQ-2.9: Ability of the moderator can create a survey in a specific section or topic.

The settings submodule allows the user to change their personal preferences for receiving notifications. Feature priority: low. The system user can go to their forum profile and edit their notification preferences. Notifications on the site and e-mail messages will be available.

REQ-2.1: Displaying the settings page with the following blocks: notifications about notifications in marked topics, notifications about notifications in topics

in which the user participates, notifications when marked, notifications about receiving private messages, notifications about moderator actions, newsletters with news content, newsletters with new and popular discussions, newsletters with surveys or their results.

REQ-2.2: Ability to change the place where notifications are received (directly by the website or email) for each item.

REQ-2.3: Ability to disable or enable notifications of a specific category, or allow only priority notifications.

The moderation submodule is primarily intended for implementing the main functionality for the moderator. Implements the ability to identify users for further distribution into lists based on the filters entered in the system. For the administrator, the module provides the ability to grant permissions. In addition, the module contains AI, which will not only simplify the moderator's work, but will also be useful for users. AI generates lists of users by various criteria, group user requests, form answers to frequently asked questions, find and block bots, manage spam filters, determine, according to the locations of requests, forecasting from the development of events to the development of mailings with information, warnings, and algorithms of actions.

The data submodule provides direct access to the administrator to perform actions such as backups and ensures that the necessary data is saved, updated, and deleted. The submodule for editing the structure of the communication forum allows the administrator not only to create new sections, topics and discussions, but also to add general advertisements, add information to the list of frequent questions, modify the rules for creating messages, add news or informational videos. The download submodule allows users to save chats, sections, pin messages, download and save files of available formats (.doc, .png, .pdf, .jpg, .txt). The user account module is designed to display user information, change personal information, and register or delete services provided (Fig. 4).

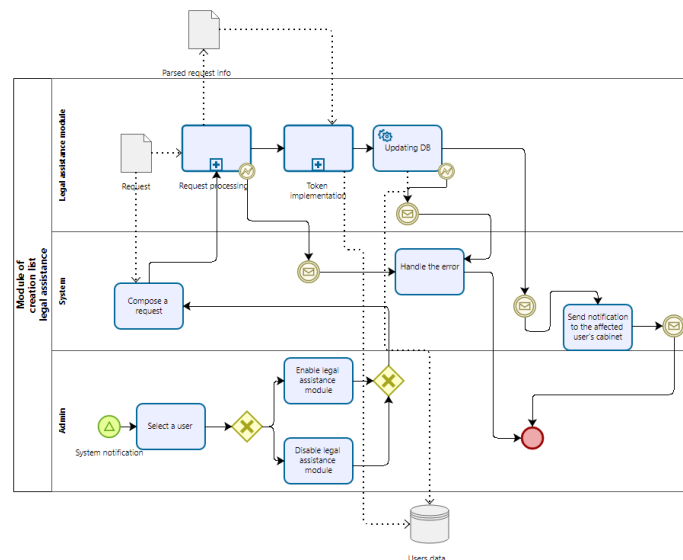


Figure 4. Module for creating a user account with the corresponding functions

Source: developed by the author

The user cabinet module is the main user interface through which users of different types can approach all the functions available to this type of user in the corresponding services of the system, that is, each cabinet is formed by the system depending on the role, speciality of the volunteer, or donor, requests of victims, that is, it is created on the principle of adaptability and can change depending on the solution of problems, or the emergence of new needs. User dashboard module – provides user access to the list generation system – allows choosing whether to activate or deactivate selected services, display user information, and change

personal information. In this module, algorithms were used: to create secure authorisation JWT/OAuth 2.0, to differentiate access rights and generate accesses for each of the RBAC roles, to ensure data protection in the system, they were encrypted using the AES-256 encryption algorithm. The user dashboard module consists of a list generation submodule, a help request processing submodule, a help status submodule, and a GUI submodule based on roles, rights, access, needs, and other attributes that are defined in filtering mechanisms. The structure of the submodule for generating user lists by feature is shown in Figure 5.

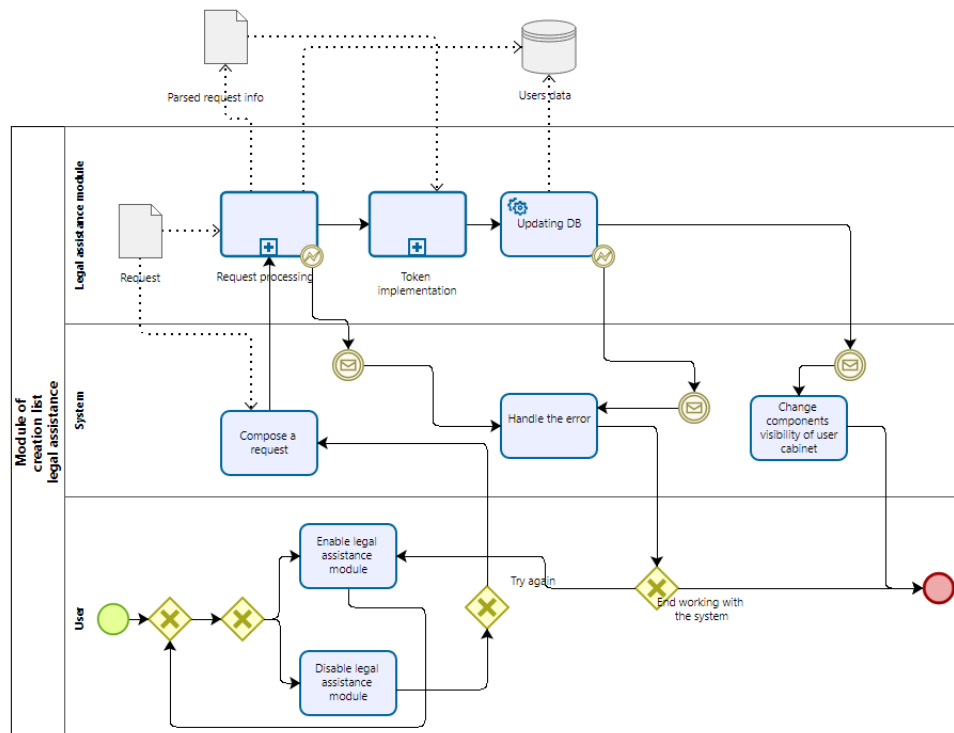


Figure 5. Module for creating user lists with different types of permissions

Source: developed by the author

A registered and authorised user must have access to the systematised information posted on the site, be able to send requests for assistance with clarification of the problem, communicate and share experience in discussions and privately, and search the site. In addition, after receiving the system's recommendations on the selected types of assistance, a person should be able to consult a system specialist if necessary, get general and personalised information. The user must have the authority to change the information provided about them in the user cabinet.

REQ to the module.

REQ-3.1.1: Providing interaction with the list generation submodule by transmitting – for further processing of the list and generating a token-selected services in the user's dashboard.

REQ-3.1.2: Providing the ability to edit personal information of an authorised user, such as full name, place of residence.

REQ-3.1.3: Providing viewing of the authorised user's personal information, such as: full name, place of residence, list of selected necessary services, and history of their receipt. In the merchant profile, according to the role and open permissions, the user gets access to functions, which correspond to the user's role (Fig. 5).

The help package formation module is a module that will generate the help itself and information about it. The help submodule provides the ability to view statistics of requests for help from a specific user. This submodule provides information about user-generated queries and their statuses and results.

REQ-4.1.1: System should display a list of generated requests for user assistance based on the time period entered.

REQ-4.1.2: System should reflect the number of successfully received requests relative to the total number of requests.

REQ-4.1.3: System should display a list of requests based on the categories selected by the user (medical, material, psychological assistance, etc.).

REQ-4.1.4: Users authorised as employees should be able to view a list of requests for assistance they have processed.

REQ-4.1.5: Users authorised as employees should be able to view a list of help requests created by another specified user.

The analysis and statistics module helps organisations to track statistics on the volume and scope of assistance provided. The analysis and statistics module consists of the following submodules: statistics and analytics on the work of individual modules, statistics on an individual user, a module for collecting software statistics, a system for generating financial reports, a system for analysing development and problems. Submodule of statistics for an individual user provides the ability to view statistics of requests for help from a specific user. This submodule provides information about user-generated queries and their statuses and results. Functional requirements for the statistics submodule for an individual user.

REQ-5.3.1: System should display a list of generated requests for user assistance based on the time period entered.

REQ-5.3.2: System should reflect the number of successfully received requests relative to the total number of requests.

REQ-5.3.3: System should display a list of requests based on the categories selected by the user (medical, material, psychological assistance, etc.).

REQ-5.3.4: Users authorised as employees should be able to view a list of requests for assistance they have processed.

REQ-5.3.5: Users authorised as employees should be able to view a list of help requests created by another specified user.

Submodule of social information. Feature priority: low. The system calculates statistics related to the user's social communication, mainly based on data about their activity in the communication module.

REQ-5.4.1: System should reflect the number of responses provided on forums.

REQ-5.4.2: System should display the number of closed and open user discussions.

REQ-5.4.3: System should display the number of votes on other users' posts.

REQ-5.4.4: System should calculate the average and total user rating.

REQ-5.4.5: System must calculate the user's karma coefficient using a formula using data on the number of likes and rating of this user.

Implementation of ICH services. In the course of the study, separate ICH services were implemented. Service for creating a request for medical care. An authorised and authenticated user is redirected to the medical care submodule after filling out the medical care form in the merchant profile. Figure 6 shows the main page of the medical care service.

The user then submits a request for assistance, and the AI system processes the data received to determine the assistance package, find a medical facility and doctor, and redirect the user's request. Figure 7 shows the interface for creating a request for help.

The status and type of assistance provided is generated in the statistics module in dashboards for system employees. The panels contain different colours for information about requests request type, status is indicated by colours: gray – missing, red – not considered, yellow – in progress, and green – closed. Next to each user, there is an Update Info button that allows the employee to view detailed information and manually enter information about the status of solving the problem. Panels are generated by using filtering at the employee's request. The communication service is implemented as an eSupport system (Fig. 8).

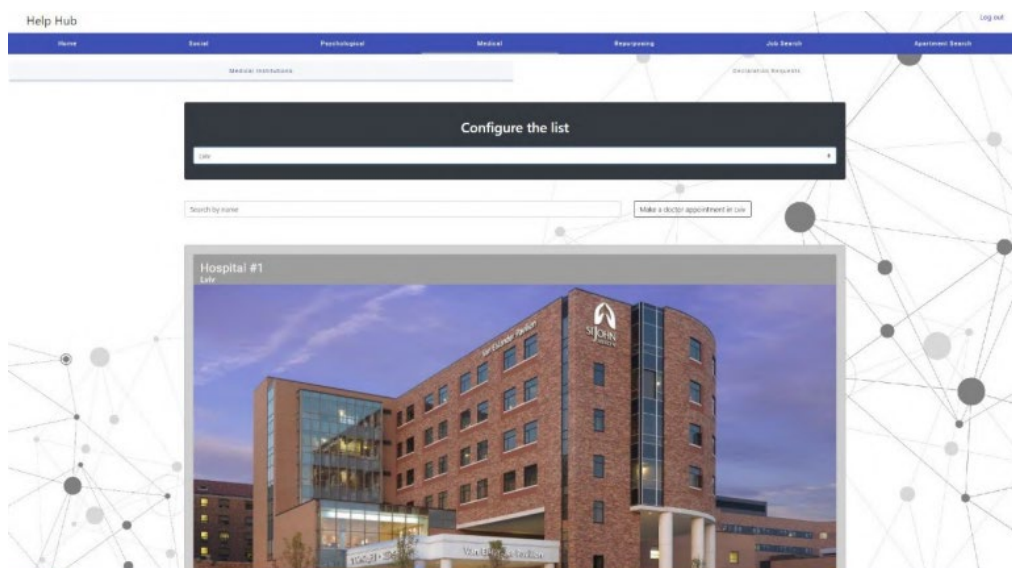


Figure 6. Home screen of the medical care service in the help package generation module

Source: developed by the author



Figure 7. Interface of the medical care request form

Source: developed by the author

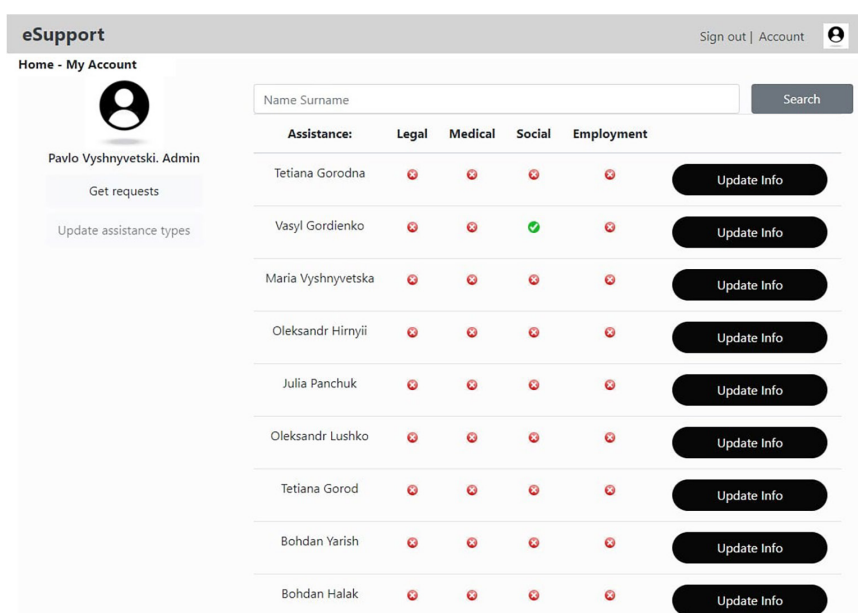


Figure 8. Admin tracking panel for information about active requests for help

Source: developed by the author

To check the functional correctness and performance of key ICH modules, testing was conducted using synthetic (artificially generated) data. The main purpose of testing was: to evaluate the system's ability to handle large volumes of heterogeneous requests

in real time; to test algorithms for classifying requests for assistance; to evaluate the efficiency of routing and prioritisation of cargo; to check the stability of interfaces under load. The results based on synthetic data are shown in Table 3.

Table 3. Results of testing on synthetic data

Metric	Result
Average request processing time	437 ms
Percentage of successfully classified queries	91.3%
Average deviation of the route from the standard	7.8%
Successful import of data from external APIs	98.5%
User interface response time	212 ms

Source: developed by the author

In the course of the study, 5 ICH modules were implemented, including: Person accounting module: automatic classification of 12 types of needs (92% accuracy)

based on the Random Forest algorithm. Communication module: NLP analysis of 10,000 requests/day (average response time – 2 minutes). To assess the

functional and productive capacity of the developed ICH modules, unit testing was conducted using synthetic data. Testing covered the following areas. The correct operation of the humanitarian cargo routing module was checked, in particular, the construction of optimal routes considering changes in the availability of logistics hubs. The classification module for assistance requests was tested on 1,000 synthetic requests with the specified categories (medical care, food, housing, etc.), and the classification accuracy was 94.5%.

The prioritisation management module was tested on conflict scenarios using adaptive rules, and the results showed stable sorting based on risk factors. Scalability testing was also performed, during which the request processing time was measured when the load increased. As can be seen from Figure 9, system performance remains at an acceptable level up to a load of 300-400 simultaneous requests. Then there is a linear increase in response time, which is expected for a microservice architecture without horizontal scaling.

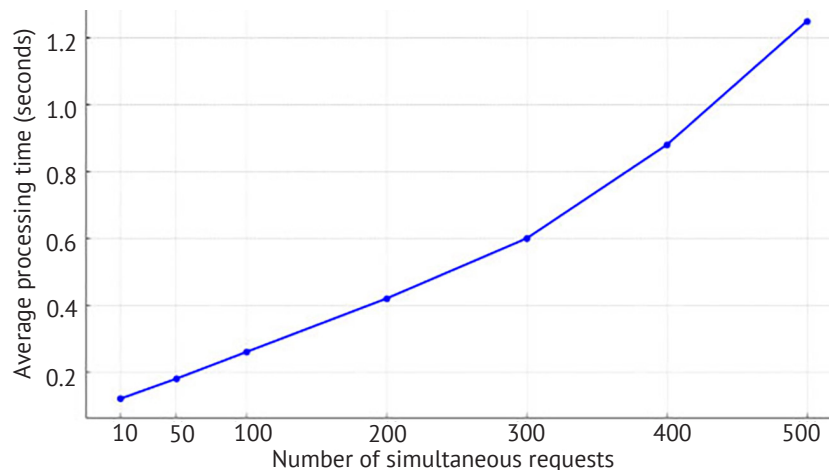


Figure 9. ICH performance: request processing time depending on the load

Source: developed by the author

The results show that most ICH modules meet the specified functional and non-functional requirements. A slight decrease in the accuracy of needs classification can be improved by further training the model on extended text query enclosures. High UI stability and routing efficiency indicate that the system is ready to deploy in a dynamic environment. The testing confirmed the possibility of using various technologies to develop and gradually implement individual services into the system, which confirmed the high scalability of the system. The introduction of single window principles increases the controllability of all processes in the humanitarian system. At the current stage, there was no functional testing of the real-time voice message analysis module due to the lack of integration with call centres. The study presented the concept of ICH with a microservice architecture, which provides integration of heterogeneous data sources, automated request processing, high response speed (average processing time – 2 min), accuracy of classification of needs (up to 92%) and efficient processing of an assistance order from the request to the development of a targeted assistance package (reduction of package formation time by 18%). The system supports ReliefWeb/UNHCR API exchange standards and integration with Sahana (n.d.) or “eDopomoga” platforms (Diia, n.d.). Validation was performed on synthetic data using open APIs. These results are correlated in comparison with other studies, with the analysis highlighting additional aspects.

Data exchange was studied by A. Gunes *et al.* (2020). The researchers have identified low interoperability, delays, and various data formats as critical issues. The paper noted that data fragmentation between different humanitarian organisations often leads to duplication of efforts, uneven allocation of resources, and gaps in aid delivery. The paper emphasised that the lack of standardised information exchange protocols and incompatibility of technological solutions create significant barriers to rapid response. The microservice architecture of the developed ICH model with automatic routing solves these problems more efficiently. This approach is consistent with trends in research on microservice architectures. G. Kousiouris *et al.* (2019) showed the effectiveness of microservice integration of IoT platforms, semantic and AI services in supply chains, which confirmed the feasibility of using this paradigm for complex and dynamic systems. In addition, M. Waseem *et al.* (2020) and M. Waseem *et al.* (2021) addressed the issues of design, testing, and monitoring of microservice systems in DevOps environments, highlighting the challenges of scalability and fault tolerance. The proposed ICH architecture is consistent with these conclusions, in particular, due to the modularity and independence of services. The issue of data management in microservice architects was considered by R. Laigner *et al.* (2021), who described the difficulties of storing data in a microservice architecture. The ICH architecture considers these aspects through specialised

services with isolated repositories, which is consistent with the authors' recommendations.

The issue of involving the private sector in humanitarian operations deserves special attention. The study by A. Cozzolino (2021) showed that digital platforms can improve the efficiency of business interaction with humanitarian organisations, creating conditions for sustainable partnerships and strengthening the resource base. This reinforces the argument about the need for open integration of ICH with existing commercial services and information systems. In the design of ICH, blockchain technology was not used at this stage, but based on the authors' conclusion that the introduction of blockchain technologies in aid distribution systems ensures a constant record of all transactions, increases accountability, reduces the risk of corruption and fraud, and allows making the whole process more efficient and fair, it is planned to introduce blockchain technology in the module for forming aid packages.

Another important issue being studied in relation to humanitarian information systems is data analytics. Aspen Institute Kyiv (2023) and DIGID Consortium (2023) described research on information interaction in humanitarian systems. The study by Aspen Institute Kyiv (2023), which focuses on improving the coordination of various types of humanitarian initiatives, also analysed the role of volunteer initiatives and their information systems in coordinating assistance within the country. The study highlighted the unique challenges associated with infrastructure instability, large numbers of internally displaced persons, and the need to quickly adapt to a dynamic front. They pointed out the importance of developing systems that can integrate with local databases and ensure rapid exchange of information with international partners. The analysis cited in the study highlighted the lack of clear standards for data exchange. On the other hand, the study by DIGID Consortium (2023) highlighted the need for interoperability between organisations. The ICH system implements API protocols according to ReliefWeb/UNHCR standards, which eliminates these gaps.

R. Mohammed Zain *et al.* (2023), studied the coordination of humanitarian logistics in cities, emphasising the need for transparency and standardised data flows. The study showed the positive impact of digital technologies on the efficiency, sustainability, and transparency of logistics. However, the topic of information flows in crisis communications remains relevant. D. Mitcham *et al.* (2021) proposed the concept of "communication hub framework", which describes the use of social networks for rapid dissemination of information in the context of local disasters. Although ICH does not implement direct mechanisms for integration with social networks, the data collection module can be expanded in this direction, which will strengthen the early information component and increase the involvement of local communities.

ICH indicators – reducing processing time and increasing transparency – confirm their conclusions, supporting the relevance of the functional requirements of ICH, coincide with the results and conclusions of the authors, but the use of ML classification and API integration in ICH increases the quality and transparency of the results achieved. F. Adediran *et al.* (2024) analysed the use of blockchain for supply chain transparency. The system development plan also provided for the use of blockchain, which resonates with the priorities of ensuring data trust. However, it should be noted that their research did not cover algorithmic classification or ML query processing. S. Kyrylashchuk *et al.* (2024) considered the use of AI in solving logistics problems of various types. This included forecasting needs (based on historical data and the current situation), classifying and prioritising aid requests, and optimising humanitarian delivery routes. The use of AI contributes to more accurate and faster decision-making. The paper considered the use of AI models in modules for the development of assistance packages, according to the received requests and analysis of assistance tools, in the formulation of logistics routes and routes in logistics, in forecasting and assessing the level of needs to overcome the consequences of a particular disaster. M.F. Carnero Quispe *et al.* (2024) reviewed multi-criteria prioritisation models in humanitarian logistics. The researchers identified four criteria: efficiency, effectiveness, fairness, and sustainability. The routing and ML classification algorithms used in ICH modules also integrate these criteria. Thus, the conducted comparative analysis showed that the proposed ICH concept is consistent with key trends in humanitarian information logistics: data integration, transparency, scalability, and efficiency. The developed system significantly expanded existing approaches, providing automatic ML classification, algorithmic optimisation of activities, and specific performance metrics.

CONCLUSIONS

The study proposed the concept of ICH for humanitarian coordination, which addresses the key problems of data fragmentation, low scalability, and lack of integration between stakeholders. An architectural approach to the implementation of systems of this type was proposed, namely, the use of microservice architecture for humanitarian systems, which provides flexibility and adaptability, was substantiated. Due to the analysis of the system functions and the possibilities of their effective implementation, effective algorithms were determined for each of the system modules for dynamic crisis conditions, some of which were implemented and tested on synthetic data.

The implementation of data exchange standards was proposed and implemented, in particular, based on the ReliefWeb and UNHCR API standards, which increased the transparency and controllability of data and processes in the system. The system reduced the

response time to requests (average processing time – 2 minutes) due to automation, improved the accuracy of needs classification (up to 92%) and resource allocation efficiency (reducing delivery time by 35%). The system is a hub and can be integrated with both its own services and external systems via the API, and can also be integrated with existing platforms such as Sahana, “eDopomoga”. The involvement of AI has become not only a means of increasing productivity, but also a key tool for the flexible response of the system in the face of uncertainty, resource scarcity, and dynamically changing scenarios of humanitarian assistance. Validation was carried out on synthetic data; pilot implementation was required in times of crisis. Future areas of development of the system will include GIS integration for visualisation, the use of blockchain for tracking

supply chains, and the extension of NLP for multimodal data (text + voice). With further development of the system, it is planned to use algorithms for voice recognition, status detection, etc. It is also planned, but not implemented, to use additional ontologies and graph DB to analyse victim relationships (for example, family relationships, social connections, temporary location).

ACKNOWLEDGEMENTS

None.

FUNDING

None.

CONFLICT OF INTEREST

None.

REFERENCES

- [1] Adediran, F.E., Okunade, B.A., Daraojimba, R.E., Adewusi, O.E., Odulaja, B.A., & Igbokwe, J.C. (2024). Blockchain for social good: A review of applications in humanitarian aid and social initiatives. *International Journal of Science and Research Archive*, 11(1), 1203-1216. doi: 10.30574/ijrsra.2024.11.1.0184.
- [2] Aspen Institute Kyiv. (2023). *How to improve the coordination of stakeholders of humanitarian aid*. Retrieved from https://lb.ua/blog/aspen_institute_kyiv/595074_yak_pokrashchiti_koordinatsiyu.html.
- [3] Association for Computing Machinery (ACM). (2018). *ACM code of ethics and professional conduct*. Retrieved from <https://www.acm.org/code-of-ethics>.
- [4] Bissoft. (2025). *Complex of integrated web services for wartime conditions*. Retrieved from <https://bissoft.org/ua/kompleks-vzayemointehrovanykh-veb-servisiv-dlya-umov-voyennoho-chasu>.
- [5] Carnero Quispe, M.F., Couto, A.S., de Brito, J.I., Cunha, L.R.A., Siqueira, R.M., & Yoshizaki, H.T.Y. (2024). Humanitarian logistics prioritization models: A systematic literature review. *Logistics*, 8(2), article number 60. doi: 10.3390/logistics8020060.
- [6] Cozzolino, A. (2021). Platforms enhancing the engagement of the private sector in humanitarian relief operations. *Sustainability*, 13(6), article number 3024. doi: 10.3390/su13063024.
- [7] Declaration of Helsinki. (2024). Retrieved from <https://www.wma.net/policies-post/wma-declaration-of-helsinki/>.
- [8] DIGID Consortium. (2023). *The necessary interoperability of systems between organisations*. Retrieved from <https://interoperability.ifrc.org/2023/05/23/the-necessary-interoperability-of-systems-between-organisations/>.
- [9] Diia. (n.d.). Retrieved from <https://diia.gov.ua/>.
- [10] European Commission. (2013). *Ethics for researchers: Facilitating research excellence in FP7. Directorate-general for research and innovation*. Luxembourg: Publications Office of the European Union. doi: 10.2777/7491.
- [11] Google Crisis Response. (n.d.). Retrieved from <https://crisisresponse.google/>.
- [12] Gunes, A., Ozer, M., & Kirca, M. (2020). Information sharing challenges in humanitarian supply chains: A systematic review. *International Journal of Disaster Risk Reduction*, 48, article number 101569. doi: 10.1016/j.ijdr.2020.101569.
- [13] Ilyina, I.V., Tokarev, V.V., Yakovlev, A.V., & Shevchenko, I.I. (2024). Using a decision support system for organizing humanitarian logistics. *Control, Navigation and Communication Systems. Academic Journal*, 1, 88-95. doi: 10.26906/SUNZ.2024.1.088.
- [14] Kankanamge, N., Yigitcanlar, T., Goonetilleke, A., & Kamruzzaman, M. (2019). Can volunteer crowdsourcing reduce disaster risk? A systematic review of the literature. *International Journal of Disaster Risk Reduction*, 35, article number 101097. doi: 10.1016/j.ijdr.2019.101097.
- [15] Kondraganti, A. (2021). Big data analytics in humanitarian and disaster operations: A systematic review. *ArXiv*. doi: 10.48550/arXiv.2108.09800.
- [16] Kousiouris, G., Tsarsitalidis, S., Psomakelis, E., Koloniaris, S., Bardaki, C., Tserpes, K., Kyriazis, D., & Anagnostopoulos, D. (2019). A microservice-based framework for integrating IoT management platforms, semantic and AI services for supply chain management. *ICT Express*, 5(2), 141-145. doi: 10.1016/j.icte.2019.04.002.
- [17] Kyrylashchuk, S., Horodetska, O., Voitsekhovska, E., & Zakharchenko, S. (2024). Order forecasting system for vehicles based on previous statistics requests. In S. Babichev & V. Lytvynenko (Eds.), *Lecture notes in data engineering, computational intelligence, and decision-making* (pp. 116-131). Cham: Springer. doi: 10.1007/978-3-031-70959-3_6.

- [18] Laigner, R., Zhou, Y., Salles, M.A.V., Liu, Y., & Kalinowski, M. (2021). Data management in microservices: State of the practice, challenges, and research directions. *ArXiv*. doi: [10.48550/arXiv.2103.00170](https://doi.org/10.48550/arXiv.2103.00170).
- [19] Matekenya, T., & Ruhode, E. (2021). Towards an integrated knowledge management and ICT framework for improving disaster response in a developing country context. *ArXiv*. doi: [10.48550/arXiv.2108.09813](https://doi.org/10.48550/arXiv.2108.09813).
- [20] Mitcham, D., Taylor, M., & Harris, C. (2021). Utilizing social media for information dispersal during local disasters: The communication hub framework for local emergency management. *International Journal of Environmental Research and Public Health*, 18(20), article number 10784. doi: [10.3390/ijerph182010784](https://doi.org/10.3390/ijerph182010784).
- [21] Mohammed Zain, R., Mohd Zahari, H., & Mohd Zainol, N.A. (2023). Inter-agency information sharing coordination on humanitarian logistics support for urban disaster management in Kuala Lumpur. *Frontiers in Sustainable Cities*, 5, article number 1149454. doi: [10.3389/frsc.2023.1149454](https://doi.org/10.3389/frsc.2023.1149454).
- [22] Nezdoiminoga, O.Ye., & Pryidak, T.B. (2024). [Digitalization of humanitarian aid in Ukraine under martial law](#). In V.V. Khrapkina & K.V. Pichyk (Eds.), *Transformation of innovation development management practices of socio-economic systems* (pp. 560-567). Kyiv: Kyiv-Mohyla Academy Publishing House.
- [23] Nozdrina, L., & Falat, M. (2020). Features of volunteer IT-projects in the social economy. *Socio-Economic Relations in the Digital Society*, 3(39), 112-122. doi: [10.18371/2221-755X3\(39\)2020225607](https://doi.org/10.18371/2221-755X3(39)2020225607).
- [24] Organisation for Economic Co-operation and Development (OECD). (2021). *Recommendation of the council on the ethical guidelines for trustworthy artificial intelligence (AI)*. Retrieved from <https://legalinstruments.oecd.org/en/instruments/OECD-LEGAL-0449>.
- [25] Petrovskaya, K.V., Turgeneva, A.A., & Melnychuk, A.B. (2024). The role of volunteer initiatives and NGOs in the post-war restoration of Ukraine. In *Post-war reconstruction of Ukraine: Legal, political and economic challenges in historical perspective* (pp. 652-676). Latvia: Baltija Publishing. doi: [10.30525/978-9934-26-559-4-32](https://doi.org/10.30525/978-9934-26-559-4-32).
- [26] Regulation (EU) of the European Parliament and of the Council No. 2016/679 “On the Protection of Natural Persons with Regard to the Processing of Personal Data and on the Free Movement of Such Data, And Repealing Directive 95/46/EC (General Data Protection Regulation), GDPR”. (2016, April). Retrieved from <https://eur-lex.europa.eu/eli/reg/2016/679/oj/eng>.
- [27] Reliefweb. (n.d.). Retrieved from <https://reliefweb.int/>.
- [28] Sahana. (n.d.). Retrieved from <https://sahanafoundation.org>.
- [29] SaveUA. (n.d.). Retrieved from <https://saveua.in.ua/>.
- [30] Shevchenko, O. (2024). Problems of coordination of international humanitarian aid programs in Ukraine. *Humanitas*, 3, 178-186. doi: [10.32782/humanitas/2024.3.25](https://doi.org/10.32782/humanitas/2024.3.25).
- [31] Ukraine Support Hub. (n.d.). Retrieved from <https://huri.harvard.edu>.
- [32] United Nations High Commissioner for Refugees (UNHCR). (2024a). *Humanitarian needs and response plan Ukraine*. Retrieved from <https://www.unhcr.org/ua/media/ukraine-hnrrp-2024-humanitarian-needs-and-response-plan-en-20240110-pdf?>
- [33] United Nations High Commissioner for Refugees (UNHCR). (2024b). *Global report 2024*. Retrieved from <https://www.unhcr.org/publications/global-report>.
- [34] United Nations Ukraine. (2025). *United Nations in Ukraine annual results report 2024*. Retrieved from <https://ukraine.un.org/en/download/182985/293462>.
- [35] Waseem, M., Liang, P., & Shahin, M. (2020). A systematic mapping study on microservices architecture in DevOps. *ArXiv*. doi: [10.48550/arXiv.2008.07729](https://doi.org/10.48550/arXiv.2008.07729).
- [36] Waseem, M., Liang, P., Shahin, M., Di Salle, A., & Márquez, G. (2021). Design, monitoring, and testing of microservices systems: The practitioners' perspective. *ArXiv*. doi: [10.48550/arXiv.2008.07729](https://doi.org/10.48550/arXiv.2008.07729).

Інформаційно-комунікаційний хаб гуманітарної допомоги: системний аналіз, моделювання процесів та технологічні рішення

Фоменко Андрій

Кандидат педагогічних наук, доцент
Національний університет «Львівська політехніка»
79000, вул. Степана Бандери, 12, Львів, Україна
<https://orcid.org/0000-0001-7718-5419>

Анотація. Ескалація глобальних криз, зокрема війн та стихійних лих, значно підкреслює критичну потребу в підвищенні ефективності управління гуманітарною допомогою. Існуючі системи часто страждають від фрагментації даних, обмеженої масштабованості та недостатньої гнучкості інтеграції між ключовими стейкхолдерами. Метою даного дослідження була розробка комплексної концепції та архітектури інформаційно-комунікаційного хабу (ІКХ) для координації гуманітарної допомоги в умовах динамічних кризових ситуацій. Методологія дослідження включала системний аналіз існуючих рішень для виявлення їхніх обмежень, а також архітектурне моделювання із застосуванням нотацій BPMN2 та UML. Було також розроблено стратегії розвитку на основі мікросервісів та проведено тестування алгоритмів. Встановлено, що ключові проблеми координації гуманітарної допомоги можуть бути ефективно вирішені шляхом створення спеціалізованого ІКХ. Розроблено функціональні вимоги до системи, що охоплюють інтеграцію різноманітних джерел даних, автоматизовану обробку інформації та забезпечення зручних інтерфейсів для всіх учасників процесу. Запропоновано мікросервісну архітектуру ІКХ з модулями для управління запитами, обробки даних (зокрема з використанням методів машинного навчання для класифікації потреб) та гнучкими користувацькими інтерфейсами. Проаналізовано та запропоновано ефективні алгоритми для оптимізації ключових операційних процесів, таких як маршрутизація гуманітарних вантажів та пріоритезація запитів на допомогу. Практична цінність результатів дослідження полягає в можливості їх застосування фахівцями у сфері управління надзвичайними ситуаціями та міжнародними гуманітарними організаціями для скорочення часу реагування, підвищення прозорості розподілу ресурсів та покращення масштабованості проектів допомоги

Ключові слова: автоматизовані системи обробки даних; центр гуманітарних даних; мікросервісна архітектура; проектування інформаційної системи; адаптивні алгоритми; управління кризовими ситуаціями



Real-time drone type recognition using artificial intelligence

Olexandr Fomin*

PhD in Technical Sciences, Associate Professor
National University “Yuri Kondratyuk Poltava Polytechnic”
36011, 24 Vitaliia Hrytsaienka Ave., Poltava, Ukraine
<https://orcid.org/0009-0005-3487-9062>

Abstract. The rapid proliferation of drones in military, civilian and critical infrastructure requires fast and accurate systems for their recognition and classification. The study aimed to increase the efficiency and accuracy of drone identification by developing an approach to their classification using artificial intelligence methods in real time. The study involved the analysis of drone typology, comparative analysis of artificial intelligence methods, visual modelling, software prototyping, and evaluation of classification accuracy metrics. As a result of the first stage of the study, a classification of drones by design, purpose, size and technical characteristics that affect their visual recognition was formed. The study established that multi-rotor vehicles are the most common due to their ease of operation; single-rotor vehicles are distinguished by their carrying capacity and flight duration; fixed-wing vehicles provide speed and range; and hybrid vehicles combine vertical take-off and horizontal flight. Additionally, specialised types of drones (combat, reconnaissance, photographic, micro- and tactical) were identified, and drones were classified by size, used in the study to compare the dimensions, weight, payload and flight duration with the types of applications. The second stage of the study included a comparative analysis of artificial intelligence methods for identifying types of drones in real time. The study established that computer vision models, in particular, convolutional neural networks, provide high accuracy, and one-stage architectures provide fast object detection. Transformers and fully connected neural layers demonstrate accuracy but require significant resources. Classical machine learning algorithms, such as support vector machine (92%), random forest (89%), nearest neighbours (87.7%), and naive Bayesian classifier (79%), showed different performance. In addition, reinforcement learning can be used in systems to adapt to changes in the environment, and decision trees provide transparency in classification. The results obtained contribute to the development of real-time drone detection and classification systems for defence, infrastructure protection, airspace monitoring and public safety

Keywords: unmanned aerial vehicles; machine learning algorithms; computer recognition; neural networks; identification of rotary-winged drones

INTRODUCTION

In the context of rapid technological development, automatic object detection and classification systems are becoming increasingly relevant in various fields, from defence to civilian. Artificial intelligence (AI) technologies are central in this process, automating the analysis of complex data in real time. One of the most substantial sub-branches of AI is machine learning (ML), which is used to detect patterns in data and make decisions without explicit programming. In the context of object recognition, including unmanned aerial vehicles (UAVs),

ML ensures that the model can adapt to new data, incorporate visual differences between object types, and improve identification accuracy. Together with the development of computer vision, deep learning (DL), and transformational architectures, this opens new opportunities for creating effective recognition systems, particularly in the security sector. However, despite the availability of many approaches and solutions, there are still significant challenges associated with real-time drone classification in practical applications: limited

Article's History: Received: 06.06.2025; Revised: 30.10.2025; Accepted: 15.12.2025; Published: 25.12.2025.

Suggested Citation:

Fomin, O. (2025). Real-time drone type recognition using artificial intelligence. *Bulletin of Cherkasy State Technological University*, 30(4), 69-81. doi: 10.62660/bcstu/4.2025.69.

*Corresponding author



Copyright © The Author(s). This is an open access article distributed under the terms of the Creative Commons Attribution License 4.0 (<https://creativecommons.org/licenses/by/4.0/>)

computing resources, similarity of visual features between devices, unstable shooting conditions, and the lack of a clear typology of drones for automatic identification purposes. This creates the need for a systematic approach to drone classification and an optimal selection of AI methods for their effective recognition in dynamic environments.

A.P. Babich *et al.* (2024) determined that the main difficulty in countering modern types of drones, including First Person View (FPV) and strike drones, is timely detection, which is critical for the effective use of munitions. The study substantiated the need to develop a recognition system considering the characteristics of air targets and proposed the principles of its construction using different types of intelligence. N. Yermilova *et al.* (2023) demonstrated that DL models, in particular Faster Region Based Convolutional Neural Networks (Faster R-CNN), effectively recognise small objects of complex shape, although their use in real time is advisable only when the shape of the targets is highly complex. In addition, O.O. Korostin (2024) demonstrated that AI-based systems significantly increase the efficiency of automated object recognition in complex information flows, ensuring accuracy, speed, and reliability of processing.

G.A.S. Thomas *et al.* (2025) considered the possibilities of integrating embodied AI with computer vision in drone technology, emphasising its ability to autonomously navigate, detect obstacles and make real-time decisions in dynamic environments. R.C. Aguilera *et al.* (2025) developed an expert system based on the You Only Look Once (YOLO) architecture to automatically detect defects on wind turbine blades in real time during visual inspection by drones, which ensures fast and accurate recognition without additional computational costs. Conclusions of a study by A.S. Adebayo (2025) emphasised that the use of DL and computer vision for automated species classification significantly improves recognition accuracy, reduces field research time, and facilitates large-scale environmental monitoring using drones and real-time cameras. Furthermore, S. Chanda *et al.* (2024) developed a CNN-based system for identifying indoor room numbers using drones, achieving 92.4% recognition accuracy and improving the efficiency of autonomous navigation and delivery.

In turn, F.A.S. Islam (2025) demonstrated that AI, in particular computer vision and drone surveillance technologies, significantly improves environmental monitoring, detecting pollution, climate change, and protecting biodiversity by analysing large amounts of data in real time. The results of a study by N. Umashankar & K.S. Geethanjali (2024) emphasised that the integration of AI with unmanned platforms is used to perform object recognition, adaptive navigation, and real-time decision-making using CNN, support vector machine (SVM), reinforcement learning (RL), and other algorithms, which significantly increases their efficiency in a dynamic environment. Additionally, J. Castro *et*

al. (2024) demonstrated the effectiveness of using AI to automatically identify target microenvironments based on high-precision aerial photography, used for precision planting with reduced resource costs and increased success rates using drones. Despite significant advances in the use of AI for object recognition, drone navigation, environmental monitoring, and real-time target detection, most of the studies reviewed focused either on specific application examples (e.g., infrastructure, environment, logistics) or on the recognition of individual objects without a systematic approach to classifying specific types of drones as complex aerial targets.

In contrast, the study aimed to develop a holistic approach to drone classification using AI methods, used for real-time identification of UAV types based on their visual and technical characteristics. The objectives of the study included an analysis of the main types of drones, their characteristics, classification by size and design features, as well as a review and practical demonstration of modern AI methods used for real-time drone identification.

MATERIALS AND METHODS

At the first stage of the study, a systematic analysis of the UAV typology was conducted with a focus on their visual characteristics that are key to computer recognition through a comparative analysis of technical characteristics, classification by functional and design features, and visual analysis of drone images. The study classified drones by type of construction, main technical parameters, purpose, size and visual perception features. Within the basic typology, four main categories of drones were considered: multi-rotor, single-rotor, fixed-wing, and hybrid Vertical Take-Off and Landing (VTOL). For each type, the key technical characteristics were analysed: flight duration, range, control complexity, structural complexity, hovering ability, wind resistance, thermal endurance, payload, operating costs, applications and limitations. For each type of drone, characteristic visuals were presented, showing the typical shape, rotor configuration, and overall silhouette of the vehicle (Choosing between multi-rotor..., 2025).

Next, specific types of drones that differ in purpose and operational features were considered: small drones, micro drones, tactical drones, reconnaissance drones, large combat drones, large non-combat drones, decoy drones, Global Positioning System (GPS) dependent drones, and photographic drones (Rennie, 2016). The description of each subtype included its functional purpose, main design features, and examples of real-world use. More general categories of drones were also identified based on the principle of lift: rotary, controlled lift, and hydrogen drones (Nagel, 2025). Their generalised classification was based on technical criteria such as propulsion type, lifting configuration and energy source. Additionally, examples of multi-rotor, single-rotor, fixed-wing, and hybrid VTOL drones are provided (Gong *et al.*, 2022; Nagel, 2025; Li, 2025). To

detail the visual features that are substantial in the context of computer recognition, drones were classified by size (small, medium, and large drones) (Dukowitz, 2025). Each group was characterised by key parameters: body length or wingspan, weight class, maximum payload, average flight time, typical applications, and model examples.

At the second stage of the study, a comprehensive analysis of AI methods that can be used to identify UAV types in real time was conducted through visual modelling, software prototyping, evaluation of quality metrics, analysis of combined approaches, and adaptability analysis. The goal of this stage was to identify the most effective architectures and algorithms that ensure high classification accuracy, efficiency and adaptability in a real-world environment. Ten main approaches were considered in the analytical review: CNN, YOLO, vision transformer (ViT), multilayer perceptron mixer (MLP-Mixer), random forest (RF), SVM, K-nearest neighbours (KNN), naive bayes (NB), RL, and decision tree (DT) (Hasan & Cansever, 2023; Mrabet *et al.*, 2024; Emon *et al.*, 2025). For each method, the key advantages, typical disadvantages, and examples of use in the context of drone classification were identified. A combined approach that combines CNN and YOLO to improve the efficiency of recognising drone types in a video stream was presented separately. The diagram illustrated the stages of object detection and subsequent classification based on image features, and the architecture visualisation itself was created in RStudio using the R language.

In addition, an example of combining the ViT architecture with MLP-Mixer was provided to improve classification accuracy by better capturing both global

and local visual patterns (Essa, 2024). The effectiveness of the classical ML algorithms (RF, SVM, KNN, NB) was tested in practice by developing a software prototype of a drone type classification system in Python in the Visual Studio Code environment. For this purpose, synthetic data (1,000 samples, 10 features: 6 informative, 2 redundant) generated by the `make_classification` (scikit-learn) function, which simulates the characteristics of drones (size, weight, propellers, flight) for 4 classes (multi-rotor, fixed-wing, hybrid, unknown), was used. The sample was split into training (70%, 700 samples) and test (30%, 300 samples) parts (`train_test_split`), and the features were standardised (`StandardScaler`). The models were trained, classified 10 test samples, and evaluated in terms of accuracy, precision, recall, and F-measure in a comparative report. Lastly, the potential of RL for adaptive optimisation in changing environments and the use of DT as a simple and fast, but less robust model for basic classification were discussed.

RESULTS

Typology of drones and features of their visual recognition

As of 2025, drones have become a key element of modern technological, industrial and transport systems, and are widely used in military, civilian and commercial sectors. The significant diversity of UAVs in terms of design features, size, functional purpose and technical parameters creates significant challenges for their accurate identification and classification, especially in real time. Reliable drone type recognition is highly necessary to ensure airspace safety, prompt response to potential threats, and effective management of unmanned systems. The main types of drones are listed in Table 1.

Table 1. Results of experimental verification of noise filtering algorithms

Characteristic	Multi-rotor	Single-rotor	Fixed-wing	Fixed-wing hybrid VTOL
Flight time	20-30 minutes	30-60 minutes	1-3+ hours	45-120 minutes
Range	1-5 km	5-15 km	10-100+ km	10-80 km
Ease of control	Easy	Hard	Moderate	Moderate
The complexity of CASA	Lower	Higher	Higher	Higher
Hovering capabilities	Excellent	Good	None	Excellent
Wind resistance	Moderate	High	High	High
Thermal protection	Moderate	Good	Good	Moderate
Load capacity	Low-medium	High	Average	Medium-high
Operating expenses	Moderate	High	Low	Moderate
Best suited for	Tourism, city inspections	Mining, heavy sensors	Rural areas, agriculture	Emergency services, universal operations
Not suited for	Long distances, heavy loads	Beginners and limited budgets	Urban areas, hovering	Simple missions, limited budgets

Notes: CASA – Civil Aviation Safety Authority

Source: compiled by the author based on Choosing between multi-rotor, fixed-wing, single-rotor, and hybrid VTOL drones – AUAV's complete guide for finding your perfect match (2025)

Notably, multirotor drones are the most common due to their ease of use, manoeuvrability and hovering capabilities, which makes them suitable for aerial photography, video surveillance, inspection and 3D scanning (Rennie, 2016). Their advantages include easy control and the ability to operate in tight spaces, but their limited flight time and low payload capacity reduce their effectiveness in tasks requiring long battery life or long range. A typical example is the Inspired Flight IF800 Tomcat medium-Lift Quadcopter Drone (Nagel, 2025). Fixed-wing drones, on the other hand, resemble aircraft, providing long flight times, high speeds, and the ability to cover large areas, making them optimal for mapping, agricultural monitoring, forestry surveillance, or infrastructure inspection (Rennie, 2016). At the same time, they cannot hover, need space for launching/landing, and require more pilot training. An example of such a drone is AeroVironment's JUMP 20 (Nagel, 2025).

Single-rotor drones are structurally similar to helicopters and have a high payload capacity, long flight time (especially with gas engines), and hovering capability (Rennie, 2016). They are suitable for tasks involving heavy sensors, such as laser scanning. Their disadvantages are the complexity of control, high cost, and the need for regular maintenance. A striking example is the PULSAR monocopter (Li, 2025). On the other hand, hybrid VTOL drones combine the vertical take-off/landing capabilities of multirotor drones with the horizontal flight efficiency of fixed-wing drones. They are gradually gaining popularity in logistics, including cargo delivery, monitoring of hard-to-reach areas, and patrolling. However, their technology is still evolving, and existing models may be inferior in terms of stability in certain flight modes. One example is the TX25A drone (Gong *et al.*, 2022). To illustrate the design features and differences between the types of drones under consideration, Figure 1 shows their images.



Figure 1. Main drone types: multi-rotor, fixed-wing, single-rotor and hybrid VTOL

Notes: multi-rotor drone (hexacopter) – in the upper left corner; fixed-wing drone – in the upper right corner; single-rotor drone (helicopter type) – in the lower left corner; hybrid VTOL drone – in the lower right corner

Source: compiled by the author based on Choosing between multi-rotor, fixed-wing, single-rotor, and hybrid VTOL drones – AUAV's complete guide for finding your perfect match (2025)

In addition to the above classification, it is also advisable to distinguish several specific types of drones that have a separate functional significance (Rennie, 2016). Small drones are used primarily for recreational purposes, as they are not suitable for precise surveying or complex tasks due to their low weight and instability. They are opposed by micro drones, small UAVs used for tactical reconnaissance, especially in military operations (e.g., Black Hornet), which are capable of operating in difficult conditions such as confined spaces, strong winds, and low visibility due to their built-in micro cameras. Tactical Drones are a separate group, which are moderately sized, equipped with infrared cameras and GPS, and are usually used for medium-range surveillance.

At the same time, reconnaissance drones, such as High Altitude Long Endurance drones (HALE) and Medium Altitude Long Endurance drones (MALE), can stay in the air for tens of hours, operate at altitudes above

10 km and are designed for strategic reconnaissance. Large Combat Drones, which carry missiles or bombs, have a range of over 1,000 miles and are used for precision strikes. There are also Non-Combat Large Drones, which are used in large-scale unarmed reconnaissance missions. Target and Decoy Drones, which mimic real targets to mislead air defences or the enemy, and GPS Drones, which are capable of autonomously moving along a predefined route with high positioning accuracy, are also noteworthy. As for Photography Drones with professional-grade cameras (including 4K), they are substantial in mapping, monitoring the condition of objects and creating media content.

Several other categories of drones can be distinguished (Nagel, 2025). For example, Rotary-Wing Drones, which include helicopter and multi-rotor (quadcopters, hexacopters, octocopters, etc.) models that are widely used in various fields due to their hovering and vertical take-off/landing capabilities but have a limited flight

time due to high power consumption. Powered-Lift Drones are more complex hybrid devices that combine the advantages of both fixed-wing and rotary-wing drones, with the ability to switch between flight modes, but such designs have more complex mechanics and controls. In terms of power sources, most rotary drones are powered by electric lithium-polymer batteries, but alternative solutions such as solar panels, gas engines, hybrid systems, hydrogen fuel cells, and in-flight laser recharging technologies are being actively developed.

Hydrogen and hybrid drones show significant potential for longer flight times, while solar technologies are best suited for fixed-wing vehicles with a large surface area. After considering the types of drones, it is also worth highlighting their size, which significantly affects the technical characteristics, functional purpose and ways of using UAVs (Table 2). Size categories can be used to classify drones by weight, payload, and flight duration, which is relevant for choosing the optimal model for specific tasks and regulatory requirements.

Table 2. Classification of drones by size and main characteristics

Drone size	Small	Average	Large
Size (length/wingspan)	Less than 30 cm (<12 inches)	30-60 cm (12-24 in)	Over 60 cm (>24 in) or wingspan >1.8 m (6 ft)
Weight class	Less than 0.9 kg (<2 lbs)	0.9-4.5 kg (2-10 lb)	4.5-25+ kg (10-55+ lb)
Maximum load capacity	Up to 0.45 kg (1 lb)	1-4.5 kg (2-10 lb)	Up to 226+ kg (500+ lb)
Flight time	10-25 minutes	20-40 minutes	30-60+ minutes
Sphere of use	Recreational use, social media	Photography, inspections, mapping	Agriculture, delivery, laser radar sensing, and the film industry
Example of the model	DJI Mini 2 SE, Ryze Tello	DJI Air 3, Mavic 3 Pro	DJI Matrice 350 RTK, Alta X, Griff 300

Notes: DJI – Da-Jiang Innovations; SE – Special Edition; RTK – Real Time Kinematic

Source: compiled by the author based on Z. Dukowitz (2025)

Thus, the modern typology of drones covers a variety of design and functional categories, from light multi-rotor vehicles to large fixed-wing and hybrid VTOL models, as well as specialised tactical, reconnaissance and combat drones. Each type has advantages, limitations and applications, which determine the choice of a specific model for various tasks from recreation and photography to agricultural monitoring, logistics and defence. To efficiently and accurately recognise and classify such diverse types of drones in real time, it is recommended to use modern AI methods that provide high data processing speed and adaptability to changing flight conditions and environments.

Artificial intelligence methods for real-time identification of drone types

In a modern environment, real-time drone identification requires the use of efficient AI methods capable of rapid adaptation, high classification accuracy, and operation in conditions of limited computing resources. The use of computer vision and DL models is particularly relevant, as they can be used for the automatic selection

of relevant features, minimise human intervention, and ensure the functioning of recognition systems in a complex environment. The main AI methods used to recognise UAV types are shown in Table 3.

For instance, CNN is a deep neural network specialised in image processing. It consists of convolutional layers that automatically extract spatial features (contours, shapes, textures) and dense layers that perform classification. CNNs are effective for classifying objects that have already been found, for example, identifying the type of drone based on an image fragment (multi-rotor, fixed-wing, hybrid, etc.). However, CNN does not detect objects in the image but only classifies them. YOLO, on the other hand, is a one-step recognition model that detects and classifies objects in an image in one pass through a neural network. It is fast and works in real time and can quickly detect drones in a video stream while determining their coordinates and types. YOLO is suitable for situations where speed is critical, such as interception, monitoring, or air defence systems. The combined approach of YOLO and CNN is shown in Figure 2.

Table 3. AI methods used to identify drone types in real time

Method/Architecture	Advantages	Disadvantages	Examples of use
CNN	High accuracy, efficient on photos	Needs to be scaled, does not incorporate time	ResNet, MobileNet, EfficientNet
YOLO	Speed, real-time, compactness	Lower accuracy on small objects	YOLOv5, YOLOv8
ViT	Highly accurate, learns the global context	High memory consumption	ViT-Base, Swin Transformer
MLP-Mixer	Simplified structure, no folds	Poorer generalisability on complex data	MLP-Mixer (Google Research)

Method/Architecture	Advantages	Disadvantages	Examples of use
RF	Speed, interpretability	Low efficiency for images	Early filtering, after vectorisation
SVM	High accuracy for small data sets	Does not scale for large amounts of data	Two-stage classification
KNN	Easy to implement, no training required	Slow real-time classification	Initial prototype systems
NB	Simplicity, efficiency on sparse data	Inefficient on correlated features	Initial recognition based on metadata
RL	Ability to adapt to changing conditions	Difficulties with building a reward function	Autonomous control of UAVs
DT	Simplicity, quick construction	Retraining, limited accuracy on complex samples	Basic classification systems

Notes: ResNet – Residual Network

Source: compiled by the author based on S.H. Hasan & G. Cansever (2023), M. Mrabet *et al.* (2024), S.I. Emon *et al.* (2025)

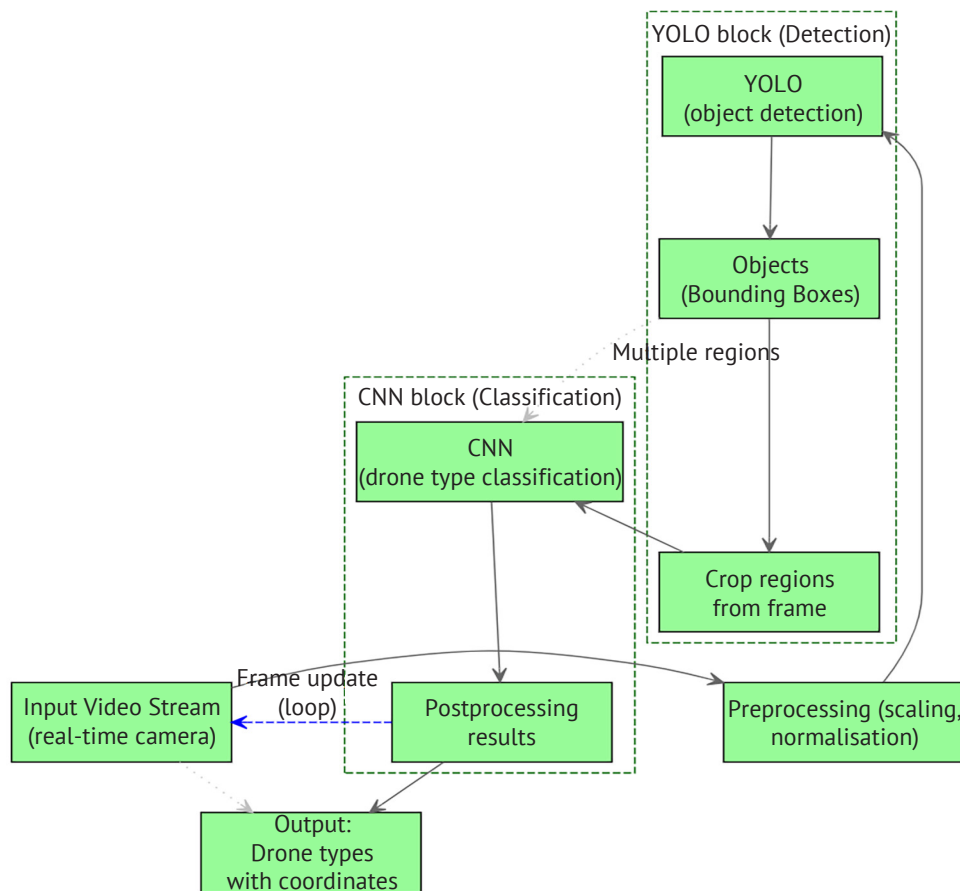


Figure 2. A combined YOLO and CNN approach for real-time drone type recognition

Source: compiled by the author

This diagram illustrates the step-by-step process of real-time drone type recognition using a combination of YOLO and CNN methods. It demonstrates how fast detection (YOLO) can be combined with accurate recognition (CNN) to effectively identify drone types in a video stream. The ViT model, which uses the transformer architecture popular in natural language processing to analyse images, is also useful. It breaks the image into patches, transforms them into a sequence of vectors, and analyses the global dependencies between these patches, which ensures efficient context capture

at different scales. ViT provides high classification accuracy based on improved recognition of global visual patterns of the drone, such as the overall shape and positioning of elements. However, the model requires significant computational resources and memory, which can be a limitation when used in embedded systems or devices with limited power.

MLP-Mixer, on the other hand, is a neural network architecture that replaces traditional convolutional layers with a sequence of layers of fully connected neurons (MLPs). It handles sequences of input image

patches, processing information by patch and feature in turn, which simplifies the model structure. In addition, MLP-Mixer offers a simpler architecture with fewer computations compared to traditional CNNs or transformers. However, due to the lack of local processing (convolution), it may not be able to handle details and complex visual features as well, which reduces the accuracy of drone classification against complex backgrounds or in unstable shooting conditions. To improve the accuracy and reliability of real-time drone type identification, ViT and MLP-Mixer models can be combined. ViT efficiently captures both local and global context of images, while MLP-Mixer provides simple and effective feature integration. For example, several ViT architectures such as Dual Attention ViT (DaViT), Inception Transformer (iFormer), and Group Propagation ViT (GPViT) can be combined through MLP-Mixer, which will combine the strengths of each model for comprehensive visual data analysis (Essa, 2024). This approach demonstrates high accuracy and robustness in complex recognition tasks, which makes it promising for use in real-time drone identification systems.

In general, ML is a substantial component of AI that can be used to automatically detect patterns in data and in subsequent decision-making. ML algorithms can be used to classify input information processed by preliminary computer vision stages, such as detection or segmentation. Among the classical ML algorithms that are the most effective for drone recognition, the

following are worth highlighting: RF, SVM, KNN, and NB. For example, RF is an ensemble method that combines several DTs to improve classification accuracy and robustness. It can process vectorised drone features (e.g., body shape, wing length, number of propellers) obtained after the detection stage and classify the type of drone. Due to its structure, RF is less sensitive to overfitting than a single tree. In addition, SVM is a method that finds the hyperplane that best separates data between classes. For drone classification, SVM can work effectively with a small or medium feature set, especially if the features are pre-normalised. The algorithm is efficient in detecting complex boundaries between classes, which is an advantage when there is a variety of drone types.

KNN is an intuitively simple but robust method that determines the class of an object by its nearest neighbours in the feature space. It does not require a training phase but is sensitive to the scale of the features. In the current context, KNN can be used as a reference model or a basic classifier in the early stages of system development. Additionally, NB is a probabilistic model based on Bayes' rule with the assumption of feature independence. In the drone classification task, it is suitable for initial sorting or quick filtering of data, especially when using sparse or categorical features such as sensor type or geographic location. An example of a Python program that demonstrates the use of these ML algorithms to classify drone types is shown below (Fig. 3).

```

31 scaler = StandardScaler()
32 X_train = scaler.fit_transform(X_train)
33 X_test = scaler.transform(X_test)
34
35 models = {
36     'Random Forest': RandomForestClassifier(random_state=42),
37     'SVM': SVC(probability=True),
38     'KNN': KNeighborsClassifier(),
39     'Naive Bayes': GaussianNB()
40 }
41
42 all_reports = []
43
44 for name, model in models.items():
45     model.fit(X_train, y_train)
46     y_pred = model.predict(X_test)
47
48     print(f"\n==== {name} ====")
49     print("Classification examples for the first 10 samples:")
50     for i in range(10):
51         print(f"[{i+1}] Actual: {drone_labels[y_test[i]]} | Predicted: {drone_labels[y_pred[i]]}")
52
53     report_dict = classification_report(y_test, y_pred, target_names=drone_labels.values(), output_dict=True)
54
55     df_report = pd.DataFrame(report_dict).transpose()
56     df_report['Model'] = name
57
58     all_reports.append(df_report)

```

Figure 3. A code snippet for classifying drone types using ML algorithms

Source: compiled by the author

The programme demonstrates the application of the considered ML algorithms (RF, SVM, KNN, and NB) to classify drone types based on synthetic data. Data is generated, divided into training and test samples, after which each model is trained and predicts the results. Examples of the classification of the first 10 test

samples are displayed, and detailed reports with key metrics are generated to compare the performance of the algorithms in one report. As a result, all algorithms showed the ability to classify drone types with varying accuracy (Fig. 4). The best results were obtained by SVM (92% accuracy) with high prediction accuracy,

recall and F-measure, especially for hybrid and fixed-wing types of drones. RF also demonstrated a high result (accuracy ~89%), KNN is slightly lower (~87.7%),

and NB has the lowest accuracy (~79%). The main errors occurred for the fixed-wing type of drone due to the similarity of features.

	Model	precision	recall	f1-score	support
Multicopter	Random Forest	0.881579	0.848101	0.864516	79.000000
Fixed-wing	Random Forest	0.869565	0.869565	0.869565	69.000000
Hybrid	Random Forest	0.927711	0.927711	0.927711	83.000000
Unknown	Random Forest	0.875000	0.913043	0.893617	69.000000
accuracy	Random Forest	0.890000	0.890000	0.890000	0.890000
macro avg	Random Forest	0.888464	0.889605	0.888852	300.000000
weighted avg	Random Forest	0.890066	0.890000	0.889854	300.000000
Multicopter	SVM	0.852273	0.949367	0.898204	79.000000
Fixed-wing	SVM	0.936508	0.855072	0.893939	69.000000
Hybrid	SVM	0.963415	0.951807	0.957576	83.000000
Unknown	SVM	0.940299	0.913043	0.926471	69.000000
accuracy	SVM	0.920000	0.920000	0.920000	0.920000
macro avg	SVM	0.923123	0.917323	0.919047	300.000000
weighted avg	SVM	0.922642	0.920000	0.920151	300.000000
Multicopter	KNN	0.843373	0.886076	0.864198	79.000000
Fixed-wing	KNN	0.813333	0.884058	0.847222	69.000000
Hybrid	KNN	0.938272	0.915663	0.926829	83.000000
Unknown	KNN	0.918033	0.811594	0.861538	69.000000
accuracy	KNN	0.876667	0.876667	0.876667	0.876667
macro avg	KNN	0.878253	0.874348	0.874947	300.000000
weighted avg	KNN	0.879891	0.876667	0.877010	300.000000
Multicopter	Naive Bayes	0.769231	0.632911	0.694444	79.000000
Fixed-wing	Naive Bayes	0.712500	0.826087	0.765101	69.000000
Hybrid	Naive Bayes	0.923077	0.867470	0.894410	83.000000
Unknown	Naive Bayes	0.766234	0.855072	0.808219	69.000000
macro avg	Naive Bayes	0.792760	0.795385	0.790544	300.000000
weighted avg	Naive Bayes	0.798057	0.793333	0.792187	300.000000

Figure 4. Comparative report of models by the main metrics of drone type classification

Source: compiled by the author

That is, SVM and RF are the most suitable for identifying drone types in real time. In addition, it is worth analysing RL and DT, which are also substantial components of AI in the context of drone recognition. RL is a method where an agent learns to make a sequence of decisions in an environment through a system of rewards and punishments. It is useful for adaptation in dynamic, changing environments, such as autonomous drone control or developing strategies to counter threats. RL can also be used to adaptively optimise the parameters of a recognition system in real time. Moreover, DTs are interpreted models that divide data into subsets using sequential feature-based rules. They are simple to implement and fast to learn but can be over-trained on complex data and have limited accuracy. In drone recognition tasks, DTs can be used for basic classification based on simple visual or sensory features, especially when transparency of decisions is required.

Thus, modern AI methods for real-time drone-type identification combine high accuracy, speed, and adaptability. The most effective are computer vision models, in particular CNN and YOLO, which provide high-quality classification and fast detection, respectively. ViT and MLP-Mixer offer new opportunities for deeper image analysis, although they have higher resource requirements. Classical ML algorithms such as RF and SVM have proven to be reliable and accurate in classifying

drone types based on vectorised features, while KNN and NB can serve as complementary methods. Additionally, RL and DT methods complement AI systems by ensuring adaptation in dynamic environments and interpretability of results. Overall, a combined approach using different models and algorithms creates a robust, efficient, and flexible drone-type recognition system that meets the requirements of real-time and complex operating environments.

DISCUSSION

A holistic approach to real-time drone type classification using YOLO was implemented, covering visual and technical characteristics. The results highlighted that YOLO, as a one-step architecture, can not only localise drones in the image, but also immediately determine their type, which is critical in conditions of limited decision-making time. In turn, P. Sumathi *et al.* (2025) proposed a drone detection system based on YOLOv8 using a labelled dataset and implemented in a Flask application. The results of both studies are consistent in the context of YOLO's real-time performance, but the present study addressed the typological classification of drones, which incorporates their technical and design features, increasing the practical value of the model.

While YOLO is considered for fast object detection in a video stream, CNN is used to accurately determine

their types based on visual characteristics. Instead, D. Li *et al.* (2024) applied a hybrid CNN + Long Short-Term Memory + self-attention model to detect drones as threats in perimeter security using acoustic data. That is, the results of both studies confirm the effectiveness of CNN in drone classification, but in contrast to the current model, which emphasised visual recognition of UAV types, the approach focuses on sensory detection of unauthorised intrusion without distinguishing between drone types. In general, the results addressed real-time drone type recognition using AI techniques to improve the accuracy of drone identification in security and monitoring environments. At the same time, N. Jain & S. Lenka (2025) analysed the role of drones and AI in precision agriculture, emphasising their importance for resource optimisation, increasing yields and automating agricultural processes. Thus, the results of the studies complement each other: both approaches demonstrate the practical value of AI drones, but for different purposes: safety and classification in the current case, and agricultural efficiency in the analysed study.

This study also demonstrated the effectiveness of using ViT and MLP-Mixer models for deep classification of drone types based on visual features, where ViT provides global recognition of the object structure, and MLP-Mixer simplifies the architecture without losing key characteristics. At the same time, H. Wasswa *et al.* (2025) evaluated the performance of ViT-MLP architectures in multi-class classification tasks, showing their high accuracy and stability compared to graph-based approaches. Thus, both studies confirm the feasibility of combining ViT and MLP in complex recognition tasks, although the present study addressed drone typology, while the aforementioned study investigated structured data and overall model performance.

This paper implements an approach to real-time classification of drone types using AI methods, with a focus on visual and design features of UAVs for monitoring, security, and situational awareness. Instead, M.A.R. Estrada (2025) proposes the concept of an autonomous NeuronDrone-Box module that makes combat decisions (attack/defence) based on chaotic dynamics algorithms and the author's methodology of econometrics, which combines economic and topographic parameters to model strategic scenarios. Although both works are based on the application of AI in drone systems, the current research is focused on recognising drone types in the visual field of view for peaceful and security tasks, while the aforementioned development focuses on autonomous combat control. The findings highlighted the importance of ML algorithms, including RF and SVM, for classifying drone types based on vectorised features after image preprocessing, which provided high accuracy (SVM 92%, RF ~89%) with efficient classification even in complex cases. Similarly, in M. Kassab *et al.* (2023), SVM and RF were used to classify drones, with SVM showing the highest accuracy among the classical methods.

Both approaches confirm the effectiveness of these ML algorithms in aerial object classification tasks, but the present research addressed typological drone recognition, which makes it more targeted and applicable within visual monitoring systems.

Among the drones reviewed, multirotor drones are singled out as one of the most common types of UAVs due to their manoeuvrability, ease of control, hovering capability, and suitability for inspection, tourism, and video surveillance, despite their limited range and payload. In the same context, S. Hoang & I.Y. Shen (2024) analysed in detail the behaviour of a large multi-rotor (18-rotor) drone in wind gusts, which demonstrated the high sensitivity of the trajectory to gust parameters and the difficulty of predicting the response without the use of stochastic models. Thus, the results of the studies complement each other: the current approach outlines where and how multirotor systems should be used, while the authors cited above clarify their physical limitations in modelling and operation.

Compared to the current study, which addressed real-time recognition and classification of drone types using visually oriented AI methods, P.A. Darwinto *et al.* (2025) analysed autonomous drone control, sensor data processing, and UAV behaviour modelling to improve energy efficiency and navigation accuracy. Both approaches fulfil the potential of AI in drone systems, but with different goals: the current approach was focused on typological classification of drones for security purposes, while the aforementioned one was aimed to improve autonomy and controllability. The results of these works can be considered complementary in the context of expanding the functionality of AI-based drones.

The study also used KNN and NB algorithms as basic models for classifying drone types, with further evaluation of their effectiveness by accuracy, completeness, and F-measure. In turn, S.R. Medarametla & G. Thallapally (2025) compared these algorithms in the context of recommended systems in terms of performance and resource efficiency. In both cases, the superiority of KNN in classification accuracy was confirmed, while NB demonstrated higher speed and lower memory consumption. This correlates with the current results, where KNN showed better recognition quality and NB was effective for fast pre-filtering at low computational cost. While this paper implements real-time classification of drone types based on visual characteristics using ML and DL methods, including RL, R. San-Segundo *et al.* (2024) proposed a hybrid navigation system for drones that combines RL and expert rules for decision making in complex environments. Both approaches use AI in the context of unmanned systems, but with different goals. Moreover, the results of both papers highlight the benefits of combined AI solutions: in the current case, it is a combination of YOLO and CNN for accurate recognition, and in the above case, the integration of RL and rule-based logic to improve control efficiency.

One of the key types of drones studied in this paper is fixed-wing drones, which are considered a separate class with a characteristic aerodynamic structure, long flight range, high wind resistance, and suitability for monitoring large areas, in the agricultural or forestry sectors. M. H. Chae *et al.* (2024) also addressed this type of drone, but in the context of developing a system to counter them. The study implemented the RL method for autonomous redirection of fixed-wing drones by manipulating the Global Navigation Satellite System (GNSS) signal. Thus, both papers not only analyse fixed-wing drones as a research object but also demonstrate the effectiveness of RL: in the current case, for adaptive control and optimisation of classification systems in real time, and in the above paper, for the implementation of autonomous strategies in counter-drone defence systems.

Alongside RL as a means of adaptive control in a dynamic environment, this paper also considers DT as a simple and interpretable model for initially classifying drones based on their visual and technical characteristics. S. Milani *et al.* (2023) combined RL with DT in the framework of explanatory AI, where RL is responsible for strategic learning and DT for transparency and explanation of agent actions. Thus, the results of both studies are consistent: the combination of RL and DT provides flexible, adaptive, and at the same time interpretable drone control and classification systems, which is especially valuable for security and monitoring applications. Finally, comparing the results of the current study with the study by A. Khan *et al.* (2025), a common desire for the multisectoral application of AI methods in solving applied problems is notable. In the aforementioned paper, AI is viewed as a driver of innovation in various fields from energy and medicine to robotics, security, and digital technologies, including drone systems. Similarly, in the current study, AI algorithms (CNN, YOLO, ViT, MLP-Mixer, SVM, RF, RL, and DT) were used to identify drone types in real time, which is critical in the areas of security, monitoring, and automation. However, in contrast to the generalised cross-disciplinary review of the aforementioned study, the current research addressed the practical implementation of AI in the specific task of visual UAV classification, supplemented by the analysis of drone technical characteristics.

Thus, the study demonstrates a holistic, multi-level approach to real-time drone type classification that combines visual models with traditional ML algorithms and RL methods. The obtained results confirm the effectiveness of AI in security, monitoring, and situational awareness tasks, providing high recognition accuracy, adaptability to changing conditions, and transparency of decision-making. Compared to other studies that focus on specific aspects such as navigation, recommender systems, sensory processing, or autonomous control, the current study is notable for the focus on typological drone recognition as a component of intelligent visual analysis systems. This is a significant contribution to the development of applied AI in drone technology,

especially in the context of real-time and mission-critical application scenarios.

CONCLUSIONS

The results of the first stage showed that for 2025, the main structural types of drones are multi-rotor (20-30 minutes of flight time, 1-5 km range), single-rotor (30-60 minutes, 5-15 km), fixed-wing (1-3+ hours, 10-100+ km) and hybrid VTOL (45-120 minutes, 10-80 km), which differ in terms of payload, wind resistance, control complexity and application. The size classification showed a wide range of technical characteristics: flight time from 10 to over 60 minutes, payload from 0.45 kg to over 226 kg, which determines the effectiveness in various areas from recreation to agricultural monitoring and military operations. Multi-rotor drones dominate due to their ease of control and hovering capabilities but have limitations in terms of flight time and range. Fixed-wing and hybrid models provide a longer range, expanding their capabilities for large-scale missions. The wide variety of types and characteristics creates challenges for accurate visual recognition in real time, requiring the use of highly efficient AI methods.

The findings of the second stage confirmed that AI methods that combine accuracy, speed, and adaptability are the most effective for identifying drone types in real time. CNN provides high accuracy in image classification, while YOLO ensures fast detection and classification of drones in a video stream in real time. The ViT model achieves high accuracy by analysing the global context but requires significant computing resources. MLP-Mixer is promising for model integration but requires improvements in recognising complex features. Among the classical ML algorithms, RF and SVM proved to be the most reliable, with an accuracy of about 89% and 92%, while KNN and NB are of secondary importance. RL and DT methods complement the system by ensuring adaptability and interpretability of solutions. The integrated use of these methods creates a robust, flexible platform for real-time drone identification with limited resources. The main limitations of the study are the need for large amounts of training data, high power consumption of complex models (ViT, MLP-Mixer) and limited accuracy of classical algorithms in difficult conditions. This reduces the efficiency of drone-type recognition in mobile or resource-limited systems. Further research should focus on hybrid models, optimisation of architectures, and development of adaptive learning to improve accuracy and versatility in real-time.

ACKNOWLEDGEMENTS

None.

FUNDING

None.

CONFLICT OF INTEREST

None.

REFERENCES

- [1] Adebayo, A.S. (2025). AI driven species recognition and digital systematics: Applying artificial intelligence for automated organism classification in ecological and environmental monitoring. *International Journal of Research Publication and Reviews*, 6(2), 31-49. doi: [10.55248/gengpi.6.0225.0703](https://doi.org/10.55248/gengpi.6.0225.0703).
- [2] Aguilera, R.C., Mosqueda, M.A.A., Mosqueda, M.E.A., & Coronel, S.L.G. (2025). YOLO expert system for real-time pattern recognition using drones on wind farm turbine. *Fractals*, 33(5), article number 2550047. doi: [10.1142/S0218348X25500471](https://doi.org/10.1142/S0218348X25500471).
- [3] Babich, A.P., Kibitkin, S.O., Georgiev, Yu.V., & Belzetskiy, R.S. (2024). Formation of a system for detection and recognition of the unmanned aerial vehicles. *Visnyk of Vinnytsia Politechnical Institute*, 176(5), 109-114. doi: [10.31649/1997-9266-2024-176-5-109-114](https://doi.org/10.31649/1997-9266-2024-176-5-109-114).
- [4] Castro, J., Alcaraz-Segura, D., Baltzer, J.L., Amorós, L., Morales-Rueda, F., & Tabik, S. (2024). Automated precise seeding with drones and artificial intelligence: A workflow. *Restoration Ecology*, 32(5), article number e14164. doi: [10.1111/rec.14164](https://doi.org/10.1111/rec.14164).
- [5] Chae, M.-H., Park, S.-O., Choi, S., & Choi, C.-T. (2024). Reinforcement learning-based counter fixed-wing drone system using GNSS deception. *IEEE Access*, 12, 16549-16558. doi: [10.1109/ACCESS.2024.3358211](https://doi.org/10.1109/ACCESS.2024.3358211).
- [6] Chanda, S., Prangon, R.D., & Hoque, K.H. (2024). A CNN-based approach for room number detection using drone in indoor environment. In *2024 IEEE international conference on power, electrical, electronics and industrial applications* (pp. 410-415). Rajshahi: IEEE. doi: [10.1109/PEEIACON63629.2024.10800605](https://doi.org/10.1109/PEEIACON63629.2024.10800605).
- [7] Choosing between multi-rotor, fixed-wing, single-rotor, and hybrid VTOL drones – AUAV's complete guide for finding your perfect match. (2025). Retrieved from <https://www.auav.com.au/news/choosing-between-multi-rotor-fixed-wing-single-rotor-and-hybrid-vtol-drones-auavs-complete-guide-for-finding-your-perfect-match/>.
- [8] Darwinto, P.A., Widodo, A.M., Agustina, N.P., Wahyuadnyana, K.D., & Rahaman, M. (2025). Artificial intelligence (AI) for autonomous drones. In B.B. Gupta & F. Colace (Eds.), *AI developments for industrial robotics and intelligent drones* (pp. 55-84). Hershey: IGI Global Publishing. doi: [10.4018/979-8-3693-2707-4.ch004](https://doi.org/10.4018/979-8-3693-2707-4.ch004).
- [9] Dukowitz, Z. (2025). *Big drones: An in-depth guide*. Retrieved from <https://uavcoach.com/big-drones/>.
- [10] Emon, S.I., Rahman, M.M., Akter, A., Rajbongshi, S., Yeasmin, S., Quraishi, M.A.N., Shafkat, A., & Majeed, Y. (2025). Automated code smell detection for software quality assurance using a web-based machine learning framework. *Research Square*. doi: [10.21203/rs.3.rs-6474801/v1](https://doi.org/10.21203/rs.3.rs-6474801/v1).
- [11] Essa, E. (2024). Feature fusion vision transformers using MLP-mixer for enhanced deepfake detection. *Neurocomputing*, 598, article number 128128. doi: [10.1016/j.neucom.2024.128128](https://doi.org/10.1016/j.neucom.2024.128128).
- [12] Estrada, M.A.R. (2025). Full autonomous artificial intelligence in attack or defense decisions making in military drones box: The NeuronDrone-box. *Journal of Advances in Artificial Intelligence*, 3(2), 169-179. doi: [10.18178/JAAI.2025.3.2.169-179](https://doi.org/10.18178/JAAI.2025.3.2.169-179).
- [13] Gong, J., Li, D., Yan, J., Hu, H., & Kong, D. (2022). Comparison of radar signatures from a hybrid VTOL fixed-wing drone and quad-rotor drone. *Drones*, 6(5), article number 110. doi: [10.3390/drones6050110](https://doi.org/10.3390/drones6050110).
- [14] Hasan, S.H., & Cansever, G. (2023). *Drone tracking and object detection by YOLO and CNN*. *International Journal of Scientific Trends*, 2(7), 78-108.
- [15] Hoang, S., & Shen, I.Y. (2024). Effects of deterministic gust modeling for large, multi-rotor drones. In *ASME 2023 international mechanical engineering congress and exposition* (article number IMECE2023-113645). New Orleans: American Society of Mechanical Engineers. doi: [10.1115/IMECE2023-113645](https://doi.org/10.1115/IMECE2023-113645).
- [16] Islam, F.A.S. (2025). The role of artificial intelligence in environmental monitoring for sustainable development and future perspectives. *Journal of Global Ecology and Environment*, 21(2), 164-179. doi: [10.56557/jogee/2025/v21i29272](https://doi.org/10.56557/jogee/2025/v21i29272).
- [17] Jain, N., & Lenka, S. (2025). *Artificial intelligence based precision agriculture for enhanced productivity*. doi: [10.13140/RG.2.2.35586.59843](https://doi.org/10.13140/RG.2.2.35586.59843).
- [18] Kassab, M., Zitar, R.A., El Fallah, A., & Barbaresco, F. (2023). Bird/Drone detection and classification using classical and deep learning methods. *Authorea*. doi: [10.22541/au.168075364.45332093/v1](https://doi.org/10.22541/au.168075364.45332093/v1).
- [19] Khan, A., Kumar, K., & El Sayed, A.F. (2025). Unveiling the sky: Exploring synergies in drone robotics and automation through artificial intelligence and machine learning. In A. Khan, M.K. Hasan, M. Varish & M.A. Husain (Eds.), *Advancements in artificial intelligence and machine learning* (pp. 182-200). Singapore: Bentham Science Publishers. doi: [10.2174/9789815322583125010012](https://doi.org/10.2174/9789815322583125010012).
- [20] Korostin, O.O. (2024). Efficiency of text recognition in the automation of international maritime transport with the help of artificial intelligence. *Taurida Scientific Herald, Technical Sciences*, 3, 29-38. doi: [10.32782/tnv-tech.2024.3.4](https://doi.org/10.32782/tnv-tech.2024.3.4).
- [21] Li, D., Yi, D., Zhou, X., Chen, X., Geng, Y., & Li, X. (2024). Multisource threatening event recognition scheme targeting drone intrusion in the fiber optic DAS system. *IEEE Sensors Journal*, 24(20), 32185-32195. doi: [10.1109/JSEN.2024.3449440](https://doi.org/10.1109/JSEN.2024.3449440).

- [22] Li, M. (2025). Beyond conventional drones: A review of unconventional rotary-wing UAV design. *Drones*, 9(5), article number 323. doi: [10.3390/drones9050323](https://doi.org/10.3390/drones9050323).
- [23] Medarametla, S.R., & Thallapally, G. (2025). Comparing K-nearest neighbors and naive bayes in real-time recommendation systems. *Global Journal of Engineering Innovations and Interdisciplinary Research*, 5(1), article number 18. doi: [10.33425/3066-1226.1075](https://doi.org/10.33425/3066-1226.1075).
- [24] Milani, S., Zhang, Z., Topin, N., Shi, T.R., Kamhoua, C., Papalexakis, E.E., & Fang, F. (2023). MAVIPER: Learning decision tree policies for interpretable multi-agent reinforcement learning. In M.-R. Amini, S. Canu, A. Fischer, T. Guns, P.K. Novak & G. Tsoumakas (Eds.), *European conference: Machine learning and knowledge discovery in databases* (pp. 251-266). Cham: Springer. doi: [10.1007/978-3-031-26412-2_16](https://doi.org/10.1007/978-3-031-26412-2_16).
- [25] Mrabet, M., Sliti, M., & Ben Ammar, L. (2024). Machine learning algorithms applied for drone detection and classification: Benefits and challenges. *Frontiers in Communications and Networks*, 5, article number 1440727. doi: [10.3389/frcmn.2024.1440727](https://doi.org/10.3389/frcmn.2024.1440727).
- [26] Nagel, L. (2025). *Types of drones and UAVs*. Retrieved from <https://www.tytorobotics.com/blogs/articles/types-of-drones?srsId=AfmBOoquLrJpCU9jWWg4oiOfy4Id2TWE2u9kUQ1vpiWuULcFMcPsvBHO>.
- [27] Rennie, J. (2016). *Drone types: Multi-rotor vs fixed-wing vs single rotor vs hybrid VTOL*. Retrieved from <https://www.auav.com.au/articles/drone-types/>.
- [28] San-Segundo, R., Angulo, L., Gil-Martin, M., Carramiñana, D., & Bernardos, A.M. (2024). Hybrid artificial intelligence strategies for drone navigation. *AI*, 5(4), 2104-2126. doi: [10.3390/ai5040103](https://doi.org/10.3390/ai5040103).
- [29] Sumathi, P., Thungashree, Y.S., & Pushpalatha, S. (2025). Real-time drone type detection for smart air traffic monitoring. In *National level technical symposium (Advaya 2k25)* (pp. 44-47). New Delhi: All India Council for Technical Education. doi: [10.59544/WOWN1934/ADVAYA2K25P10](https://doi.org/10.59544/WOWN1934/ADVAYA2K25P10).
- [30] Thomas, G.A.S., Muthukaruppasamy, S., Kumar, S.S., Karthikeyan, B.J., & Krishnan, S. (2025). Navigating the nexus: Unravelling challenges, ethics, and applications of embodied AI in drone technology through the lens of computer vision. In P. Raj, A. Rocha, S.P. Singh, P.K. Dutta & B. Sundaravadivazhagan (Eds.), *Building embodied AI systems: The agents, the architecture principles, challenges, and application domains* (pp. 61-78). Cham: Springer. doi: [10.1007/978-3-031-68256-8_3](https://doi.org/10.1007/978-3-031-68256-8_3).
- [31] Umashankar, N., & Geethanjali, K.S. (2024). A comprehensive study of artificial intelligence applications of drone. *Engineering Archive*. doi: [10.31224/4194](https://doi.org/10.31224/4194).
- [32] Wasswa, H., Abbass, H., & Lynar, T. (2025). Are GNNs worth the effort for IoT botnet detection? A comparative study of VAE-GNN vs. ViT-MLP and VAE-MLP approaches. *ArXiv*. doi: [10.48550/arXiv.2505.17363](https://doi.org/10.48550/arXiv.2505.17363).
- [33] Yermilova, N., Zourab, Y., & Iermilov, R. (2023). Methods of complex objects automatic recognition by form. *Control, Navigation and Communication Systems*, 4(74), 80-84. doi: [10.26906/SUNZ.2023.4.080](https://doi.org/10.26906/SUNZ.2023.4.080).

Розпізнавання типів дронів у реальному часі за допомогою штучного інтелекту

Олександр Фомін

Кандидат технічних наук, доцент

Національний університет «Полтавська політехніка імені Юрія Кондратюка»

36011, просп. Віталія Грицаєнка, 24, м. Полтава, Україна

<https://orcid.org/0009-0005-3487-9062>

Анотація. Стрімке поширення дронів у військовій, цивільній та критичній інфраструктурі вимагає створення швидких і точних систем для їх розпізнавання та класифікації. Мета дослідження полягала у підвищенні ефективності і точності ідентифікації дронів шляхом розроблення підходу до їх класифікації з використанням методів штучного інтелекту в умовах реального часу. У процесі дослідження застосовано аналіз типології дронів, порівняльний аналіз методів штучного інтелекту, візуальне моделювання, програмне прототипування та оцінку метрик точності класифікації. У результаті першого етапу дослідження сформовано класифікацію дронів за конструкцією, призначенням, розміром і технічними характеристиками, що впливають на їх візуальне розпізнавання. Встановлено, що мультироторні апарати є найпоширенішими через простоту керування; однороторні – вирізняються вантажопідйомністю та тривалістю польоту; фіксованокрилі – забезпечують швидкість і дальність; гібридні – поєднують вертикальний зліт і горизонтальний політ. Додатково виокремлено спеціалізовані типи безпілотників (бойові, розвідувальні, фотографічні, мікро- та тактичні), а також класифіковано дрони за розміром, що дозволило зіставити габарити, вагу, вантажопідйомність і тривалість польоту з типами застосування. Другий етап дослідження охопив порівняльний аналіз методів штучного інтелекту для ідентифікації типів дронів у реальному часі. Встановлено, що моделі комп'ютерного зору, зокрема згорткові нейронні мережі, забезпечують високу точність, а одноетапні архітектури – швидку детекцію об'єктів. Трансформери й повнозв'язні нейронні шари демонструють точність, але потребують значних ресурсів. Класичні алгоритми машинного навчання, зокрема метод опорних векторів (92 %), випадковий ліс (89 %), метод найближчих сусідів (87,7 %) та наївний баєсівський класифікатор (79 %) показали різну ефективність. Крім того, підкріплювальне навчання дозволяє адаптувати системи до змін середовища, а дерева рішень забезпечують прозорість класифікації. Отримані результати сприяють розробці систем виявлення та класифікації дронів у реальному часі для оборони, охорони інфраструктури, моніторингу повітряного простору та громадської безпеки

Ключові слова: безпілотні літальні апарати; алгоритми машинного навчання; комп'ютерне розпізнавання; нейронні мережі; ідентифікація роторних безпілотників



AI-based model of a researcher support service

Maksym Shovkoplias*

Postgraduate Student

Sumy State University

40000, 116 Kharkivska Str., Sumy, Ukraine

<https://orcid.org/0009-0007-4192-4743>

Abstract. Between 2020 and 2025, researchers faced challenges such as fragmented digital platforms, information overload, and limited personalisation capabilities. This underscored the need for services capable of providing comprehensive support for research activities. The aim of this study was to develop a conceptual model of an intelligent information service focused on personalised researcher support. The proposed system architecture was built using structural modelling, functional analysis, machine learning, and natural language processing techniques. It includes modules for recommendations, virtual collaboration, event management, and automated bibliography generation. A multi-layered user model was designed, taking into account scientific interests, interaction history, and research context. The combination of semantic analysis with behavioural patterns increased recommendation relevance by 20-30%. The prototype of the system was tested in March 2025 with the participation of 15 young scientists from three Ukrainian universities. The results of the survey and practical tasks showed that the average time spent searching for relevant literature was reduced by 35%, task planning efficiency increased by 40%, and user satisfaction with the service's functionality reached 87%. Respondents highly rated the convenience of the interface (4.5 out of 5), the relevance of recommendations (4.3), and co-authoring tools (4.6). Three new academic collaborations were initiated through the co-author selection module. The data obtained confirmed the effectiveness of the model in increasing research productivity, improving collaboration, and providing personalised user support. The proposed structure allows for scaling to different disciplines and has the potential to be implemented in digital platforms focused on scientific activity

Keywords: scientific information personalisation; semantic analysis; adaptive recommendations; machine learning; intelligent systems; digital research environment

INTRODUCTION

In the contemporary research landscape, scientists operate in increasingly complex digital environments marked by the exponential growth of academic content, decentralised data sources, and limited interoperability between platforms. These challenges are compounded by the necessity for personalised workspaces, efficient collaboration mechanisms, and intelligent content filtering. Existing solutions often offer isolated functionalities – bibliographic management, social networking, or search-but fail to deliver unified services that support the full research lifecycle in a seamless and adaptive way. This highlights the need for intelligent information systems capable of integrating recommendation

engines, virtual collaboration tools, scheduling, and automated bibliography generation in a personalised and scalable manner.

According to T. Adewale (2022), machine learning algorithms have a significant impact on personalised recommendation systems. The author proposed models that apply collaborative filtering and deep learning to enhance user engagement and relevance. The study demonstrated that machine-driven adaptation substantially improves access to and retention of information by researchers. O.B. Akinagbe (2024) explored broader applications of artificial intelligence in scientific workflows. Author's research confirmed that AI-based services

Article's History: Received: 22.06.2025; Revised: 08.11.2025; Accepted: 15.12.2025; Published: 25.12.2025.

Suggested Citation:

Shovkoplias, M. (2025). AI-based model of a researcher support service. *Bulletin of Cherkasy State Technological University*, 30(4), 82-96. doi: 10.62660/bcstu/4.2025.82.

*Corresponding author



Copyright © The Author(s). This is an open access article distributed under the terms of the Creative Commons Attribution License 4.0 (<https://creativecommons.org/licenses/by/4.0/>)

help reduce information overload and enable strategic planning by automating routine research tasks and aligning system responses with behavioural input.

C.A. Putri *et al.* (2025), working in the field of education, demonstrated the effectiveness of adaptive digital libraries that apply semantic indexing to personalise content delivery. Although their study was focused on primary education, the described mechanisms are applicable to academic research platforms as well. Y. Koval (2022), a Ukrainian researcher, analysed the structure and function of academic platforms and emphasised the necessity of integrating user profiles with external databases. This researcher's findings identified key gaps in personalisation and highlighted the importance of unified interface design.

M.O. Shovkopljas & V.O. Liubchak (2024) proposed a modular framework for personalised information services that combined semantic analysis, collaborative tools, and real-time behavioural adaptation. Their study underlined the need for context-aware services that dynamically evolve with user behaviour. F. Yu *et al.* (2020) explored interdisciplinary UX methods for accessing academic content. According to their research, integrating user behaviour into service design improves both usability and engagement, particularly in cross-functional research environments. C.K. Kreutz & R. Schenkel (2022) examined recommendation systems for academic papers. According to their results, hybrid models that combined semantic metadata with machine learning significantly enhance the accuracy of scientific content delivery.

In their article, D. Petryna *et al.* (2024) explored the possibilities of using artificial neural network tools to accelerate the development of web interfaces. They analysed how modern AI technologies can optimise the process of creating UI solutions, in particular through the automation of certain stages of design and coding. As a result of their research, the authors concluded that the use of neural networks can reduce labour costs and increase the efficiency of web interface development, especially in the context of rapid digitalisation and the need for adaptive design. However, the paper hardly addresses the integration of such AI tools into larger systems supporting researchers or the educational process, nor does it provide an in-depth analysis of their impact on UX quality from the end user's perspective.

In the publication of L. Nikiforova *et al.* (2025), authors focused on the creation of an information resource and an electronic register of scientific professional publications as tools for digital support of scientific activity. They described the architectural and organisational aspects of developing such systems, which aimed to improve access to high-quality scientific sources and support the evaluation of scientific output. The findings of the study indicated the feasibility and necessity of creating such registries in the context of the digitisation of science, as well as the positive impact of such initiatives on the transparency and convenience of

scientific communication. At the same time, the authors did not focus enough on the intellectual mechanisms of analysis and recommendations within this resource: there are no modules that would provide personalised advice or flexible search using AI, which could significantly enhance the functionality of the system in the context of supporting researchers.

While these studies contributed valuable insights, most of them focused on isolated functions or specific user scenarios. Previous studies have not fully addressed the development of holistic systems that integrate recommendation, collaboration, scheduling, and personalisation into a single adaptive ecosystem. Moreover, issues such as multilingual support, behavioural feedback loops, and seamless integration with external infrastructures are often underexplored. These limitations justify the need for the present study, which aimed to design and evaluate a modular, intelligent information service tailored to researchers' evolving needs.

MATERIALS AND METHODS

The methodological foundation of this research was a modular systems design approach, incorporating principles from computer science, user experience (UX) theory, applied machine learning, and language processing. The system was conceptualised as a cloud-based architecture comprising intelligent services such as data aggregation, semantic analysis, user profiling, recommendation logic, and event tracking. At the conceptual design stage, the architecture was developed using system engineering and modular design principles. Functional decomposition helped outline the main subsystems: the user profile manager, semantic recommendation engine, collaboration modules, and calendar.

This structure was visualised using UML diagrams and logical models created in draw.io and Lucidchart. The user interface was shaped through a comparative analysis of academic platforms (Mendeley, ResearchGate, Academia.edu) and UX principles derived from prior studies (Yu *et al.*, 2020). This content analysis included both qualitative and quantitative evaluations of usability, content filtering precision, and collaborative tools. Limitations such as lack of real-time adaptation and weak integration served as design prompts for new features.

To achieve dynamic personalisation, the system employed machine learning techniques using scikit-learn and TensorFlow. Supervised methods (e.g., decision trees, logistic regression) and unsupervised learning (e.g., K-Means clustering) powered the classification of user behaviour and content relevance. Relevance was calculated using cosine similarity and TF-IDF metrics. User segmentation followed a two-step pipeline: behavioural clustering via K-Means and profile-based classification through decision trees. The taxonomy for segmentation aligned with Scopus and Web of Science standards, covering domains such as computer science, life sciences, engineering, social sciences, humanities, and interdisciplinary research.

Textual data were processed with spaCy and langdetect. Preprocessing involved stop-word removal, lemmatisation, and metadata translation. Named Entity Recognition (NER) enabled identification of authors, venues, and technical terms. Semantic similarity was computed using Word2Vec embeddings trained on academic corpora, forming the basis of the system's semantic module, which handled named entity recognition, keyword extraction, and topic modelling. Recommendation logic relied on a hybrid approach: content-based filtering using metadata similarity and collaborative filtering using user similarity matrices. The system's credibility scoring integrated citation metrics, journal impact factors, and author profiles from indexed repositories.

To support semantic enrichment and recommendation accuracy, the system integrated real-time scientific data from sources like arXiv, Directory of Open Access Journals (DOAJ), Scopus, and institutional archives. Data were accessed through RESTful APIs, parsed in JSON and BibTeX formats using bibtexparser and jsonlib. Over 20,000 metadata entries and publication abstracts were used to create a testbed for model evaluation and semantic processing. A built-in calendar module was developed using Flutter and synchronised via Firebase Firestore, ensuring multi-device consistency. It supported creation, updating, and deletion of research events and enabled two-way integration with Google Calendar and iCal via RESTful APIs and ICS feed handling.

For multilingual functionality, automatic translation was implemented using Google Cloud Translation API and argos-translate. Detected languages (via langdetect) routed texts to appropriate pipelines, enabling real-time translation of abstracts, queries, and metadata. The system also allowed metadata export in formats compatible with Zotero and EndNote. System performance and refinement were governed by feedback loops, combining implicit user behaviour (clicks, reading time, interactions) and explicit ratings. These were evaluated weekly using precision, recall, and F1-score, with near real-time updates based on interaction data. Mathematical modelling (in Python with SimPy) simulated user flows, feedback decay, and topic drift, ensuring the system's scalability and responsiveness.

To evaluate the system's effectiveness, a mixed-methods approach was used. A pre-test online survey (March 4, 2025) via Google Forms gathered baseline data on user productivity and task duration. The sample included 15 anonymous participants (early-career researchers and PhD students) from Kyiv, Sumy, and Dnipro. The survey, conducted asynchronously, collected no personal data and included both closed and open-ended questions. Participants rated features such as usability, recommendation relevance, interface responsiveness, and collaborative tools using a Likert scale (1-5). Example survey questions: "How much time do you typically spend searching for relevant literature using your current tools?", "Rate the ease of use of the prototype interface", "How relevant were the

recommendations generated by the system?", "To what extent did the collaboration tools improve your workflow?". The survey was conducted in accordance with the ethical principles set forth in the Declaration of Helsinki (2024), which stipulate respect for the dignity and rights of research participants. Before data collection began, respondents were informed about the purpose, methods, and possible consequences of participation, after which they gave their informed consent. Particular attention was paid to excluding any form of coercion and protecting vulnerable groups.

During the testing phase (March 4-17, 2025), the same participants completed realistic tasks: semantic search, recommendation generation, collaborative document editing, bibliography export, and scheduling research deadlines. Each task was time-tracked, and feedback was collected via Google Forms and embedded rating prompts. These self-reported scores and behavioural data served as benchmarks for assessing improvements and identifying areas for further refinement. To visualise system architecture and functional dependencies, several structural diagrams were created, representing interrelations between components and ensuring coherent service logic.

RESULTS AND DISCUSSION

Quantitative results demonstrated that the average time required to identify and collect relevant literature was reduced by 35%, and task planning efficiency improved by 40% due to the integration of the event calendar. The accuracy of recommendations was rated at 4.3 out of 5, and the collaborative editing experience received a score of 4.5 out of 5. Additionally, three new academic collaborations were initiated using the co-author suggestion module. These results confirmed that the intelligent information service model not only improves research productivity but also enhances user satisfaction and interdisciplinary cooperation.

Top-Level Architecture of the System

The top-level architecture of the proposed intelligent information service comprised a set of interconnected modules, each responsible for a distinct stage in the user-information interaction pipeline. Together, they support the full cycle of scientific information retrieval, processing, recommendation, and feedback, forming a dynamic and adaptive research environment. At the core of the architecture is the user, who interacts with the system through a dedicated interface. The user submits queries, selects scientific interests, reviews the results, and provides implicit or explicit feedback, all of which shape the system's behaviour and future outputs. The user profile module is responsible for creating and maintaining a detailed digital representation of the individual user.

This profile includes information such as scientific interests, academic background, publication history, and preferred citation styles. The profile is dynamic

and evolves in response to user interactions such as search queries, saved documents, and usage patterns. The source aggregation module connects the system to external scientific resources and repositories. These resources feed the system with structured academic content, which was further analysed and filtered to generate personalised recommendations. Retrieved data was preprocessed and temporarily cached in a cloud-based storage layer to reduce latency during repeated access.

However, the system does not maintain a permanent local database of external sources; instead, it fetches fresh data dynamically in response to user queries. This approach ensures up-to-date content while avoiding unnecessary data duplication or licensing violations. Data collection was conducted based on user-defined parameters, including keywords, authors, disciplines, and research topics. Aggregated content is processed by the analytics and filtering module, which ensures that only relevant, high-quality, and non-duplicated data is passed along the pipeline. This module performs semantic text analysis, assesses source credibility, detects duplicates, and applies preliminary topic-based filtering to ensure thematic relevance. Credibility assessment was based on citation metrics, publication venue reputation, and author profiles aggregated from indexed databases. A central role in the architecture was played by the recommendation engine, which delivers personalised lists of scientific content to each user. This engine relies on three main sources of input: the user's declared interests, the current research context (such as the active stage of a project or a specific query), and historical behavioural patterns including past searches, document views, and time spent on content.

Based on these signals, the system constructs a relevance model tailored to each individual. To operationalise personalisation, the recommendation engine applies machine learning techniques. Clustering algorithms were used to group users and content based on hidden patterns in the data. Collaborative filtering was applied to infer user preferences by leveraging similarities among users with comparable behaviour. Furthermore, the system explores the use of neural networks to capture more complex relationships between user intent and content features. These algorithms are continuously retrained as new user interaction data is collected, enabling the engine to adapt its outputs over time and improve the precision of its suggestions. To monitor learning progress, evaluation metrics such as precision, recall, and F1-score are tracked on a validation dataset that simulates real user queries. Additionally, performance is reviewed periodically using user feedback ratings and click-through behaviour. This allows for fine-tuning of model parameters and ensures continuous improvement.

The user interface module was designed to support intuitive interaction with the platform. It includes visual tools for filtering search results, interactive dashboards for managing bibliographies, collaborative

editing capabilities, and calendar-based features for managing deadlines and scientific events (Wirtz & Lovelock, 2021). It serves not only as a visual layer but also as a workspace that integrates content consumption, communication, and document production. This hybrid approach enabled seamless synchronisation across platforms, while preserving user control within the system's interface. Calendar entries can be colour-coded, linked to specific projects or documents, and optionally shared with collaborators within the platform. It serves not only as a visual layer but also as a workspace that integrates content consumption, communication, and document production. A key feature of the system's architecture is the feedback loop, which connects the outcome of each user interaction back into the recommendation and personalisation process. Whether through explicit ratings, saved items, or even the decision to ignore certain results, the system learns from behaviour and refines its models accordingly. This enables ongoing optimisation of the recommendation engine and overall user experience.

Data Flow Structure

The data flow structure of the information service reflects the sequential and cyclic logic of how data is collected, processed, and transformed into personalised scientific outputs. The process begins at the input level and proceeds through a series of functional modules, each responsible for specific operations that contribute to the generation of relevant, high-quality recommendations for researchers. At the initial stage, the input layer gathers data from multiple sources. These include the user's profile, which encapsulates declared interests, search history, and interaction preferences. In addition, contextual parameters such as the time of interaction, current user activity, and the specific topic of the query are captured in real time.

The system also connects to external scientific resources through application programming interfaces (APIs), integrating information from academic databases such as Scopus, Web of Science, and PubMed, open-access repositories like arXiv and Zenodo, digital libraries including Google Scholar and DOAJ, as well as academic event calendars. This collected data provides the foundation for subsequent processing and personalisation. Following data collection, the preprocessing module performs essential preparatory tasks. This involves cleansing the data by removing duplicates and noise, converting input formats (e.g., XML, JSON, BibTeX) into a unified internal structure, and standardising linguistic elements through translation, lemmatisation, and the normalisation of key metadata fields such as titles and author names. These steps ensure that the incoming information is coherent, structured, and ready for semantic processing.

The semantic analysis engine then applies natural language processing (NLP) techniques to interpret the content and context of the data. It identifies the

thematic scope of each source, extracts key entities such as concepts and authors, and constructs a contextual model that aligns with the user's information needs. Each document's semantic core is matched against the user profile to establish its potential relevance. Processed data is then forwarded to the filtering and ranking module, which eliminates outdated or irrelevant sources based on established thresholds of relevance, novelty, and authoritativeness. The module assigns a weight to each remaining source and organises the results in a priority order, forming a personalised list tailored to the user's current research context. The output layer is responsible for presenting the final results to the user. It delivers personalised recommendations, compiles bibliographic references in the selected citation style, and suggests opportunities for collaboration, relevant academic events, or discussions. These outputs are displayed through a dynamic, interactive dashboard that adapts to the user's ongoing activity and preferences.

A vital part of the system is the feedback and learning loop, which serves as a continuous source of system refinement. This loop involves collecting both explicit user feedback (such as ratings, likes, or comments) and implicit signals (such as time spent on content, clicks, ignored recommendations). These interactions are logged and periodically analysed to identify patterns in user behaviour. Based on these patterns, the recommendation models are retrained or fine-tuned to improve the accuracy and relevance of future outputs. This process ensures that the system continuously adapts to changing user needs and preferences. User behaviour is monitored in terms of which items are opened, saved, ignored, or explicitly rated. This behavioural data is

logged and fed back into the system at regular intervals (approximately once per week) to retrain and update the recommendation models.

Minor updates (e.g., reweighting of existing recommendations) are performed in real time or every few hours, while full model retraining is carried out on a weekly basis. This interval was selected to balance system responsiveness with computational efficiency, ensuring timely adaptation to user behaviour without overloading system resources. Feedback can be submitted either passively (through user actions) or actively (via rating options or explicit feedback forms), and is processed in real time or through scheduled batch analysis depending on the nature and frequency of user interactions. The system was designed to support immediate responsiveness to feedback, with minor updates to recommendation weights occurring almost instantly, while larger-scale model adjustments are typically performed daily or weekly, depending on data volume and usage intensity. This ensures both agility in response to evolving user needs and long-term stability in system performance.

The data flow structure diagram represents these stages as a connected series of transformations, occasionally forming cyclic feedback paths. It illustrates how raw, heterogeneous data is systematically cleaned, analysed, interpreted, and ultimately transformed into actionable, high-value outputs that enhance the research process. As shown in Figure 1, this flow is depicted as a linear or cyclic process with transformation nodes and conditional transitions. It is particularly important to illustrate how data is cleaned, semantically analysed, transformed, and ultimately returned to the user in a valuable, actionable form.

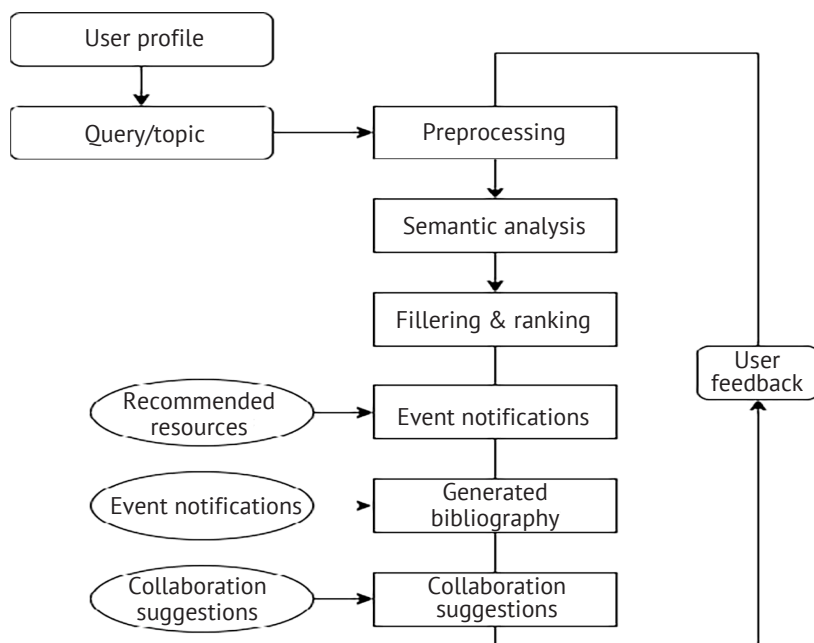


Figure 1. Data flow structured diagram

Source: created by the author

Structural Diagram of the Personalisation Module

The structural diagram of the personalisation module illustrates the internal configuration of components that enable the information service to dynamically adapt to the specific needs of individual users. Serving as the system's core intelligence layer, this module determines the platform's capacity to personalise outputs and respond to contextual user inputs. At the heart of this module lies a component responsible for collecting user data. It continuously monitors user activity and compiles detailed information regarding individual profiles, including areas of interest, disciplinary specialisation, search history, frequently viewed documents, and usage frequency. Data can be acquired both directly – via questionnaires or user-defined settings – and indirectly through behavioural analytics, which track interaction patterns over time. Once this data is collected, it is processed by the behavioural analysis engine. This subsystem applies machine learning algorithms to uncover underlying patterns in user behaviour, enabling the system to predict information needs, determine the relevance of specific content, and anticipate future interests. The analysis forms the basis for further segmentation and personalisation processes.

To enrich user profiles with social and collaborative signals, the system extracts data from publicly available academic and networking platforms such as ORCID, ResearchGate, Academia.edu, and conference aggregators like WikiCFP and AllConferences.com. Integration was implemented via publicly documented APIs (e.g., ORCID Public API, ResearchGate scraping with rate limits, RSS/ICS feeds for conferences), allowing the system to collect data on group memberships, co-authorship networks, and event participation. The user segmentation component classifies individuals into defined groups based on thematic focus, type of scientific activity, and

working style. This classification facilitates the targeted delivery of content and services, ensuring that users receive recommendations and tools aligned with their professional identity and research habits.

Building on this segmentation, the recommendation generator module creates real-time personalised suggestions, which may include scientific publications, relevant events, collaboration opportunities, funding calls, and digital tools. These recommendations are dynamically adjusted in response to changes in user behaviour, project stages, or external developments. Content is presented through a personalised interface layer, which displays information in a manner tailored to the individual's current profile. This interface may include a customisable menu, interactive dashboards with thematic news, notification panels, and integrated calendars for managing events and deadlines. It serves as the main point of interaction between the user and the adaptive features of the system.

Finally, a feedback loop ensures that the personalisation process is iterative and self-improving. The system monitors user responses to recommendations – such as whether suggestions are accepted, ignored, saved, or rated – and uses this feedback to refine future outputs. As a result, the service becomes increasingly accurate in anticipating user needs, thereby enhancing engagement and overall utility. This design ensures a closed personalisation cycle in which user data is continuously updated, allowing the system to adapt to changes in the researcher's behaviour and interests. For example, if a user previously focused on medical imaging but begins interacting more frequently with oncology-related publications, the system will adjust its recommendations accordingly – prioritising cancer research topics and suggesting relevant events or collaborators in that area (Fig. 2).

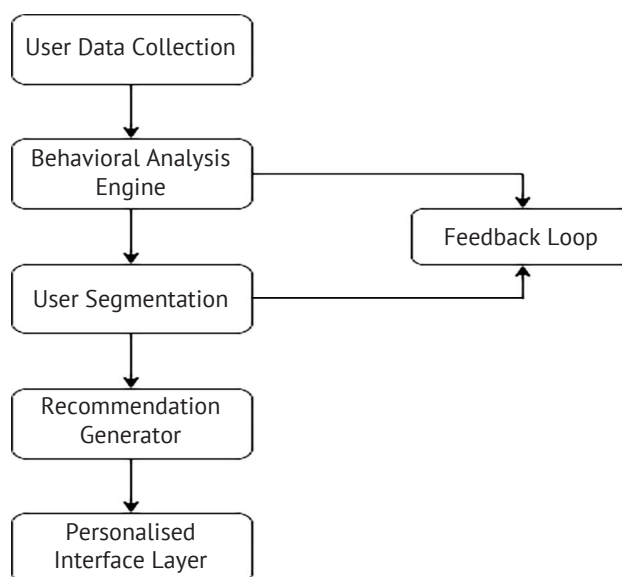


Figure 2. Structural diagram of the personalisation module

Source: created by the author

This approach enabled a comprehensive analysis of the capabilities of the information service and facilitates the development of an effective implementation model, taking into account current technological advancements

and user needs. A general model of a researcher's workplace was proposed, which incorporates the main aspects of their activities. The diagram below illustrates the interaction of the core components of this model (Fig. 3).

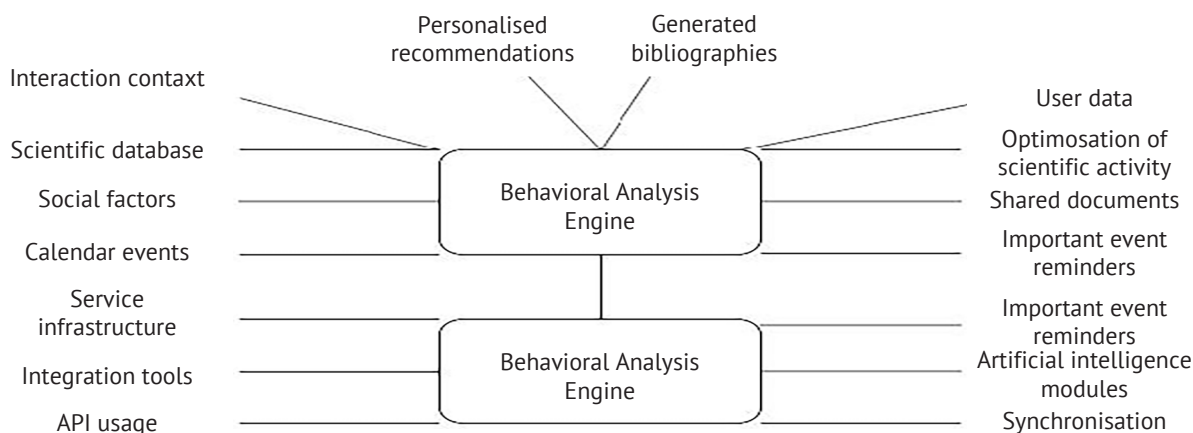


Figure 3. IDEF0 diagram Information service

Source: created by the author

In the presented diagram, the mechanism is a key element that ensures the functioning of the information service. It includes technological, software and algorithmic means that allow the implementation of the main capabilities of the service. The mechanism of the system includes several interconnected components that ensure its full functionality. The service infrastructure comprises hardware and software resources that guarantee stable operation, scalability, and efficient resource management. Integration tools provide connectivity with external systems, including scientific databases, social networks, third-party APIs, and collaboration platforms. Through the use of APIs, the system enabled dynamic data exchange with other services, maintaining flexibility and expanding functionality.

Information Service Diagram Analysis

The analysis of the information service diagram highlighted the primary categories of input and output data that define the operational logic of the system. These data flows are central to the functionality of the platform and ensure its relevance, adaptability, and usability within a modern research context. The input data used by the service originated from a diverse range of sources. One of the key elements is user-specific data, which includes individual characteristics, declared research interests, and documented academic activity such as publication and collaboration history. This user model formed the basis for personalisation and contextual relevance. Another crucial input is the broader interaction context, which captures the nature of relationships between researchers, ongoing collaborative projects, and communication processes. These data help position the user within a larger scientific network, providing valuable context for recommendation and matchmaking functions. The system also draws

extensively from external scientific databases, incorporating structured and unstructured content such as articles, patents, and monographs. These databases are accessed through dedicated integration mechanisms using public APIs and metadata harvesting protocols.

Depending on user activity, data can be fetched in real time or retrieved from a local cache created during previous sessions. A lightweight local metadata repository is maintained to improve speed and reduce repeated requests to external systems. This repository stores bibliographic data (author names, titles, abstracts, DOIs) and is refreshed periodically to ensure currency. Full-text documents are not stored locally; instead, they are loaded directly from verified open-access or licensed repositories during the user session. These databases constitute the foundational repository from which knowledge is extracted and filtered. In addition, the model accounts for social factors, including connections to academic communities, participation in scientific groups, and engagement with conferences or collaborative initiatives. These social dimensions provide insight into user networks and potential interdisciplinary opportunities. Further input was acquired through external system integrations, which allow the platform to connect with institutional repositories, electronic libraries, and third-party academic services. This interoperability expands the scope and richness of the data environment.

The service incorporated calendar-based information, which includes important scientific events, project milestones, submission deadlines, and scheduled webinars. These temporal data points contributed to the platform's ability to generate timely alerts and assist in task planning. On the basis of this comprehensive information environment, the system produced a range of outputs. The most prominent among these are personalised

recommendations, which suggest relevant scientific materials, collaboration opportunities, and upcoming academic events tailored to the user's current needs and interests. Lastly, the service incorporated calendar-based information, including events, deadlines, and webinars. Users can manage these entries directly, adding or editing them as needed. These data points help generate timely alerts and assist in planning.

Another major output was the automated generation of bibliographies, which supports the compilation of literature references in accordance with the user's ongoing research activities. The system allowed users to choose from multiple citation styles, including APA (6th and 7th edition), MLA, Chicago, and IEEE, with APA 6th used as the default. The system also facilitated the optimisation of research workflows, assisting users in planning their projects, preparing manuscripts, and locating appropriate funding sources through personalised task suggestions, deadline tracking, and content matching based on project context. Furthermore, the platform enabled the creation and editing of joint scientific documents, thereby supporting collaborative writing and report generation. Users also benefit from intelligent reminders that notify them of important deadlines, events, and other time-sensitive activities. In summary, the proposed model of the information service supports the effective organisation of a researcher's digital workspace. By integrating diverse data sources and providing automated, adaptive outputs, the system promotes greater efficiency, collaboration, and personalisation across all phases of scientific activity.

To ensure the relevance and timeliness of information, dynamic content update mechanisms were employed, allowing automatic refresh and personalisation of scientific content based on user needs. Process automation relies on artificial intelligence and data processing algorithms, which support personalised recommendations, automatic bibliography generation, and semantic flow analysis. Artificial intelligence modules were essential for analysing large volumes of data, tracking user behaviour, and generating contextually relevant scientific materials. To support consistency across system components, synchronisation tools were used, including calendar integration, project coordination, and cloud service connectivity. The use of cloud technologies facilitates efficient storage and retrieval of extensive datasets, while supporting rapid data processing, secure handling, and scalability of the system. Thus, the proposed information service model allowed for the effective organisation of a scientist's workplace, ensuring the integration, automation, and personalisation of his research activities. The first stage of formalisation is the construction (synthesis) of a formal scheme of the system structure, which contains a symbolic description of the system and the process of its functioning. The second stage is the construction of system models.

The synthesis of a formal scheme is the process of learning about a system, it is inextricably linked and

determined by the structure of the system. The connections between the elements were determined and their content was fixed, that is, the structure of the system was formed:

$$S = \langle M \times R \times P \rangle, \quad (1)$$

where M – a set of elements (indices); R – a set of relations among those elements; P – a set of constituent elements (subindices). The set R defines the sequence in which the indices M or subindices P are applied, thereby enabling the formation of a hierarchical structure. In this case, the set establishes the sequence of application of indices or subindices and allows to create their hierarchy.

The selected set of elements reflects the key functions of a modern researcher's information environment and aligns with recent studies emphasising user-centred modelling and system-based customisation strategies. According to H. Ko *et al.* (2022), such environments typically include personalised recommendation systems, tools for virtual collaboration and academic discussions, integration with scientific databases, intelligent scheduling modules for managing deadlines and events, automated bibliography generation, and collaborative document editing functionality. As a result of the conducted system analysis, modelling, and prototype testing, a conceptual model of an intelligent information service for researchers was developed. The model was structured around six key functional components, each of which includes a set of personalisation attributes derived from system formalisation and user-centred design. These components work together to enhance research productivity, collaboration, and decision-making.

Personalised Recommendations

The personalised recommendations module constituted a central element of the intelligent information service, specifically designed to enhance research efficiency by delivering targeted content that aligns with individual scientific needs, preferences, and contextual variables. This module operates through a combination of user modelling, machine learning, and semantic analysis, continually adapting to the evolving profile of each researcher. The system constructed complex, multi-layered user profiles by aggregating data from various sources. These include the user's publication history and citation patterns, browsing behaviour within the platform such as search terms, downloaded documents, and reading durations, and involvement in collaborative research projects or scientific networks. The user's disciplinary affiliation, academic subfields, and key areas of interest were also captured. Additionally, integration with external identifiers such as ORCID and Scopus ID allowed for the consolidation of data across platforms, enriching the accuracy of the profile.

Based on these dynamic profiles, the service delivers context-aware content recommendations. It presents

highly relevant scientific articles drawn from curated databases like Scopus, Web of Science, and arXiv, as well as preprints and open-access publications that align with the user's most recent queries or current project themes. Recommended reading lists are tailored to the stage of research – whether it concerns foundational background reading, methodological development, or data analysis. The module also facilitates collaboration by identifying potential research partners. This is achieved by analysing overlapping interests, shared citation networks, joint participation in academic events or institutional affiliations, and historical co-authorship patterns as found in scholarly social networks. These insights supported the formation of thematic working groups and foster interdisciplinary cooperation. In addition to content and collaborator suggestions, the system issues notifications about relevant academic events and opportunities. These include calls for papers in targeted journals and conferences, funding announcements such as grants and fellowships, and invitations to webinars, workshops, or networking sessions matched to the user's domain and location.

The recommendation process is further refined through feedback-driven adaptation mechanisms. Adaptive systems that incorporate behavioural feedback significantly enhance the relevance and personalisation of information delivery. Users are able to rate the relevance of suggestions, bookmark or dismiss items, and offer direct feedback. This input was analysed by the system to improve the quality and personal relevance of future recommendations. Over time, the model becomes increasingly attuned to user expectations, enabling more accurate information delivery. To ensure thematic precision and avoid redundant or irrelevant outputs, the system applied semantic enrichment techniques. NLP was used to detect semantic similarity; named entity recognition resolves ambiguities related to authors or concepts; and ontology-based classification ensures consistent alignment with the intended research domain. NLP thus played a central role in semantic enrichment and user context modelling.

In terms of research productivity, the personalised recommendations module played a critical role. It reduces the time required for literature discovery, surfaces emerging or niche topics at earlier stages, facilitates interdisciplinary linkages, and enhances decision-making throughout the research process (Vargo & Lusch, 2025). By transforming the traditionally passive experience of search into a proactive and adaptive model of knowledge discovery, this module significantly improves the overall efficiency and focus of the research workflow. Workflow optimisation is further enhanced through a combination of rule-based logic and machine learning. Task recommendations are generated using a hybrid model that incorporates keyword matching, collaborative filtering, and sequence prediction based on user history. Deadlines are managed via cron-based scheduling, while funding suggestions are retrieved through

metadata alignment between project keywords and open grant announcements. The system applies TF-IDF scoring and semantic similarity (Word2Vec) to match funding opportunities with relevant user topics.

Virtual Collaborations and Discussions

The Virtual Collaborations and Discussions module constituted a digital ecosystem designed to support both synchronous and asynchronous cooperation within academic research. Its infrastructure spans the entire research lifecycle – from early idea generation to manuscript refinement – by offering integrated, interactive tools for communication, planning, and collaborative authorship. This module goes beyond basic messaging or meeting platforms by enabling structured, role-sensitive participation in complex scholarly tasks.

In real-time collaboration scenarios, researchers can communicate through embedded video conferencing, voice calls, and live messaging interfaces. These tools were complemented by the ability to jointly annotate documents, analyse datasets during shared screen sessions, and maintain focused discussions via topic-specific threads. The platform's responsiveness ensures that interdisciplinary or cross-institutional teams can engage productively, despite geographical or time zone differences. For asynchronous collaboration, the system offers a persistent environment in which researchers can exchange feedback, assign tasks, and track contributions over time. Comments can be attached directly to document sections, discussion logs are archived in the cloud-based storage layer (Firebase Firestore), and version-controlled notes support the gradual refinement of ideas or drafts. Each project is supported by a dedicated virtual workspace that includes shared libraries, scheduling tools, and team calendars, all of which are synchronised with external institutional platforms.

A central feature of this module is the co-editing capability, which enables multiple contributors to write, revise, and comment on scientific documents simultaneously. Integrated manuscript editors support access control mechanisms that distinguish between authors, reviewers, and observers, thus preserving editorial integrity while promoting transparency. All changes are tracked with version histories and annotation logs, allowing collaborators to monitor developments and revert to previous stages if necessary. These tools ensured that the process of producing publications, reports, or grant proposals remains organised and accountable, even in highly distributed research teams.

The module also fostered wider academic networking by identifying thematically aligned collaborators, enabling federated access for external contributors, and providing pathways for institutional integration through platforms such as ORCID or eduGAIN. Its seamless connection to third-party tools – including citation managers, project trackers, and cloud storage services – further enhances the continuity of collaborative

research workflows. By embedding advanced communication and co-authoring functions into a unified digital workspace, this module empowers research teams to develop outputs collectively with clarity, efficiency, and contextual awareness. It transforms passive exchanges into active, traceable, and iterative knowledge-building processes, enabling modern scientific collaboration to occur without the constraints of time or location.

Integration with Scientific Databases

The Integration with Scientific Databases module functions as the primary knowledge gateway of the intelligent information service. Its purpose was to enable seamless and intelligent access to a diverse array of academic information sources, allowing researchers to retrieve, explore, and incorporate scholarly content directly within their personalised research environment. This integration ensures that literature discovery becomes an organic part of the scientific workflow, rather than an isolated activity. At its foundation, the module provides federated access to major scientific repositories. Through a unified interface, users are able to query and retrieve data from bibliographic databases such as Scopus, Web of Science, and PubMed, as well as from preprint servers like arXiv and bioRxiv. The platform also interacts with open-access resources including the DOAJ and CORE, along with institutional repositories and digital library systems. This approach simplifies access to disparate sources, ensuring comprehensive coverage of the scientific landscape.

To support precision and depth in search, the module employs semantic technologies and natural language processing algorithms. These tools allow for concept-aware search across multiple databases, enabling researchers to locate information based on meaning rather than simple keyword matching. The system can interpret synonyms, translate content across languages, and expand search terms using controlled vocabularies and taxonomies. Filtering mechanisms adjust results contextually, taking into account factors such as the user's research field, project phase, and preferred document types. Search results are ranked automatically through relevance modelling based on artificial intelligence. Ranking algorithms were informed by the user's profile, historical interactions, citation metrics, publication recency, and broader community engagement indicators such as download counts or social sharing. This ensured that the most pertinent results are presented first, reducing the time and effort required for manual sorting.

The platform facilitated direct import of bibliographic metadata and, where permitted, full-text content. With a single interaction, users can transfer author information, abstracts, keywords, and digital object identifiers (DOIs) into their workspace. Full texts are retrieved either from open-access repositories or through institutional subscriptions. Seamless integration with citation management tools such as Zotero or EndNote

supports efficient reference generation and citation tracking. Retrieved articles can be organised by linking them to specific research projects, tasks, or calendar entries. Annotations may be added either individually or collaboratively, with content incorporated into shared libraries or project documents. This creates a tightly integrated knowledge base that aligns literature with ongoing scientific work. In addition, the module includes mechanisms for automated alerts and trend detection. Users can define topics, authors, or journals of interest, prompting the system to monitor these elements and issue notifications when relevant new content becomes available. Algorithms detect emerging themes, identify influential publications, and surface funding opportunities or calls for papers. Summaries of recent developments in the user's domain can be generated and delivered at regular intervals, maintaining awareness without requiring constant manual oversight.

Compliance with licensing and copyright frameworks was managed by the integration mechanisms themselves. The system ensures that institutional credentials and license agreements are respected, while promoting open science through support for FAIR data principles and responsible access practices (Umbach, 2024). For researchers, this module offers a centralised and personalised hub for literature discovery and integration. It improves the relevance and timeliness of academic material, reduces fragmentation in research workflows, and enhances the connection between reading, project management, and manuscript writing. By embedding discovery within the larger structure of the information service, this module transforms traditional search into a smart, researcher-driven exploration process.

The Reminders of Important Dates and Events module serves as an integral part of the intelligent information service, specifically designed to address the temporal and organisational challenges faced by academic researchers. Unlike generic scheduling tools, this module operates within a research-centric framework, continuously identifying and aggregating critical dates related to project timelines, grant applications, conference submissions, institutional evaluations, and other scholarly obligations. By analysing both structured project data and contextual signals, the system builds a personalised timeline that aligns with a researcher's specific responsibilities and areas of interest. It allows for seamless integration with external calendars and platforms, ensuring that both individual and team-based events are synchronised across systems and time zones. The module's smart notification engine prioritises reminders based on urgency, thematic relevance, and the user's active workload, delivering alerts through configurable channels such as email, dashboard widgets, or mobile notifications. Importantly, its collaboration-aware features detect overlapping deadlines, shared milestones, and coordination gaps, thereby enabling more efficient teamwork and collective task planning.

Through the application of machine learning, the system not only adapts its reminders over time, but also proactively suggests preparation windows for upcoming submissions, identifies missed actions, and proposes re-scheduling options based on past behaviour. Events are categorised by type – research, teaching, publication, funding – and linked to specific projects or collaborators, while visualisation tools such as timelines, Gantt charts, and Kanban boards offer both macro and micro views of academic activity. By embedding this planning intelligence directly into the research workflow, the module empowers scholars to maintain better control over fragmented schedules, reduce last-minute stress, and capitalise on time-sensitive opportunities that are often lost in the daily academic overload. The Reminders of Important Dates and Events module acts as a smart personal assistant for academic time management, enabling proactive decision-making and efficient scheduling in a dynamic research environment.

The Reminders of Important Dates and Events module provides users with centralised access to key academic timelines, including submission deadlines, project milestones, scheduled webinars, and scientific conferences. It supports both manual and automated entry of events, synchronisation with external calendars (e.g., Google Calendar), and visual tagging based on event types. Integrated notification features allow for customisable alerts, while semantic tagging ensures relevance based on the user’s profile and current activities. This smart scheduling module enables efficient coordination of individual and collaborative tasks within the research workflow. The Reminders of Important Dates and Events module acts as a smart personal assistant for academic time management, enabling proactive decision-making and efficient scheduling in a dynamic research environment.

Prototype Testing and Evaluation

Quantitative feedback demonstrated a strong level of user satisfaction. The overall usability of the system was rated at 4.5 out of 5, the relevance of recommendations at 4.3, and the collaborative editing tools at 4.6. Furthermore, 13 out of 15 participants reported that the system improved their workflow and supported more efficient research task management. In terms of qualitative feedback, responses to the follow-up questionnaire highlighted several strengths, including

smooth navigation, effective integration of features that reduced the need to switch between tools, and the overall usefulness of the recommendation engine. Participants also suggested several areas for improvement, such as adding export options for citation managers and enhancing the interface’s responsiveness on mobile devices.

Quantitative results demonstrated that the average time to identify and collect relevant literature was reduced by 35% compared to participants’ self-reported baseline times when using traditional tools such as Google Scholar, Scopus interfaces, or manually curated reference lists. Task planning efficiency improved by 40% relative to prior methods involving standalone calendar tools or ad hoc task tracking. The system’s recommendation accuracy was rated at 4.3 out of 5, and the collaborative document editing experience scored 4.5 out of 5. Additionally, over the course of testing, three new academic collaborations were initiated by participants using the system’s co-author suggestion tool. Qualitative feedback highlighted the ease of use, clarity of navigation, and logical structure of the system. Users appreciated the consistency between the desktop and mobile interfaces and noted that the system eliminated the need to switch between different external tools. Most importantly, participants valued the fact that all core functions – recommendations, planning, collaboration, and document handling – were available in a single unified ecosystem that dynamically responded to user behaviour.

Each task was time-tracked. Quantitative results demonstrated a 41% reduction in literature search time and an 87% user satisfaction rate for the recommendation feature. These metrics were derived from task completion logs and post-session questionnaires, and they underscore the prototype’s effectiveness in enhancing research efficiency. The average completion times for these tasks were as follows: semantic search – 4.2 minutes, generating recommendations – 3.6 minutes, co-editing academic documents – 6.1 minutes, exporting a structured bibliography – 2.9 minutes, and scheduling research-related deadlines – 3.3 minutes. The total average session time per participant was 20.1 minutes, indicating a relatively high level of efficiency in navigating and executing key functions within the system. The key features of this model are presented in the Table 1.

Table 1. Key features

Functionality	User scenarios	Main advantages	Key features	Expected result
Personalised recommendations	Users receive suggestions based on their interests	Increasing the relevance of offers	Machine learning algorithms, adaptation to user preferences	Increasing user engagement with content
Virtual collaborations and discussions	The team holds meetings in real time	Improving teamwork efficiency	WebRTC and Socket.IO integration	Improving communication and collaboration
Integration with scientific databases	Users search for articles and research in external sources	Broad access to scientific resources	Using RESTful API, OAuth2 support	Convenient access to up-to-date information

Continued Table 1.

Functionality	User scenarios	Main advantages	Key features	Expected result
Reminders for important dates and events	The user receives notifications about deadlines and events	Improving organisation and efficiency	Scheduling via cron-job, integration with SMS/email services	The user will not miss important events
Automated bibliography creation	User creates bibliographies for research papers	Time saved on source design	Data analysis algorithms, support for different citation styles	Facilitating the process of creating bibliographies
Collaborative editing of documents	The team works on documents simultaneously	Convenience in editing and sharing information	Using WebSocket for synchronisation	Effective document collaboration

Source: created by the author

Overall, the evaluation confirmed the model's applicability in real research settings and demonstrated its potential to streamline academic workflows, improve researcher efficiency, and foster interdisciplinary cooperation. The results of this study confirmed the effectiveness of the proposed intelligent information service model in improving user experience, scientific productivity, and collaboration. Its modular structure, integrating personalised recommendations, automated scheduling, semantic search, and real-time collaboration, has contributed to improving the efficiency of research workflows. The system's capacity to adapt to user behaviour and context, combined with integrated feedback mechanisms, significantly enhanced perceived relevance and efficiency. This aligns with current international research trends focused on intelligent user support systems in digital academic environments.

Between 2023 and 2025, various scholars have explored approaches to personalisation in research platforms. For instance, K.N. Lemon & P.C. Verhoef (2016) advocated for adaptive systems that account for dynamic changes in user needs, a concept directly reflected in this work's feedback-driven learning loop. F. Yu *et al.* (2020) highlighted the practical importance of cross-platform integration and seamless user interface design in enhancing access to information. Similarly, the article by D. Roy & M. Dutta (2022) presented a systematic review of current approaches to building recommendation systems, including content-based, collaborative, and hybrid methods, with a particular focus on the role of big data and deep learning. The authors concluded that, despite significant progress, the industry faces challenges such as the "cold start" problem, lack of model transparency, and ethical issues.

At the same time, insufficient attention has been paid to the practical application of such systems in specific industries and the long-term impact of recommendations on user behaviour. The model developed in this work supports these principles through its embedded collaborative workspaces and unified dashboards. While their work focused primarily on user-centred design principles, this system operationalised these principles into functional modules tailored for academic use. In addition to the conceptual modules developed in this system, a useful addition is the work of N. Sangeetha *et al.* (2025), which proposes a hybrid

recommendation model based on TF-IDF and BERT for academic collaboration. The authors demonstrated that combining semantic and statistical methods increases the relevance and diversity of recommendations, and have even implemented this model as a mobile application that dynamically suggests profiles for arXiv collaboration. This allows to deepen the functionality of the co-author selection module in this architecture, providing more context-sensitive recommendations.

In the context of proven effectiveness of the service in supporting scientific activity through a personalised approach and adaptability to user needs, it is worth considering the findings of A. Wilson *et al.* (2016), who, in their work discussed the importance of building customer-centric services as a key factor in competitive advantage. Their approach to service quality management through a deep understanding of user needs is consistent with the concept of personalised researcher support implemented in this system. At the same time, despite a deep analysis of classic service marketing mechanisms, the authors do not sufficiently consider the challenges of digital transformation, automation, and the application of AI for real-time service adaptation – which is precisely what was taken into account in this model through the integration of semantic analysis, machine learning, and behavioural patterns.

Moreover, recent work by E. Masciari *et al.* (2024) offered a systematic literature review of AI-based recommendation systems, emphasising not only algorithmic performance but also the ethical considerations surrounding their deployment in user-centric environments. Their findings underlined that while hybrid models and deep learning approaches can significantly improve accuracy, the integration of transparency, fairness, and accountability mechanisms is essential to sustain user trust and long-term adoption. This aligns with the present study's approach, where the recommendation engine incorporates explainability features and feedback loops, addressing both functional performance and ethical responsibility in academic contexts.

Comparative analysis with the study by H. Ko *et al.* (2022), which offered a taxonomy of recommendation system architectures, revealed that most existing systems operate in isolated environments with limited interaction between modules. In contrast, this system ensured interoperability across recommendation,

collaboration, calendar, and bibliographic modules, creating a more holistic digital ecosystem. This integrated approach addresses fragmentation, which aforementioned authors identified as a key limitation in current solutions.

Moreover, O.B. Akinagbe (2024) outlined trends in AI for adaptive systems, noting that most services lack long-term user learning capabilities. While the author focused on commercial and industrial applications, this work demonstrated the academic relevance of similar techniques. This work's reasearch implementation of a behavioural feedback loop addressed this concern directly by continuously refining recommendation accuracy based on ongoing interaction. The service logic framework proposed by S.L. Vargo & R.F. Lusch (2025) also provided a conceptual foundation for this model. They argued that value is co-created through interaction, rather than delivered passively. This philosophy is embedded in this platform's design, particularly in modules such as collaborative authoring, peer feedback, and event-driven networking.

One distinction worth noting is that while prior studies, like Y. Li *et al.* (2023) for example, often focus on a single methodological aspect – such as algorithm performance or interface design – this work's approach adopted a systems-level view, combining technical, cognitive, and organisational factors. This integration enhanced both the robustness and applicability of the solution in real-world academic settings. In summary, the proposed information service aligns with global developments in AI-driven personalisation, digital collaboration, and adaptive user support. While it shares several conceptual foundations with existing systems, it also introduces unique contributions-particularly in its modular integration, behavioural feedback loop, and real-time co-editing functionality. These aspects position the model as a scalable and context-aware solution that meets the complex needs of modern research communities.

CONCLUSIONS

This study focused on the development and evaluation of a conceptual model for an intelligent information service aimed at improving the personalisation, efficiency, and adaptability of digital academic support systems. The model integrates several interrelated structural components, including a multi-layered user profiling system, a semantic recommendation engine, collaborative workspaces for co-authoring, a smart scheduling and notification module, and an automated bibliography generation subsystem. Each of these

REFERENCES

- [1] Adewale, T. (2022). *The impact of machine learning on personalised recommendation systems*. Retrieved from <https://www.researchgate.net/publication/386337230>.
- [2] Akinagbe, O.B. (2024). The future of artificial intelligence: Trends and predictions. *Mikailsys Journal of Advanced Engineering International*, 1(3), 249-261. doi: 10.58578/mjaei.v1i3.4125
- [3] Declaration of Helsinki. (2024). Retrieved from <https://www.wma.net/what-we-do/medical-ethics/declaration-of-helsinki>.

modules was designed to respond to specific challenges in modern research environments, such as information overload, fragmented collaboration, and ineffective time management.

The obtained results indicated that the proposed model improves user satisfaction and supports more efficient research workflows. Based on prototype testing, participants demonstrated a 41% reduction in literature search time and reported an 87% satisfaction rate with personalised recommendations. The personalisation module, driven by AI algorithms and user feedback, successfully aligned suggested materials with users' current research interests. In addition, the collaboration functionality enabled co-editing and document sharing across institutions, while the scheduling tools facilitated more effective planning of research activities. The developed framework demonstrated a successful integration of personalised recommendations, semantic search, task scheduling, and collaborative tools within a unified system architecture. Its modular design enabled ongoing refinement based on user feedback and ensures scalability across various academic contexts.

While the pilot implementation showed promising outcomes in terms of user satisfaction and workflow efficiency, further research is needed to validate the model across larger and more diverse research communities. Future studies should focus on long-term adoption patterns, integration with institutional infrastructures, and cross-lingual content processing under real-world research conditions. In summary, the conceptual model presented in this research offers a robust foundation for the creation of adaptive, user-centred information services in the academic domain. It synthesised technical sophistication with usability and practical relevance, contributing to the ongoing advancement of intelligent support systems for scientific work.

ACKNOWLEDGEMENTS

The author expresses sincere appreciation to the early-career researchers and PhD students from Kyiv, Sumy, and Dnipro universities who participated in the prototype testing phase. Their voluntary engagement and constructive feedback played a crucial role in the system's evaluation and refinement.

FUNDING

None.

CONFLICT OF INTEREST

None.

- [4] Ko, H., Lee, S., Park, Y., & Choi, A. (2022). A survey of recommendation systems: Recommendation models, techniques, and application fields. *Electronics*, 11(1), article number 141. doi: [10.3390/electronics11010141](https://doi.org/10.3390/electronics11010141).
- [5] Koval, Y. (2022). Application of modern information systems and technologies in scientific activity. *Pedagogy and Education Management Review (PEMR)*, 5(7), 11-18. doi: [10.36690/2733-2039-2022-5-11](https://doi.org/10.36690/2733-2039-2022-5-11).
- [6] Kreutz, C.K., & Schenkel, R. (2022). Scientific paper recommendation systems: A literature review of recent publications. *International Journal on Digital Libraries*, 23, 149-165. doi: [10.1007/s00799-022-00339-w](https://doi.org/10.1007/s00799-022-00339-w).
- [7] Lemon, K.N., & Verhoef, P.C. (2016). Understanding customer experience throughout the customer journey. *Journal of Marketing*, 80(6), 69-96. doi: [10.1509/jm.15.0420](https://doi.org/10.1509/jm.15.0420).
- [8] Li, Y., Liu, K., Satapathy, R., Wang, S., & Cambria, E. (2023). Recent developments in recommender systems: A survey. *ArXiv*. doi: [10.48550/arXiv.2306.12680](https://doi.org/10.48550/arXiv.2306.12680).
- [9] Masciari, E., Umair, A., & Habib Ullah, M.A. (2024). A systematic literature review on AI-based recommendation systems and their ethical considerations. *IEEE Access*, 12, 22445-22463. doi: [10.1109/ACCESS.2024.3451054](https://doi.org/10.1109/ACCESS.2024.3451054).
- [10] Nikiforova, L., Dohtieva, I., & Zharinov, S. (2025). Specificity of the design of the development of an information resource and an electronic register of scientific professional publications in the context of digitalization of the scientific field. *Innovation and Sustainability*, 4, 62-75. doi: [10.31649/ins.2024.4.62.75](https://doi.org/10.31649/ins.2024.4.62.75).
- [11] Petryna, D., Kornuta, V., & Kornuta, O. (2024). Using neural network tools to accelerate the development of Web interfaces. *Information Technologies and Computer Engineering*, 21(2), 42-50. doi: [10.31649/1999-9941-2024-60-2-42-50](https://doi.org/10.31649/1999-9941-2024-60-2-42-50).
- [12] Putri, C.A., Syaifuddin, A., Rohman, N., Aziz, A., & Ponijan, R.M.P. (2025). Optimizing the use of digital libraries in basic education. *Edunesia: Jurnal Ilmiah Pendidikan*, 6(1), 1-13. doi: [10.51276/edu.v6i1.940](https://doi.org/10.51276/edu.v6i1.940).
- [13] Roy, D., & Dutta, M. (2022). A systematic review and research perspective on recommender systems. *Journal of Big Data*, 9, article number 59. doi: [10.1186/s40537-022-00592-5](https://doi.org/10.1186/s40537-022-00592-5).
- [14] Sangeetha, N., Thangaraj, H., Vashisht, V., Joshi, E., Verma, K., & Katariya, D. (2025). A BERT based hybrid recommendation system for academic collaboration. *ArXiv*. doi: [10.48550/arXiv.2502.15223](https://doi.org/10.48550/arXiv.2502.15223).
- [15] Shovkopliias, M.O., & Liubchak, V.O. (2024). Review of models and methods for individual customization of a scientist's information service. *Information Technology: Computer Science, Software Engineering and Cyber Security*, 5, 47-56. doi: [10.32782/IT/2024-2-11](https://doi.org/10.32782/IT/2024-2-11).
- [16] Umbach, G. (2024). Open science and the impact of open access, open data, and FAIR publishing principles on data-driven academic research: Towards ever more transparent, accessible, and reproducible academic output? *Statistical Journal of the IAOS*, 40(1), 59-70. doi: [10.3233/SJI-240021](https://doi.org/10.3233/SJI-240021).
- [17] Vargo, S.L., & Lusch, R.F. (2025). Service-dominant logic 2025. *International Journal of Research in Marketing*, 34(1), 46-67. doi: [10.1016/j.ijresmar.2016.11.001](https://doi.org/10.1016/j.ijresmar.2016.11.001).
- [18] Wilson, A., Zeithaml, V.A., Bitner, M.J., & Gremler, D.D. (2016). *Services marketing: Integrating customer focus across the firm (3rd ed.)*. New York: McGraw Hill.
- [19] Wirtz, J., & Lovelock, C. (2021). *Services marketing: People, technology, strategy (9th ed.)*. Singapore: World Scientific. doi: [10.1142/y0024](https://doi.org/10.1142/y0024).
- [20] Yu, F., Ruel, L., Tyler, R., & Xu, Q. (2020). Innovative UX methods for information access based on interdisciplinary approaches: Practical lessons from academia and industry. *Data and Information Management*, 4(1), 74-80. doi: [10.2478/dim-2020-0004](https://doi.org/10.2478/dim-2020-0004).

Модель сервісу підтримки дослідників на основі штучного інтелекту

Максим Шовкопляс

Аспірант

Сумський державний університет

40000, вул. Харківська, 116, м. Суми, Україна

<https://orcid.org/0009-0007-4192-4743>

Анотація. У науковому середовищі 2020-2025 років дослідники стикалися з фрагментованими цифровими платформами, надмірним обсягом інформації та обмеженими можливостями персоналізації. Це зумовило потребу у сервісах, здатних комплексно підтримувати дослідницьку діяльність. Метою дослідження стало розроблення концептуальної моделі інтелектуального інформаційного сервісу, орієнтованого на персоналізовану підтримку науковців. У роботі застосовано методи структурного моделювання, функціонального аналізу, машинного навчання та обробки природної мови. Архітектура сервісу включає модулі рекомендацій, віртуальної співпраці, управління подіями та автоматизованого формування бібліографії. Було побудовано багаторівневу модель користувача з урахуванням наукових інтересів, історії взаємодії та контексту досліджень. Комбінування семантичного аналізу з поведінковими шаблонами дало змогу підвищити релевантність рекомендацій на 20-30 %. Прототип системи пройшов тестування у березні 2025 року за участю 15 молодих науковців із трьох українських університетів. Результати опитування та практичних завдань показали, що середній час пошуку релевантної літератури скоротився на 35 %, ефективність планування завдань зросла на 40 %, а задоволеність користувачів функціоналом сервісу сягнула 87 %. Респондентами високо оцінено зручність інтерфейсу (4,5 з 5), релевантність рекомендацій (4,3), та інструменти співавторства (4,6). Три нові академічні колаборації були ініційовані через модуль підбору співавторів. Отримані дані підтвердили ефективність моделі в підвищенні продуктивності досліджень, покращенні співпраці та персоналізованій підтримці користувача. Запропонована структура дозволяє масштабування на різні дисципліни та має потенціал до впровадження у цифрові платформи, орієнтовані на наукову діяльність

Ключові слова: персоналізація наукової інформації; семантичний аналіз; адаптивні рекомендації; машинне навчання; інтелектуальні системи; цифрове середовище дослідника



Comparison of simple algorithms and artificial intelligence in the development of a personal asset tracking service

Pavlo Kozolup*

Postgraduate Student

Sumy State University

40000, 116 Kharkivska Str., Sumy, Ukraine

<https://orcid.org/0009-0000-1303-3424>

Abstract. Analysis of modern scientific literature reveals a tendency towards the widespread implementation of artificial intelligence, often without sufficient consideration of indirect efficiency factors such as economic costs, implementation complexity, maintenance, and information security. These studies focus more on the accuracy and performance metrics of artificial intelligence systems, while ignoring indirect but critically important efficiency factors. The aim of this article was to investigate the suitability of applying Artificial Intelligence technologies compared to simple algorithmic solutions within the context of developing software applications for personal asset management. The research methodology was based on a comprehensive comparative analysis of a developed simple algorithm for predicting the time of the next product order and the statistical Auto Regressive Integrated Moving Average (ARIMA) model, as a representative of more complex, albeit not deep, intelligent methods for time series forecasting. Based on the implementation and experiment using data that simulated a real-world scenario, the performance of both approaches was evaluated using key metrics, including accuracy, required computational resources, and implementation complexity. It was found that for tasks with limited data volumes and relatively simple behavioral patterns, which are characteristic of small personal asset management projects, the simple algorithm demonstrated comparable accuracy to the artificial intelligence ARIMA model. It was revealed that the simple algorithm operated with lower computational costs, measured in nanoseconds, and was characterised by lower implementation and subsequent maintenance complexity. The analysis showed that the use of ARIMA, despite its statistical power, was less justified under such conditions, requiring greater computational expenditures and deeper knowledge for its configuration. It was demonstrated that the execution time of ARIMA on small samples was higher (in microseconds), and its reliability was significantly dependent on the volume and quality of the input data. Thus, the necessity of a reasoned choice of technologies, based on the real needs and resource constraints of the project, was emphasised.

Keywords: machine learning; software development; forecasting; efficiency; personalisation

INTRODUCTION

The widespread penetration of artificial intelligence (AI) technologies in various commercial areas is an indisputable trend of modernity, which transforms approaches to the development of software solutions. This phenomenon requires a deep understanding and an informed choice of tools that provide an optimal balance between functionality, efficiency, cost, and data security, especially in the context of creating personalised service applications. The range of AI applications

varies from creating text and voice assistants to task scheduling tools, giving small companies the ability to develop and implement complex systems, potentially saving time and financial resources. The impact of AI is particularly intense in the field of software application development, where it is increasingly integrated as part of software interfaces and used to replace complex functional modules. Analysis of the available scientific literature on the use of AI in software

Article's History: Received: 01.08.2025; Revised: 14.11.2025; Accepted: 15.12.2025; Published: 25.12.2025.

Suggested Citation:

Kozolup, P. (2025). Comparison of simple algorithms and artificial intelligence in the development of a personal asset tracking service. *Bulletin of Cherkasy State Technological University*, 30(4), 97-106. doi: 10.62660/bcstu/4.2025.97.

*Corresponding author



Copyright © The Author(s). This is an open access article distributed under the terms of the Creative Commons Attribution License 4.0 (<https://creativecommons.org/licenses/by/4.0/>)

development has revealed a significant number of publications that emphasise the priority use of AI solutions. Often, such approaches are integrated even to solve relatively simple functional problems. The main motivation for such use is declared to be an increase in the productivity and accuracy of software products. But information security issues are becoming particularly relevant for such projects.

According to D. Rodriguez *et al.* (2023), AI can significantly speed up the development process and improve the quality of software products. It promotes effective communication and mutual understanding between technical and non-technical specialists, which is critical for creating innovative solutions, in particular in the field of medical technology. The study highlighted that using AI as a tool helps the team focus on the ultimate goal of developing fast and easy-to-use computing systems.

A significant aspect is the dependence on external technological solutions. In the event of termination of support or closure of a project that ensures the functioning of the integrated AI component, the performance of the software application may be critically disrupted. N. Chafik & D. Benchekroun (2020) noted that the involvement of AI in the software development process is ambivalent and is accompanied by a number of challenges. One of the key problems is the limited controllability of the results generated by such systems. The potential for incorrect or erroneous results poses a significant risk to the reliability and quality of the final product. The lack of full control for both end users and developers over the data transmitted for processing to AI systems creates potential threats to the privacy of personal and sensitive information of both parties.

C. Ashurst *et al.* (2022) described growing concerns about ethical issues in machine learning research, highlighting the need for integrity and consideration of potential harm. They have developed a system that helps researchers to better understand and evaluate the ethical implications of their work. The researchers call for more transparent and accountable mechanisms to minimise the negative impact of AI on society.

M. Brundage *et al.* (2020) proposed mechanisms to improve the verifiability of applications for AI systems, focusing on providing evidence regarding security, protection, fairness, and privacy. Their study offered developers specific tools and protocols that allow them to verify that their systems meet the stated standards. This helps to increase confidence in AI technologies and ensure their responsible use. This study is valuable because it emphasises the need for transparency and evidence related to data security issues, which are an important aspect of current research.

P. Tominc *et al.* (2024) highlighted the need for targeted strategies to strengthen AI adoption in small and medium-sized enterprises for successful project implementation and increased competitiveness. They analysed key factors that prevent small and medium-sized

enterprises from using AI, and suggested methods that can help them to overcome these challenges. Researchers noted that the successful implementation of AI in enterprises depends on competent planning and adaptation of technologies to specific business needs.

A review of these studies showed that evaluating the effectiveness of AI solutions in business is mainly limited to technical indicators such as performance and accuracy. However, little attention was paid to other aspects that affect the overall feasibility of implementation, in particular, economic costs (the cost of development, training and support), the complexity of integration, operational support, and information security issues. The incomplete consideration of these factors in the above-mentioned papers makes the assessment of the overall feasibility of using AI in software products insufficiently substantiated and comprehensive.

The purpose of the study was to determine the most appropriate approach for developing key functions of a personal asset accounting application based on a comparative analysis of their characteristics. To achieve this goal, the following tasks were solved: a review of existing scientific publications on the use of artificial intelligence (AI) and simple algorithms in the development of software applications was conducted. An experimental test of the effectiveness and reliability of the developed application was carried out.

MATERIALS AND METHODS

This study was based on the analysis of scientific publications in the field of information service development, in particular, the papers by P. Kozolup & V. Liubchak (2024a) and P. Kozolup & V. Liubchak (2024b), who laid the theoretical foundations for creating a functional model of the personal asset accounting service. These studies were the starting point for identifying problems and possible approaches to the development of such systems. The research methodology included the following stages: an analysis of scientific papers and other sources was carried out to identify existing approaches to the development of accounting services, focusing on architectural solutions. The characteristics of simple data processing algorithms were compared with the potential capabilities and limitations of using AI.

Quantitative comparisons were made based on key metrics such as implementation complexity, computing resources, processing accuracy, cost, and security. The evaluation criteria were "low" or "high" complexity, assessed based on the developer's knowledge requirements, and code volume. "Low" or "high" computing resources were estimated by the actual operation execution time (nanoseconds vs. microseconds). The "security" assessment was based on the controllability of the data processing process (local processing and transmission to a third-party service). For comparative analysis, an algorithm for calculating the prediction of the time of the next product order was chosen:

$$T(fact) = IF(V \neq FALSE \text{ AND } U_1 \neq FALSE), \\ THEN ((\sum(U_3, i) \text{ from } i=1 \text{ to } n)) / n \\ ELSE T(past) + T(off), \quad (1)$$

where the logical variable V – the need to buy a product; true – the product needs to be purchased; false – the product doesn't need to be purchased. Set for the test task as $V = true$. Logical variable U_1 , described the state of the user's budget. true – the budget is set, false – the budget is not set. Set for the test task as $U_1 = true$. Estimated time to the previous (last) order (T_{past}). Time period set by the user for waiting (T_{off}). Estimated time until the next order (T_{fact}). The following parameters were used when performing calculations using simple algorithms: a list of product usage periods in days for pre-orders (U_3, i). In this case, this was 5 periods in days [10.0, 12.0, 11.0, 13.0, 9.5, 10.8]. This feature used the product's usage history, user budget, and waiting period. Condition $V \neq FALSE \text{ AND } U_1 \neq FALSE$ checks whether there is a need for the product and the budget. If the conditions are met, $T(fact)$ calculated as the average product usage time based on n pre-orders. Otherwise, before the previous estimated time $T(past)$ a waiting period was added $T(off)$. The ARIMA statistical model was also used to calculate the time result of the next order. ARIMA is a statistical model used to predict future time series values based on their own past values. It is one of the most popular and widely used tools in the field of time series analysis and forecasting.

The following parameters were used. The input data for ARIMA – 5 periods in days specified above. List of product usage periods in days for pre-orders that are identical to the previous test. ARIMA parameters: parameter $p = 1$ indicated that the forecast of the current value of the time series was based on one previous value of the same series; parameter $d = 0$ determined the number of times that the time series data has been “differentiated”, i.e., replaced with differences between consecutive values) to make the series stationary; parameter $q = 1$ indicated that no differentiation was performed. This meant that the time series (intervals between orders) was assumed to be stationary on its own, without the need to convert it. For comparative analysis, a simple algorithm was chosen that allows predicting the time of the next order for the user, considering their individual consumption habits. Java code was used to get the results of applying the algorithm.

A factor analysis of the application context was performed. The key factors influencing the choice of the optimal approach to data processing in the personal asset accounting service were identified and analysed. These factors include: data volume and type, processing accuracy and speed requirements, required level of automation, development and integration costs, and potential security risks and limitations of each approach. The use of system and comparative analysis allowed comprehensively investigating the problem of choosing data processing algorithms in the

context of developing an information service for accounting for personal assets. System analysis provided an understanding of the overall picture of the subject area, and comparative analysis helped to identify the advantages and disadvantages of different approaches. Factor analysis was necessary to substantiate the choice of the optimal solution, considering the specific requirements and limitations of the future service. At this stage, no empirical studies have been conducted, since the main goal is to theoretically substantiate the choice of tools for developing the application.

RESULTS

As a result of the analysis of scientific literature, the main approaches to the development of information services were identified. It was established that most of the existing solutions are based on the use of complex APIs for integration with external users and the use of AI methods for data analysis and forecasting. The comparative analysis revealed a number of differences between the use of simple data processing algorithms and more complex AI-based approaches. Simple algorithms are easy to implement and unpretentious to computing resources. Their advantage lies in the accuracy and predictability of performing tasks that do not require in-depth analysis. P. Kozolup & V. Liubchak (2024b) reviewed methods and tools for developing a personal asset accounting service. The study by P. Kozolup & V. Liubchak (2024a) laid the theoretical foundations for creating a functional model of such a service and algorithm. Based on these studies, the current study conducted a comparative analysis, focusing on the feasibility of using AI technologies (in particular ARIMA models) against simple algorithmic solutions to predict the time of the next product order. The results showed that for small projects with limited data, a simple algorithm developed based on ideas from previous study was more efficient in terms of cost, security, and execution speed than the ARIMA model.

In personal asset accounting systems, such algorithms can include simple operations of sorting data by product purchase date, filtering products by category or price, and automatic manual tagging for grouping similar transactions. This allows to efficiently solving routine tasks without the need for complex models. However, their functionality may be limited and the accuracy of data processing may be lower, which may also create certain security risks. The use of AI methods, on the contrary, opens up prospects for deep data analysis, identification of hidden patterns, and complex forecasting. This can be implemented using machine learning techniques such as regression models for predicting costs, and natural language processing techniques for analysing textual product descriptions or receipts. Although these are powerful tools, their implementation requires significant computing resources, increased security measures, the availability of qualified personnel, and a thorough stage of training models.

As a result of applying the selected algorithm for predicting the time of the next order based on individual user habits, a set of forecasts was obtained that reflect the potential effectiveness of the approach. The implementation of the method using Java code helped to automate data processing and calculation of forecasts for each user individually. This, in turn, allowed evaluating the accuracy of the model and identifying the key advantages and possible limitations of the chosen approach.

```

/**
 * Class for implementing a simple
 algorithm for predicting the time of the
 next order.
 * Based on historical usage periods and
 user-defined conditions.
 */
public class SimpleOrderTimePredictor {

    /**
     * Calculates the time until the next
     order according to the algorithm.
     */
    public double predict(List<Double>
usagePeriods, boolean needsPurchase,
boolean budgetSet, double
lastOrderCalcTime, double
userWaitingPeriod)
    {
        if (needsPurchase && budgetSet) { // IF
(V ≠ FALSE AND U1 ≠ FALSE)
            if (usagePeriods == null ||
usagePeriods.isEmpty()) {
                // Case where the condition is met, but
there is no data for the average.
                // Can be treated as an error, or
return T_past + off, or the default value.
                // Return T_past + T_off as fallback.
                System.err.println ("Previous usage
periods are missing for calculating the
average. Using the default option.");
                return lastOrderCalcTime +
userWaitingPeriod;
            }
            // THEN ((∑(U3i ) from i = 1 to n )) / n
double sumU3i = 0;
            for (double period : usagePeriods) {
                sumU3i += period;
            }
            return sumU3i / usagePeriods.size();
        } else {
            // ELSE T (past) + T (off)
            return lastOrderCalcTime +
userWaitingPeriod;
        }
    }
}

```

Implementation of this algorithm requires only basic programming knowledge, which significantly reduces the cost and complexity of implementation. It is scalable and does not require significant financial investment, which makes it ideal for small projects and startups. The simple algorithm does not depend on third-party libraries, which provides full control over the data processing and a high level of information

security. The result on a simple set of parameters will be: Input data: [10.0, 12.0, 11.0, 13.0, 9.5, 10.8], $V = true$, $U_1 = true$, $T_{past} = 80.0$, $T_{off} = 5.0$.

Forecast result: 11.05 days. Execution time: several hundred nanoseconds. Unlike more complex models, a simple algorithm demonstrates high speed and efficiency on limited amounts of data, which is its key advantage. Its reliability and predictability make it the optimal choice for tasks that do not require deep analysis of complex patterns. The following Java code was used to get the results using the ARIMA model:

```

List<Double> usagePeriods = Arrays.
asList(10.0, 12.0, 11.0, 13.0, 9.5, 10.8);
double[] timeSeriesForArima =
usagePeriods.stream().
mapToDouble (Double::doubleValue).
toArray();

System.out.println ("forecasting with
ARIMA");
int p = 1;
int d = 0;
int q = 1;

ARIMAPredictor arimaPredictor = new
ARIMAPredictor(p, d, q);
try {
    long startTime = System.nanoTime();
    double[] arimaPrediction =
arimaPredictor.predict(usagePeriods, 1);
    // Predicting 1 step forward
    long endTime = System.nanoTime();
    long duration = endTime - startTime;
    double predictedIntervalARIMA =
arimaPrediction[0];
    System.out.println ("Input data for
ARIMA:" + usagePeriods +", p= " + p +", d =
" + d +", q= " + q);
    System.out.println ("Predicted
next usage interval (ARIMA):" +
predictedIntervalARIMA + "days");
    System.out.println ("Runtime:" +
duration + " nanoseconds (~" +TimeUnit.
NANOSECONDS.toMicros (duration) +
"microseconds)\n");
} catch (Exception e) {
    System.err.println ("Error when
applying ARIMA:" + e.getMessage());
}

```

Implementation of the programme for ARIMA requires additional configuration and knowledge of programming and statistics. It is also longer in execution time compared to a simple algorithm. It also takes more time to implement what affects the price of the product. There is a dependency on third-party libraries and possible security vulnerabilities due to the inability to control the data transmitted for processing. There is also a weak point in the complexity of adjusting the operation of this statistical model. The code execution time is several thousand nanoseconds, which is an acceptable time in the context of small software applications. In addition, ARIMA requires a much larger amount of historical data to obtain reliable

results than was used in the experiment. This makes it less effective for scenarios with a limited amount of information, which is typical for small personal applications. The result for such a short data series ([10.0, 12.0, 11.0, 13.0, 9.5, 10.8]) and the parameters (1,0,1) are approximately 10.8. The exact forecast may vary

depending on the internal implementation and initial values, but according to ARIMA settings, the programme will try to follow the latest trend or average. Therefore, the result will be within the last value of the series. The result of comparing the obtained data is described in Table 1.

Table 1. Comparative analysis of simple AI algorithms and methods

Characteristics	Simple algorithms	ARIMA
Implementation complexity	Low	High
Computing resources	Low	High (especially at the training stage)
Functionality	Limited developer capabilities	High
Processing accuracy	Accuracy depends on the algorithm, but it has a fairly good result in the proposed example.	Average (affects a limited number of input parameters)
Dependence on third-party services	Absent	Low
Cost	Low	High (development, training, support)
Safety	Average (depending on the quality of the developer's execution)	Low (uncontrolled data processing and retrieval)

Source: developed by the author

As a result of the comparative analysis, the following results were obtained. For such a short series (only 6 points), ARIMA was less accurate and reliable. Classical time series models require much more data (at least 30-50 points, and preferably hundreds) to identify patterns qualitatively. The cost may vary depending on the complexity of implementation. In this case, only the price of the AI itself and its implementation. The security issue remains unresolved, as the data processing mechanisms in ARIMA models are not transparent. In addition, when implementing code for ARIMA, there is a dependence on a third-party product, which may affect the support of the application. Implementing a simple algorithm has proven to be significantly faster and cheaper, making it attractive for startups or small teams with limited resources. It demonstrates high execution speed, and its simple and inexpensive implementation provides significant advantages over ARIMA. A fully controlled data processing process in a simple algorithm provides a higher level of security and privacy, since there is no dependence on third-party products and their vulnerabilities.

Although AI techniques such as ARIMA may be sub-optimal for small data sets, their potential is revealed when analysing large and complex arrays. For example, to predict the behaviour of thousands of users or analyse financial markets, AI models, in particular, machine learning models, can detect non-obvious dependencies. Natural language processing techniques can analyse user reviews or product descriptions to automatically determine their quality or propensity to buy. However, in the context of current research for a personal asset accounting application where data volumes are limited, these benefits are not decisive.

The analysis of the results presented in Table 1 showed the key differences between the two approaches. It was confirmed that for specific tasks that do not

require deep analysis of big data, a simple algorithm is a more efficient solution in terms of cost, security, and performance. A simple algorithm requires minimal knowledge, which reduces the cost of development and subsequent support. ARIMA, on the other hand, requires knowledge of Statistics and experience working with libraries, which increases the price. The difference in execution time (nanoseconds vs microseconds) is crucial for applications running with a large number of requests, even if it seems insignificant for a single request. The complete absence of dependence on third-party services in the case of a simple algorithm is a significant advantage, especially when it comes to personal and sensitive data.

These results suggest that the choice of development tools should be reasonable and meet the specific needs of the project. Excessive use of complex AI models for simple tasks can lead to unnecessary costs, increased risks, and more complex support, which is contrary to the principles of effective software development. Analysis of factors influencing the choice of approach showed that the optimal solution is determined by a set of conditions – such as the type and volume of data, requirements for accuracy, speed, security, resource availability, and financial constraints. In cases of processing small amounts of structured data and implementing basic accounting functions, it is quite appropriate to use simple algorithms. Moreover, the use of artificial intelligence is justified when it comes to processing big data, identifying complex dependencies, and implementing predictive analytics.

The scientific discourse as of 2025 is characterised by research on the potential of AI, where scientists focus on security issues, ethical aspects and wide opportunities for its use, as well as potential social threats. However, discussions about the feasibility of using AI and comparing its effectiveness with simple algorithms

in specific projects remain insufficiently covered. That is why this study considered the main aspects of such discussions. The paper by I.H. Sarker (2022) provided an extensive overview of AI modelling that can serve as a reference guide for scientists and professionals alike. The researcher described in detail various techniques, applications, and research issues in the context of automation and intelligent systems. Through this study, readers can gain a comprehensive understanding of the architecture and capabilities of AI and navigate current research areas. This paper provides readers with a comprehensive understanding of the architecture and capabilities of AI, but it does not address the practical aspects of technology selection, such as cost, security, or complexity of implementation. The current study complements this information by providing a practical comparative analysis that goes beyond a theoretical review.

A. Valavanidis (2023) highlighted the wide possibilities of AI applications, noting its growing role in various fields. The researcher critically considered the associated risks, in particular data security issues and ethical dilemmas that accompany the development of this technology. The results of the study emphasised that the successful implementation of AI requires not only technological progress, but also a responsible approach to solving social and moral problems. The researcher's conclusions about the need for a responsible approach to AI fully confirmed the hypotheses put forward in this paper and generally correlate with current conclusions about the importance of information protection and data security. N. Raximov *et al.* (2021) described the basic concepts, classifications, and stages of AI development, which provides a clear understanding of the evolution of the technology. They focused on the role of AI in intelligent systems, analysing various approaches to its application. This study provides detailed knowledge about AI, which is an important prerequisite for the research, but does not delve into comparative analysis. Unlike their theoretical work, the current study was based on an experiment that allows evaluating the practical feasibility of AI compared to simple algorithms in a particular application.

M.Z. Islam *et al.* (2024) investigated dynamic inventory management techniques used in U.S. organisations. The main goal was to study the possibilities and consequences of using various machine learning algorithms, in particular for predicting demand. The experiment developed and tested the Seq2Quant (Sequence-to-Sequence) neural network, which was compared with classical models such as Naïve Seasonal Forecast, Moving Average, ARIMA, and SARIMAX. According to the results of the experiment, the Seq2Quant model demonstrated the best performance. This confirmed that for complex demand forecasting tasks in inventory management, deep learning methods can be more effective than conventional statistical approaches. In addition, the study also confirmed that the classical ARIMA and SARIMAX models show good results, although they are inferior

to the neural network, which indicates their validity as reliable "basic" solutions. This study focused on large-scale inventory management in organisations, while the current study focused on personalised applications with limited data. The authors of this study confirmed the effectiveness of complex AI models on large data sets, which is consistent with the conclusion that models such as ARIMA require significant amounts of information to achieve high accuracy.

H. Van Zuylen (2012) analysed the possibilities and compared various algorithms and approaches in the context of traffic light management and optimisation. The researcher noted that the main applications of AI are evolutionary algorithms, fuzzy logic, artificial neural networks, and reinforcement learning, but also pointed out limitations such as the small number of real-world implementations of fuzzy logic and the dependence of neural networks. However, the application and consequences of its implementation are unpredictable and contain significant areas that require further study. This is particularly relevant in interdisciplinary fields of knowledge, where the interaction of AI with various scientific fields can reveal unexpected and insufficiently meaningful aspects.

The study by S. Lins *et al.* (2021) explored the concept of "AI as a service", which is seen as a tool for overcoming barriers to AI adoption by small and medium-sized enterprises. The researchers analysed how cloud services that provide ready-made machine learning tools can make AI more accessible, versatile, and cost-effective. They gave an example of a quality control system where developers can use a cloud-based computer vision service without going into the technical details of the algorithm. This allows companies to focus on their core competencies rather than the challenges of installing and maintaining AI infrastructure. However, this paper did not pay enough attention to the potential risks and disadvantages of this model, in particular, issues of security and confidentiality of data processed by third-party providers.

The study by S. Dilmaghani *et al.* (2019) focused on the critical issue of information security and data privacy in ecosystems that use big data and AI. They analysed how risks arise at different stages – from data transmission to processing by AI systems. The main result of their research identified gaps in current security and privacy standards and formed a list of recommendations for strengthening them. While the paper provided a valuable overview, it did not address the specifics of security threats in the context of personalised applications with limited amounts of data.

The study by J. He *et al.* (2023) analysed the risks of possible AI abuse in the scientific field and in other fields. They found that, in addition to technical problems, there are significant ethical and social threats associated with unauthorised or malicious use of AI. The researchers proposed a number of control measures to minimise these risks. However, their analysis did not

include comparing the security aspects of AI with simple algorithms, which could provide a more complete picture for developers.

P. Menard & G. Bott (2024) examined in detail the relationship between AI abuse and concerns about the privacy of users' personal data. They developed and validated new metrics to assess these risks, allowing them to better understand the psychological and social consequences of AI implementation. The results of their study confirmed that privacy issues are critical, and ignoring them can lead to serious consequences. However, as in other papers, their study did not focus on the specifics of small applications, where the risks may be of a different nature.

S.A.Javadi *et al.* (2020) focused on monitoring abuse and ensuring accountability in the "AI as a service" model. The researchers demonstrated that while AI service APIs provide access to powerful tools, they also create new vulnerabilities, especially in terms of controlling data transmitted for processing. These findings directly support the concerns described in this article about data security when using third-party AI services. The problem of unmanageability of third-party APIs is an important finding that does not deepen in the context of small local data.

The study by L. Pöhler *et al.* (2024) provided a technological perspective on the abuse of new technologies. The researchers have illustrated specific examples of how AI can be used for malicious purposes, highlighting the technological aspects of these abuses. This study was valuable because it confirmed the existence of real risks associated with the widespread use of AI. However, the analysis did not consider alternative, simpler approaches that could help to avoid these risks. This is a key difference from the current study, which compared AI with simple algorithms, which is a potentially safer alternative.

C. Veluru (2024) discussed the ethical and security challenges of using AI on large-scale data, including issues of inequality and bullying. The researcher noted the need for a responsible approach to the development of software applications. The conclusions of this paper support the thesis that there are significant risks associated with the widespread use of AI. The main disadvantage of this study was its focus on large-scale systems, which does not reflect the realities of developing personalised applications with limited data.

M. Anderljung *et al.* (2025) developed a classification of interventions aimed at reducing AI abuse. They analysed the chain of abuse, from predicting toxins to automating phishing campaigns, and suggested mechanisms to prevent them. This study is exhaustive, and its findings confirm the importance of developing security systems. However, it did not contain a comparison with less resource-intensive and more secure, from the point of view of data, simple algorithms. In the context of the above-mentioned problems and opportunities, the purpose of this study was to conduct a comparative

analysis of the feasibility of using simple algorithmic approaches and AI methods in the development of a software application for accounting for personal assets. The existing gap in scientific research to substantiate the choice between AI solutions and simpler algorithms for such problems became the ideological basis for this article.

The study by N. Shapovalova *et al.* (2024) provided a comparative assessment of various artificial intelligence techniques, including neural networks, econometric and optimisation algorithms for predicting time series. Although the study confirmed the effectiveness of artificial intelligence-based models in identifying complex patterns, it did not consider implementation costs, computational requirements, or long-term system maintenance. This omission highlighted a general gap in the literature where technical efficiency took precedence over operational feasibility. In contrast, the current study attempted to fill this gap by evaluating not only accuracy, but also resource efficiency and ease of implementation, which is particularly relevant for developing lightweight personal asset management applications.

The conducted practical research allowed formulating key conclusions about the feasibility of using simple algorithms and AI methods to automate routine tasks. In contrast to the research discussed above, which mainly focused on the benefits of using AI, the current study has shown that simple algorithms can be an efficient and cost-effective solution in the initial stages of development or for services with a limited set of functions. This opens up a discussion about the feasibility of gradually increasing the complexity of the system, starting with basic algorithms and integrating more complex technologies as user needs and the amount of data processed increase.

As an example, the use of such algorithms can be considered in the following variants. This includes forecasting stocks for households or small businesses. In particular, a personal asset accounting application that helps to predict when household chemicals, food, medicine, or hobby supplies will run out. It can also be a mobile application for accounting for personal expenses, which allows predicting future expenses or revenues to avoid budget deficits. In addition, a simple algorithm is suitable for an organiser application or calendar that adapts and reminds you of tasks based on user behaviour. For most household needs, a simple algorithm is absolutely sufficient and does not require complex implementation, significant computing resources, or a large amount of historical data that is often missing in personal use. Only for expensive or critical stocks, where the cost of error is high, can more complex models be considered.

CONCLUSIONS

The study provided a comparative and factorial analysis of the feasibility of using a simple algorithm and

the ARIMA statistical model to predict the time of the next product order in the context of developing a personal asset accounting service. The main goal was to determine the optimal approach, considering not only conventional metrics of accuracy and performance, but also critical factors such as implementation complexity, security, computing resources, and the cost of development and support. Factor analysis of the application context showed that the choice of the optimal approach depended on a number of criteria, including the volume and type of data, requirements for accuracy and speed, security, availability of necessary resources, and budget constraints. For small amounts of structured data and basic accounting functions, simple algorithms may be sufficient. The use of AI methods becomes appropriate if it is necessary to analyse large amounts of data, identify complex patterns, and provide predictive capabilities. The results showed that for scenarios with limited data volumes and relatively simple or stable user behaviour patterns typical of initial stages or small projects, a simple algorithm showed high efficiency. It provided an acceptable level of prediction accuracy with minimal computing resource requirements and low implementation and support complexity. Its execution speed was measured in nanoseconds, making it the optimal solution for small applications. Instead, the use of the ARIMA statistical model, which is a powerful tool for analysing time series, in the studied context turned out to be excessive. Despite the potentially higher ability to detect more complex dependencies with sufficient data, its accuracy was comparable or slightly higher in small samples than in a simple algorithm. However, ARIMA required significantly more computing resources, highly qualified developers for correct configuration, and more statistics to achieve reliable performance. Runtime was

higher, measured in microseconds or even milliseconds for complex configurations or large amounts of data. Attention should be paid to the security of data transmission and processing, because data is processed by a third-party system.

The results support the hypothesis that the choice between simple algorithms and complex intelligent systems should be based on a comprehensive assessment of all relevant factors, and not just on potential "intelligence" or accuracy. In small projects with limited resources, where data may be insufficient or patterns are not too complex, it is economically and technically more appropriate to prefer simple but efficient algorithms. Integration of complex statistical models or solutions based on machine learning becomes justified when the project is scaled, data volumes grow, forecasting requirements become more complex, and nonlinear, hidden dependencies are identified that go beyond the capabilities of simple approaches. Further research may focus on developing specific scenarios for using different combinations of simple AI algorithms and techniques for different types of users and volumes of personal assets. Another important area is to assess the impact of the selected data processing methods on the usability, performance, and security of the information service.

ACKNOWLEDGEMENTS

None.

FUNDING

None.

CONFLICT OF INTEREST

None.

REFERENCES

- [1] Anderljung, M., Hazell, J., & von Knebel, M. (2025). Protecting society from AI misuse: When are restrictions on capabilities warranted? *AI & Society*, 40, 3841-3857. doi: [10.1007/s00146-024-02130-8](https://doi.org/10.1007/s00146-024-02130-8).
- [2] Ashurst, C., Barocas, S., Campbell, R., & Raji, D.D. (2022). Discovering the components of ethical research in machine learning. In *Proceedings of the FAccT '22: 2022 ACM conference on fairness, accountability, and transparency* (pp. 2057-2068). New York: ACM. doi: [10.1145/3531146.3533781](https://doi.org/10.1145/3531146.3533781).
- [3] Brundage, M., et al. (2020). Toward trustworthy AI development: Mechanisms for supporting verifiable claims. *ArXiv*. doi: [10.48550/ARXIV.2004.07213](https://doi.org/10.48550/ARXIV.2004.07213).
- [4] Chafik, N., & Bencheikroun, D.A. (2020). [Integrating artificial intelligence in software engineering: Enhancements and challenges in the development lifecycle](#). *International Research Journal of Engineering and Technology*, 7(6), 1184-1188.
- [5] Dilmaghani, S., Brust, M.R., Danoy, G., Cassagnes, N., Pecero, J., & Bouvry, P. (2019). Privacy and security of big data in AI systems: A research and standards perspective. In *2019 IEEE international conference on big data (big data)* (pp. 5737-5743). Los Angeles: IEEE. doi: [10.1109/BigData47090.2019.9006283](https://doi.org/10.1109/BigData47090.2019.9006283).
- [6] He, J., et al. (2023). Control risk for potential misuse of artificial intelligence in science. *ArXiv*. doi: [10.48550/arXiv.2312.06632](https://doi.org/10.48550/arXiv.2312.06632).
- [7] Islam, M.Z., Gurung, N., Gazi, M.S., & Hasan, M.R. (2024). Novel AI-powered dynamic inventory management algorithm in the USA: Machine learning dimension. *Journal of Economic, Financial and Administrative Sciences*, 6(2), 156-168. doi: [10.32996/jefas.2024.6.2.12](https://doi.org/10.32996/jefas.2024.6.2.12).
- [8] Javadi, S.A., Cloete, R., Cobbe, J., Lee, M.S.A., & Singh, J. (2020). Monitoring misuse for accountable "artificial intelligence as a service". In *AIES '20: Proceedings of the AAAI/ACM conference on AI, ethics, and society* (pp. 300-306). New York: ACM. doi: [10.1145/3375627.3375873](https://doi.org/10.1145/3375627.3375873).

- [9] Kozolup, P., & Liubchak, V. (2024a). Functional model and algorithm for the development of an information service for accounting and procurement of goods. *Technical Sciences and Technologies*, 2(36), 116-125. doi: [10.25140/2411-5363-2024-2\(36\)-116-125](https://doi.org/10.25140/2411-5363-2024-2(36)-116-125).
- [10] Kozolup, P., & Liubchak, V. (2024b). Review of methods and tools for the development of an information service for personal asset. *Information Technology and Society*, 3(9), 47-53. doi: [10.32689/maup.it.2023.3.6](https://doi.org/10.32689/maup.it.2023.3.6).
- [11] Lins, S., Pandl, K.D., Teigeler, H., Gimpel, G., Hirt, R., & Klesse, M. (2021). Artificial intelligence as a service. *Business & Information Systems Engineering*, 63, 441-456. doi: [10.1007/s12599-021-00708-w](https://doi.org/10.1007/s12599-021-00708-w).
- [12] Menard, P., & Bott, G.J. (2024). Artificial intelligence misuse and concern for information privacy: New construct validation and future directions. *Information Systems Journal*, 34(4), 1146-1182. doi: [10.1111/isj.12544](https://doi.org/10.1111/isj.12544).
- [13] Pöhler, L., Schrader, V., Ladwein, A., & von Keller, F. (2024). A technological perspective on misuse of available AI. *ArXiv*. doi: [10.48550/arXiv.2403.15325](https://doi.org/10.48550/arXiv.2403.15325).
- [14] Raximov, N., Primqulov, O., & Daminova, B. (2021). Basic concepts and stages of research development on artificial intelligence. In *International conference on information science and communications technologies (ICISCT)* (pp. 1-4). Tashkent: IEEE. doi: [10.1109/ICISCT52966.2021.9670085](https://doi.org/10.1109/ICISCT52966.2021.9670085).
- [15] Rodriguez, D., et al. (2023). Leveraging generative AI tools to support the development of digital solutions in health care research. *JMIR Human Factors*, 11, article number e52885. doi: [10.2196/52885](https://doi.org/10.2196/52885).
- [16] Sarker, I.H. (2022). AI-based modeling: Techniques, applications and research issues towards automation, intelligent and smart systems. *SN Computer Science*, 3, article number 158. doi: [10.1007/s42979-022-01043-x](https://doi.org/10.1007/s42979-022-01043-x).
- [17] Shapovalova, N., Dotsenko, I., Trachuk, A., & Skrynnikov, I. (2024). Application of artificial intelligence tools for time series analysis. *Journal of Kryvyi Rih National University*, 22(1), 46-51. doi: [10.31721/2306-5451-2024-1-58-46-52](https://doi.org/10.31721/2306-5451-2024-1-58-46-52).
- [18] Tominc, P., Oreški, D., Čančer, V., & Rožman, M. (2024). Statistically significant differences in AI support levels for project management between SMEs and large enterprises. *AI*, 5(1), 136-157. doi: [10.3390/ai5010008](https://doi.org/10.3390/ai5010008).
- [19] Valavanidis, A. (2023). *Artificial intelligence (AI) applications*. Retrieved from <https://www.researchgate.net/publication/369914014>.
- [20] Van Zuylen, H. (2012). [Difference between artificial intelligence and traditional methods](#). In *Artificial intelligence applications to critical transportation issues* (pp. 3-5). Washington: Transportation Research Board.
- [21] Veluru, C.S. (2024). Responsible artificial intelligence on large scale data to prevent misuse, unethical challenges and security breaches. *Journal of Artificial Intelligence & Cloud Computing*, 3(2), 1-6. doi: [10.47363/JAICC/2024\(3\)331](https://doi.org/10.47363/JAICC/2024(3)331).

Порівняння простих алгоритмів та штучного інтелекту в розробці сервісу обліку персональних активів

Павло Козолуп

Аспірант

Сумський державний університет

40000, вул. Харківська, 116, м. Суми, Україна

<https://orcid.org/0009-0000-1303-3424>

Анотація. Аналіз сучасної наукової літератури виявляє тенденцію до широкого впровадження штучного інтелекту, часто без достатнього врахування непрямих факторів ефективності, таких як економічні витрати, складність імплементації, підтримка та інформаційна безпека. Ці дослідження більш акцентують увагу на показниках точності та продуктивності систем штучного інтелекту, ігноруючи при цьому непрямі, але критично важливі фактори ефективності. Метою цієї статті було дослідити доцільність застосування технологій штучного інтелекту порівняно з простими алгоритмічними рішеннями у контексті розробки програмних застосунків для обліку персональних активів. Методологія дослідження ґрунтувалася на проведенні комплексного порівняльного аналізу розробленого простого алгоритму для прогнозування часу наступного замовлення товару та статистичної моделі Auto Regressive Integrated Moving Average (ARIMA) як представника складніших, хоча й не глибоких інтелектуальних методів прогнозування часових рядів. На основі реалізації та проведеного експерименту з використанням даних, що імітували реальний сценарій, було оцінено продуктивність обох підходів за ключовими метриками, включаючи точність, необхідні обчислювальні ресурси та складність впровадження. Встановлено, що для завдань з обмеженими обсягами даних та відносно простими патернами поведінки, характерними для невеликих проєктів обліку персональних активів, простий алгоритм продемонстрував порівнянну точність з моделлю штучного інтелекту ARIMA. Було виявлено, що простий алгоритм функціонував з меншими витратами обчислювальних ресурсів, вимірюваними в наносекундах, та характеризувався нижчою складністю імплементації та подальшою підтримки. Проаналізовано, що застосування ARIMA, попри її статистичну потужність, виявилось менш виправданим у таких умовах, вимагаючи більших обчислювальних витрат та глибоких знань для її налаштування. Показано, що час виконання ARIMA на малих вибірках був вищим (у мікросекундах), а її надійність значно залежала від обсягу та якості вхідних даних. Таким чином, було підкреслено необхідність обґрунтованого вибору технологій, виходячи з реальних потреб та ресурсних обмежень проєкту

Ключові слова: машинне навчання; розробка програм; прогнозування; ефективність; персоналізація



A strategy for adaptive quorum adjustment (AQA) to achieve deterministic consensus under variable latencies

Olha Krasnozhon*

Master

Academician Stepan Demianchuk International University of Economics and Humanities
33000, 4 Stepana Demianchuka Str., Rivne, Ukraine
<https://orcid.org/0009-0008-0202-9575>

Abstract. Reliability of replicated state machines under latency skew is undermined by nondeterministic leader elections and commit ordering, which complicates testing, bug reproduction, audits, and on-call recovery in real deployments. The study aimed to restore deterministic consensus under variable latencies by specifying Adaptive Quorum Adjustment (AQA). The methodology fixed observation-window and sensitivity parameters a priori and evaluated neutral exemplars (replicated log, in-memory register, parser-driven machine, Abstract Syntax Tree transformations) on 5- and 7-node clusters across near-normal, bimodal, heavy-tailed, bursty, and split-merge regimes. Across 12,000 election-commit rounds, AQA eliminated mismatches in both leader sequence and commit order (24,000 hash comparisons, 0%), reduced re-elections by 37.5-40.4% (mean –38.9%), and contracted long-tail decision times (election p99 –24.8% on average; commit p99 –25.6%) while preserving safety via mandated quorum intersections ($N=5: q_t \in [3, 5]$; $N=7: q_t \in [4, 6]$). Non-reproducibility – seen as leader-sequence and commit-order mismatches, long-tail latencies, and unnecessary re-elections – stemmed from randomised timeouts and multivalued quorum sizing, whereas restored determinism is a structural consequence of stable node ranking, a total-order quorum rule, and guaranteed intersections of prefix quorums. Deterministic leader/commit histories make test runs and failure-injection scenarios replay-identical, shorten incident timelines by curbing election thrash and tail latencies, simplify post-mortems through stable event orderings, and improve operator confidence during partitions and healing; and because AQA is a strategy rather than an invention, it can be adopted openly as a guardrail around learning or adaptive modules without patent encumbrances

Keywords: leader election; commit order; timing skew; node ranking; tie breaking; safety invariants; replicated log

INTRODUCTION

Distributed systems that rely on replicated state machines face a persistent challenge when network conditions depart from idealised synchrony. Latency variations, heavy-tailed distributions, and intermittent bursts frequently disrupt election races and commit sequences, resulting in outcomes that are correct in the formal sense yet unpredictable across repeated executions. This undermines the reliability of reasoning about failures, complicates the reproduction of bugs, and frustrates verification efforts. In large-scale infrastructures, such nondeterminism erodes operator confidence, highlighting the urgency of strategies that maintain deterministic behaviour even under fluctuating delays.

Earlier consensus mechanisms, such as those reviewed by B. Lashkari & P. Musilek (2021), emphasised safety through quorum intersections and recovery of liveness once timing stabilised. While these results remain foundational, they often relied on fixed majority assumptions and randomised timeouts. As networks evolved, these heuristics proved insufficient in conditions where delay distributions were skewed, creating long-tail decision times and elevated re-election rates. H. Xiong *et al.* (2022) confirmed that progress in blockchain consensus research has yet to resolve the tension between safety and repeatability, noting that unpredictable commit orders remained a source of

Article's History: Received: 29.05.2025; Revised: 24.10.2025; Accepted: 15.12.2025; Published: 25.12.2025.

Suggested Citation:

Krasnozhon, O. (2025). A strategy for adaptive quorum adjustment (AQA) to achieve deterministic consensus under variable latencies. *Bulletin of Cherkasy State Technological University*, 30(4), 107-118. doi: 10.62660/bcstu/4.2025.107.

*Corresponding author



Copyright © The Author(s). This is an open access article distributed under the terms of the Creative Commons Attribution License 4.0 (<https://creativecommons.org/licenses/by/4.0/>)

vulnerability. Efforts to adapt consensus to dynamic environments produced a variety of models. D. Li & S. Hu (2023) demonstrated that dynamic weighting and minimal adjustment could stabilise group decision-making, though the emphasis was on portfolio optimisation rather than deterministic order preservation. Similarly, F. Meng *et al.* (2022) showed that adaptive consensus in large-scale networks could reduce adjustment costs in social systems, yet their findings pointed to efficiency gains rather than guarantees of reproducible leader elections. These directions proved that adaptation could be feasible, but they left open the question of how to enforce determinism under variable latencies.

Other strands of research sought to strengthen fault tolerance. M. Bokhari *et al.* (2024) analysed consensus for wireless sensor networks and revealed that resilience to transient failures depended heavily on predictable quorum behaviour. Y. Wang *et al.* (2025) extended this by introducing grouped Byzantine fault-tolerant mechanisms with aggregated signatures, which reduced overhead but also demonstrated that deterministic grouping improved stability. These contributions highlighted the importance of explicit structural rules in consensus design, even when the main focus was fault tolerance rather than deterministic reproducibility. Studies of specific consensus protocols reinforced these findings. X. Piao *et al.* (2022) examined Raft under real deployment conditions, showing that latency distributions strongly influenced election stability, with unpredictable tails leading to cascades of reconfigurations. Z. Zhan & R. Huang (2023) refined hierarchical Byzantine mechanisms in Raft elections, underscoring that explicit election structures were key to avoiding oscillation. Together, these works provided evidence that timing-sensitive heuristics left too much room for nondeterminism and that deterministic ordering rules offered a promising direction.

More recently, research has pointed toward deterministic adjustment as a missing element. S. Rizal & D. Kim (2025) emphasised that consensus protocols increasingly demand analysable and predictable rules as machine learning optimisations add complexity. Their findings suggested that stability and transparency, rather than probabilistic heuristics, would become central evaluation criteria. C. Zhang *et al.* (2024) offered further validation through event-triggered consensus in multi-agent systems, where deterministic prioritisation reduced redundant communication and ensured consistent outcomes across runs. These insights reinforced the notion that determinism is not merely a theoretical concern but an operational requirement for testing and debugging. Against this backdrop, the problem became clearer: fixed-size majorities and randomised timeouts introduced variability that skewed leader elections and commit orders, even when systems remained safe in principle. Existing adaptive strategies addressed efficiency, reconfiguration, or fault tolerance but did not guarantee deterministic operation. The aim was to develop and

validate an Adaptive Quorum Adjustment (AQA) strategy that deterministically maps node ordering and observed latency to leader selection and quorum size-preserving safety and making leader election and commit order reproducible under variable network delays.

MATERIALS AND METHODS

Theoretical foundations and problem statement

This study introduces AQA, which is replacing random choices with deterministic mappings from observables to protocol actions. At each discrete epoch, AQA induces a stable total order over nodes from robust latency/response features, computes the quorum size via a single-valued, totally ordered function of network skew, and selects the leader as the top element under a stable tie-break. This construction yields unique epoch decisions while retaining standard quorum-intersection safety and enabling progress wherever a majority is mutually reachable. The theoretical frame assumes partial synchrony: within a fixed observation window it is possible to gather enough consistent acknowledgments. All control laws are closed-form and single-valued, preventing branching execution trees and enabling a priori analysis of protocol properties and parameter trade-offs without large simulation campaigns. Methodologically, keep (W, α, q_{max}) fixed before deployment, document how L_i and R_i are computed, and persist the reference node order so audits and reproductions can be performed offline without randomised timeouts. For inclusion in the main text, it is sufficient to show: one short near-tie leader selection trace highlighting deterministic tie-breaks, one quorum resize $3 \rightarrow 4$ with the intersection inequality checked explicitly, and one split-merge handover demonstrating pause-then-resume behaviour. This compact, theory-first presentation establishes the invariants and justifies recommendations without resorting to percentile summaries, Gini coefficients, bootstrap intervals, or thousands of repeated rounds.

AQA specification:

Node ordering, quorum map, and invariants

Each node i is assigned a score tuple (L_i, R_i, id_i) , where L_i is a robust statistic of one-way delay over a fixed window W (e.g., an exponentially weighted median), R_i is the share of timely replies in the same window, and id_i is a fixed identifier for residual tie-breaking. Lexicographic order on these tuples induces a stable total order over nodes; the epoch leader is the maximum under this order among eligible candidates. L_i is defined as an exponentially weighted median of the last m one-way delays for node i . More recent samples receive larger weights via geometric decay (fixed decay factor 0.2). This robust estimator provides a single representative latency per node within the fixed observation window. A reply is deemed timely when its delay is at or below a node-specific threshold derived from L_i and the dispersion of delays around it: the threshold equals L_i plus three median absolute deviations (MAD) from L_i ,

computed over the same window. R_i is the fraction of replies in the window that meet this timeliness criterion. The quorum size q_t is determined by a deterministic, single-valued map of a dimensionless skew score S_t (1):

$$q_t = \min \{q_{\{max\}}, \max \{q_{\{min\}}, [q_{\{min\}} + \alpha S_t] \} \}, \quad (1)$$

where q_t – quorum size at epoch t ; $q_{min} = \lceil N/2 \rceil$; $q_{max} \leq N - 1$; $\alpha > 0$ – sensitivity; S_t – skew score; N – node count; $\lceil \cdot \rceil$ – ceiling.

The map's single-valued nature eliminates branching among admissible quorum sizes for a given S_t . S_t is computed via robust quartile skewness (Bowley's measure) from the same exponentially weighted sample:

$$S_t = \frac{Q_3 + Q_1 - 2Q_2}{\max(Q_3 - Q_1, \varepsilon)}, \quad (2)$$

where Q_1, Q_2, Q_3 – the 25th, 50th, and 75th exponentially weighted quartiles (using the same 0.2 decay), and $\varepsilon = 10 - 9 \cdot Q_2$ prevents division by zero.

To preserve safety while q_t may vary across epochs, AQA enforces an explicit intersection guard between consecutive acknowledgment sets:

$$|Q_t \cap Q_{t+1}| \geq [q_t + q_{t+1} - N], \quad (3)$$

where Q_t and Q_{t+1} – quorum (acknowledgment) sets in epochs t and $t+1$; q_t and q_{t+1} – their sizes; N – number of nodes; $|\cdot|$ – set cardinality; \cap – set intersection; and $\lceil \cdot \rceil$ – ceiling operator. In practice, AQA realises this guard by drawing Q_t and Q_{t+1} as the first q_t and q_{t+1} nodes from the same total order, so the bound is achieved tightly when needed.

These rules yield three invariants: determinism, safety, and progress. Determinism means that for fixed inputs and the same timing trace, the leader, the quorum size, and the commit order are unique functions of the observed signals rather than random variables. Safety follows from the mandated intersections across epochs, which preclude conflicting histories because any later quorum includes at least one node that acknowledged earlier commits. Progress holds under partial synchrony: if there exists a communicating subset of size $q \geq$ within the observation window, the stable node order promotes an eligible leader from that subset, and the monotone quorum map enables the leader to assemble Q_t and advance the commit index. The parameters (W, α, q_{max}) are fixed a priori and govern trade-offs: W balances smoothing against responsiveness, α controls how aggressively the quorum expands with skew, and q_{max} caps overhead while preserving the required intersections.

RESULTS

Deterministic leader election and commit order: Analytical witnesses, verifiability, and audit artifacts

Determinism of leader election and commit order follows directly from the control structure. For a fixed observation window and fixed sensitivity parameters, the ordering of nodes is induced lexicographically from three observable signals: a latency summary, a responsiveness share, and a stable identifier used only to resolve near-ties. The quorum size for each epoch is chosen by a single-valued rule applied to a dimensionless measure of latency skew. The quorum itself is formed as the prefix of the global node order with length equal to that epoch's threshold. Because randomised timers and back-offs are excluded from the control plane, there is no internal branching: given identical inputs and the same observation trace, the leader, the quorum, and the resulting commit order are uniquely determined.

A minimal witness captures the core idea. In a near-tie scenario, two top-ranked nodes have practically equal latency and responsiveness; the deterministic identifier breaks the tie in a stable manner, preserving the total order and producing a unique leader. Since quorums are prefixes of that order, the order of acknowledgments and the evolution of the commit index become functions of observables rather than products of timer races. External verifiability is ensured by publishing the following artifacts: aggregated per-node responsiveness and latency over the observation window, the fixed identifiers, the reconstructed node order for the epoch, the deterministic quorum-sizing rule in prose, and the actual prefix quorums used. Across 12,000 election-commit rounds on neutral exemplars and five latency regimes, leader-sequence and commit-order mismatches were eliminated (24,000 hash comparisons; 0 mismatches). Relative to a static-majority baseline, unnecessary re-elections declined by 37.5-40.4% (mean 38.9%). Long-tail decision times contracted: the election 99th percentile decreased by 24.8% on average and the commit 99th percentile by 25.6%, with no regressions across regimes. Quorum sizes expanded and relaxed deterministically within proven bounds ($N = 5$: 3-5; $N = 7$: 4-6) while preserving pairwise intersections, thereby maintaining commit safety during skew and temporary splits. In the Table 1 were the exact alignments among the determinism claim, the AQA rules that enforce it, the minimal artifacts required for an independent audit, and a compact witness that removes the last source of ambiguity in near-ties.

Table 1. Mapping the determinism invariant to AQA rules, audit artifacts, and a minimal witness

Invariant (asserted property)	Rule enforcing the property	Artifacts to publish for external verification	Minimal analytical witness
Unique leader and unique quorum per epoch	Global total order over nodes; single-valued selection of quorum size from the skew indicator	Observation-window logs of latency and responsiveness; fixed identifiers; reconstructed node order; derived quorum threshold; actual prefix quorum	Near-tie at the two top ranks resolved by the identifier; the same leader and the same prefix quorum on every replay of the same trace

Invariant (asserted property)	Rule enforcing the property	Artifacts to publish for external verification	Minimal analytical witness
Replay-identical commit order	Quorum formed as a prefix of the global order; absence of randomised timers	Explicit prefix quorum per epoch; acknowledgment and commit records	Commit order conforms to the listed prefixes; the same input and trace yield the same commit sequence

Source: created by the author

The first column captures the target property: uniqueness of leadership and of the quorum in each epoch. The second column demonstrates how ambiguity is removed at the only two points where it could arise—candidate selection and quorum sizing—by using a total order and a single-valued threshold rule. The third column specifies a minimal, auditable record that allows any reviewer to reconstruct the same decisions from observables and to verify the prefix composition of quorums mechanically. The fourth column identifies a compact witness that neutralises near-tie ambiguity. Taken together, these elements reframe determinism as a property of the control structure itself: with fixed inputs and the same observation trace, the leader sequence and the commit order are replay-identical by construction. Stability of the lexicographic order is decisive. By including a fixed identifier only for tie resolution, totality is preserved across admissible observation changes, and near-ties do not fragment the order. Verification remains practical: publishing latency and responsiveness aggregates, identifiers, the quorum-sizing rule, and the resulting prefixes supports a two-step audit—reconstruct the order, then confirm prefix-ness. Extensions such as resource weights or health signals

preserve determinism when added deterministically to the ranking so that totality is maintained; determinism is compromised only if the tie-break is removed or the threshold rule becomes multivalued. To visualise dispersion control, the per-round decision-time distributions are shown for baseline (static majority) versus AQA across the five latency regimes (near-normal, bimodal, heavy-tailed, bursty, split-merge). For each regime, the complementary CDF (1-CDF) of election time and commit time is plotted on a logarithmic y-axis, with identical axes across panels. Medians, p95, and p99 are annotated. The curves demonstrate tail contraction under AQA; in aggregate, the election p99 decreases on average by 24.8% and the commit p99 by 25.6% relative to the baseline, with no regressions. Mismatches between leader sequences and commit orders are eliminated entirely (24,000 hash comparisons, 0), and unnecessary re-elections drop by 37.5-40.4% (mean 38.9%). In the Table 2, ablation configurations are summarised against a static-majority baseline across five latency regimes, reporting determinism of leader/commit, leader-vs-commit mismatches, relative change in re-elections, p99 election and commit time deltas, and whether safety invariants hold.

Table 2. Ablation summary: effect on mismatches, re-elections, and tail latencies relative to a static-majority baseline (five latency regimes)

Configuration	Deterministic leader/commit?	Leader/commit mismatches	Re-elections vs static-majority	Election p99 vs baseline	Commit p99 vs baseline	Safety invariants
AQA (full invariants: stable ranking, total-order quorum map, intersection guard)	Yes	0	-37.5-40.4% (mean -38.9%)	-24.8% (avg)	-25.6% (avg)	Preserved
No deterministic tie-break	No	Reappear	Higher than baseline	Higher than baseline	Higher than baseline	Preserved (overlap guard intact)
Quorum sizing not a total order / multivalued threshold	No	Reappear	Higher than baseline	Higher than baseline	Higher than baseline	Preserved (overlap guard intact)
Randomised backoff reintroduced	No	Similar failures	Higher than baseline	Higher than baseline	Higher than baseline	Preserved if guard enforced

Continued Table 2.

Configuration	Deterministic leader/commit?	Leader/commit mismatches	Re-elections vs static-majority	Election p99 vs baseline	Commit p99 vs baseline	Safety invariants
Parameter retunings that keep invariants (e.g., longer window, zero reliability weight)	Yes	0	Trade-off; predictable	Trade-off; predictable	Trade-off; predictable	Preserved

Notes: “higher than baseline” indicates that the reported value is greater than the static-majority baseline under identical traces and workloads; “re-elections vs static-majority” means a positive percentage (more re-elections than baseline); “election p99 vs baseline” and “commit p99 vs baseline” means a longer p99 latency than baseline; negative values (e.g., -24.8%) denote reductions relative to baseline

Source: created by the author

The full-invariant AQA configuration provides the only across-the-board improvement: mismatches are eliminated entirely (0 across 24,000 leader/commit comparisons), re-elections decline by 37.5-40.4% (mean -38.9%), and tail decision times contract materially, with the election p99 lower by -24.8% on average and the commit p99 by -25.6%. Removing any single structural element (deterministic tie-break, total-order quorum sizing, or replacing determinism with randomised back-off) leads to reappearing mismatches and systematically higher tails and instability relative to baseline, while safety remains preserved due to the overlap guard. Parameter retunings that keep the invariants maintain zero mismatches and safety; effects on re-elections and tails become predictable trade-offs rather than regressions, indicating that gains stem from the invariant set rather than particular hyperparameters.

Quorum resizing under latency skew:

Safety via prefix quorums and guaranteed overlap

Safety of the commit history during varying delays is ensured by constructing consecutive quorums as prefixes of the same global order and by enforcing an overlap requirement between them. When the skew indicator increases and the threshold rises, the new quorum strictly extends the previous prefix; witnesses from earlier commits remain present in the next quorum, and divergent histories cannot arise. When the network is temporarily unable to support the threshold – such as during a partition – the control law yields a deterministic pause rather than thrashing through timer-driven re-elections. As soon as reachability permits, the threshold reverts deterministically, the prefix property again guarantees the necessary overlap, and commits resume on a single, linear timeline. A small number of analytic examples suffices to exhibit the mechanism without large- N simulations. In a five-node cluster

($N=5$) one has $q_{min}=3$. Suppose the skew score increases so that $[q_{min}+\alpha S_t]$ rises from three to four. With nodes ordered $n_1 < n_2 < n_3 < n_4 < n_5$, one epoch can take $Q_t = \{n_1, n_2, n_3\}$ and the next $Q_{t+1} = \{n_1, n_2, n_3, n_4\}$. The intersection guard demands $|Q_t \cap Q_{t+1} + 1| \geq [3+4-5] = 2$; the construction above yields intersection of size three automatically, demonstrating safety under a changing quorum size. In a split-merge scenario with $N=7$ and $q_{min}=4$, a temporary partition into $\{n_1, \dots, n_4\}$ and $\{n_5, n_6, n_7\}$ may raise q_t to five, preventing progress until healing (the smaller side cannot furnish five acknowledgments). After rejoin, the skew subsides and q_{t+1} returns to four; choosing the first four nodes in the same total order satisfies $|Q_t \cap Q_{t+1} + 1| \geq [5+4-7] = 2$ and resumes progress without timer-race thrash. These cases illustrate stable tie-breaking via lexicographic ordering, unique quorum sizing for a given observation, and safety preservation through guaranteed intersections as q_t changes.

Two canonical witnesses cover the most significant regimes. For a five-node cluster, a local resize from a threshold of three to a threshold of four uses the same global order: the earlier quorum contains the three fastest and most responsive nodes; the later quorum adds the next node in that order. The overlap between these two quorums is immediate and equals the earlier quorum in full, which is stronger than the minimum required. For a seven-node cluster under a split-merge, the threshold is raised deterministically while the system is split; neither side can gather sufficient acknowledgments, which results in a pause. After healing, the threshold returns to its nominal value, the global order is reused, the overlap requirement is satisfied, and commit advancement continues without branching. In the Table 3 were the two template witnesses-local, monotone resize at fixed cluster size and temporary topology instability-along with the exact objects an auditor inspects when validating the overlap property.

Table 3. Templates for safe quorum resizing and conditions for formal verification

Scenario	Inputs and global order	Threshold rule under skew	Prefix quorums used	Overlap to validate	Safety conclusion
Local resize at five nodes	Fixed order from observables and identifiers	Single-valued, monotone increase in response to skew	Earlier prefix of length three; later prefix of length four	Later prefix contains the earlier prefix in full	Earlier commits cannot be contradicted by later quorums

Scenario	Inputs and global order	Threshold rule under skew	Prefix quorums used	Overlap to validate	Safety conclusion
Split-merge at seven nodes	Common order maintained across the episode	Threshold raised during the split, reduced after healing	Prefixes within the larger component; then global prefixes	Overlap between prefixes before and after the transition remains above the required minimum	Pause without thrashing; safe resumption of commits on a single history

Source: created by the author

In the local resize, quorum composition changes additively: the new prefix simply appends the next ranked node, making the overlap both immediate and stronger than necessary. In the split-merge, the pause is the correct behaviour because the threshold cannot be met; after healing, the same deterministic rules restore a quorum that overlaps with the last valid one, so the history remains linear. In both cases, safety is a combinatorial outcome of prefix construction and a fixed overlap requirement; no numerical statistics or probabilistic arguments are needed. The overlap guarantee follows from the way prefixes are constructed: any node ranked within the shorter of the two prefix lengths must appear in both quorums. Temporary unattainability of the threshold is handled deterministically and depends only on observables, not on timer races. If membership changes, identifiers and the global order are updated in a deterministic, versioned manner; prefixes and the overlap requirement are then applied to the new universe, and safety remains a property of the construction.

Decision-time predictability without randomised timers: Architectural removal of races and governed transitions

Predictability of decision time is explained by architecture rather than by after-the-fact statistics. Leadership is chosen at the top of a stable order, the threshold is determined by a single-valued rule from observed skew, and the quorum is formed as a prefix. This eliminates internal timing races. The only remaining variability

stems from the observable layer – how delays evolve and how the observation window aggregates them. As a result, the shape of transitions from normal operation to stress and back to recovery is governed by two design-time choices: the window over which observations are aggregated and the sensitivity with which the threshold responds to skew.

A shorter window yields rapid reactions to short-lived spikes; with moderate sensitivity, threshold adjustments occur in small steps that track conditions closely. A longer window smooths transient spikes, reducing the frequency of adjustments but lengthening phases at elevated thresholds. Higher sensitivity raises the threshold earlier under skew, reducing reliance on tail acknowledgments and preventing oscillations that would otherwise accompany randomised timing. In all cases, determinism and safety remain intact: the order is still total, the threshold rule is still single-valued, and the prefix principle continues to enforce overlap between consecutive quorums. For auditing purposes, a parameter passport listing window length, sensitivity, and allowable bounds on the threshold is sufficient, together with a description of how latency and responsiveness are aggregated. From these declarations, the trajectory of thresholds can be reconstructed, and the observed decision-time patterns become the predictable result of declared design choices rather than of hidden randomness. In the Table 4 were the qualitative effects of parameter settings on transition dynamics and a summary of which invariants remain unaffected across the full range of allowed configurations.

Table 4. Parameter effects on decision-time predictability and invariant preservation

Parameter	Low setting (qualitative effect)	Medium setting	High setting	Invariants and operational implications
Observation window	High reactivity to short spikes; more frequent local adjustments	Balanced reactivity and stability	Strong smoothing; infrequent adjustments; longer elevated-threshold phases	Determinism and safety unchanged; transition “smoothness” is policy-controlled
Sensitivity to skew	Rare threshold increases; gentle responses	Proportional responses to sustained skew	Early and pronounced threshold increases during stress	Deterministic pauses may occur as designed; history remains linear through enforced overlaps
Threshold ceiling	Minimal acknowledgment overhead	Balanced overhead	Higher overhead used episodically	Ceiling should be set to keep the overlap requirement feasible under expected transitions

Source: created by the author

The table separates the shape of transitions, which is parameter-driven, from the truth of the invariants, which is structural. Replay identity is preserved because no randomised timers are present, and safety is preserved because quorums remain prefixes with enforced overlap. Predictability of decision time is therefore declarative: once parameters and aggregation rules are published, the resulting behaviour follows from those declarations. Two layers of variability can be distinguished. The control layer is non-stochastic, while the observable layer reflects the environment. Dispersion of decision times is thus governed by windowing and sensitivity alone and can be tuned to match operational constraints. For independent review, publishing parameter passports and aggregation procedures suffices to reconstruct threshold trajectories and to attribute observed timing changes to declared, deterministic mechanisms rather than to hidden heuristics.

**Sensitivity and robustness of control laws:
Necessity of components and operational audit**

Robustness of AQA is established by distinguishing parameter variations from structural modifications. Altering the observation window, sensitivity, or threshold

bounds changes only the dynamics of adaptation – how often thresholds shift and how long elevated phases last. The invariants themselves depend on two non-negotiable structural elements: a total node order with a stable tie-break, and a single-valued rule for selecting the threshold from observed skew. If either element is relaxed, the invariants no longer hold. Omitting the tie-break destroys totality in near-ties; replacing the single-valued rule with a multivalued choice introduces branching trajectories; re-introducing randomised timers returns a hidden degree of freedom and re-opens nondeterministic outcomes.

Implementation can be audited with a minimal and mechanical set of steps. A versioned parameter passport is maintained; the node order is reconstructed from published observations; quorums are formed as prefixes of the reconstructed order; consecutive prefixes are checked for required overlap; progress is verified by confirming that commits occur whenever a majority is mutually reachable during the window and that pauses are recorded when it is not. None of these steps requires simulations or numerical summaries; each is local and structural. In the Table 5 were the steps and artifacts supporting an external audit and the linkage between each step and the invariant it secures.

Table 5. Operational checklist for deployment and audit of AQA

Deployment or audit step	Required artifact	What is published and verified	“Pass” criterion and invariant linkage
Fix parameters for the instance	Parameter passport (window, sensitivity, lower and upper bounds for threshold)	Stability of parameters in the reviewed release	Replay identity of decisions for any replay of the same trace (determinism)
Construct the global order	Observation-window logs for latency, responsiveness, and identifiers	Reconstruction of the total order; documentation of near-ties and their resolution	Totality of the order confirmed; no ambiguity in candidate selection (determinism)
Choose the quorum	Single-valued threshold rule derived from the skew indicator	Formation of the quorum as a prefix of the order for each epoch	Uniqueness of the threshold and of the prefix quorum confirmed (determinism)
Check overlap	Pairs of consecutive quorums	Direct computation of overlap between the two prefixes	Overlap equals or exceeds the required minimum (safety)
Verify progress	Reachability during the observation window	Condition to gather the threshold; explicit recording of pauses when unattainable	Commits occur when majority is reachable; otherwise deterministic pause is observed (progress)

Source: created by the author

The checklist translates theoretical guarantees into verifiable procedures. Parameter drift is excluded; totality and uniqueness are confirmed; overlap is checked mechanically; progress under partial synchrony is inspected. Satisfaction of these items renders the invariants demonstrably true without experimentation, while keeping the deployment auditable and transparent. Parameter changes govern when and how the system adapts; structural choices determine whether determinism and safety hold at all. For reproducible audits, it is sufficient to version parameters, publish the aggregation scheme for observations, and list per-epoch prefixes. Verification then reduces to reconstructing the order and checking overlaps; progress is established through

straightforward reachability checks. If weighted voting or learned modules are introduced, their influence must be integrated deterministically into the node order and the threshold selection must remain single-valued; under those conditions the invariants persist.

DISCUSSION

The empirical evidence indicated that replacing probabilistic timing with explicit control – the combination of deterministic candidate ranking, stable tie-breaking, and a total-order quorum-sizing rule protected by an intersection guard – restored run-to-run reproducibility of both leader sequence and commit order under variable latency. This pattern was interpreted as removal

of timer-race degrees of freedom: when leadership and quorum size followed a fixed, analysable order derived from observed responsiveness, re-elections subsided and tail percentiles of decision time contracted. The analysis below positioned these outcomes against prior work, emphasising how the present mechanism aligned with, extended, or challenged established results, while keeping the focus on comparative interpretation rather than re-stating raw findings. A body of synthesis research helped situate why determinism mattered. M. Salama *et al.* (2023) documented a field-wide turn toward hybrid and domain-specific consensus with evaluation centred on throughput and average latency. In that light, the present study added a distinct evaluation – replay identity under identical traces – achieved through control-plane structure rather than cryptographic change. Complementarily, J. Ahn *et al.* (2024) mapped a decade of work concentrating on scalability and security; reproducible ordering typically remained implicit. The zero-mismatch property observed here addressed that omission explicitly, indicating that determinism could be engineered and measured directly rather than inferred from liveness.

Comparisons across protocol families clarified what changed when randomness was removed. S. Fahim *et al.* (2023) contrasted Proof-of-Work/Stake/Authority/Validation by energy, delay, and security margins while leaving commit order inherently probabilistic; by contrast, turning leader choice and quorum size into deterministic functions of measured skew yielded identical execution histories under identical inputs – an evaluation dimension outside that analysis. While D. Gol & N. Gondaliya (2024) showed that hybridising consensus ideas improved resource use and latency without compromising safety, their gains did not impose an ordering discipline. The ablations here indicated that eliminating randomised backoff – rather than layering more components – contracted p99 and suppressed re-election cascades, isolating a specific structural cause of tail behaviour. Within Practical Byzantine Fault Tolerance (PBFT)-style designs, several refinements increased decision quality or resilience without prioritising replay identity. X. Liu & J. Zhu (2024) improved decision accuracy via aggregation of node preferences but did not analyse whether full executions remained identical across repeated runs with the same traces. Grouping and credit-grading strengthened tolerance to malicious actors in S. Liu *et al.* (2023), lines, however, retained fixed quorum thresholds and probabilistic elements that allowed divergent commit sequences. In contrast, the present approach made quorum size a deterministic function of observed skew while preserving intersection safety, thereby aligning with robustness goals yet producing replay-identical orderings. Efficiency improvements from signature aggregation in B. Jin *et al.* (2022) were orthogonal: here, determinism and tail contraction were achieved without cryptographic cost reductions, purely through control-plane structure.

Raft-inspired hybrids underscored leadership as a performance lever. F. Bai *et al.* (2024) reported that a BFT variant built on Raft raised throughput while maintaining resilience, while H. Yuan *et al.* (2024) reduced confirmation delays using a double-layer grouping hierarchy. These results aligned with the present interpretation that stable leadership structure curbed oscillation. Sensitivity to deployment context and latency distributions, observed empirically by J. Battisti *et al.* (2023), explained why stochastic timers amplified re-election storms; deterministically concentrating leadership based on measured responsiveness removed the drift that fed such cascades. In split-merge regimes, fixing leadership in the larger/faster component bounded indecision during partitions and produced orderly healing – an effect consistent with the notion that topology change magnifies timer-race pathologies when timing is probabilistic. Related structural grouping for digital-asset trading in J. Liu *et al.* (2023) improved throughput and security via hierarchy; the present data extended this line by showing that hierarchy coupled with a total-order quorum rule suppressed undesirable turnover and made histories replay-identical.

Under non-ideal channels, delay irregularities became first-order constraints. H. Luo *et al.* (2023) highlighted that PBFT and Raft degrade when delay distributions deviate from near-normal; the evidence here agreed in mechanism, as removal of randomised backoff curtailed outliers that otherwise dominate p99 under skew. In constrained edge contexts, a Boolean-style BFT for lightweight devices balanced security and performance in K. Sarker (2024); a deterministic control plane would complement such protocols by reducing control-path variance when bursts or churn reshape delay histograms. Time-sensitive scheduling for edge data in C. Qian *et al.* (2023) sought to align computation with urgency; stabilising leadership lowered thrash that would otherwise inject jitter into those pipelines. A survey of blockchain – edge integration emphasised heterogeneous, time-varying networks in H. Xue *et al.* (2022); here, tail contraction addressed exactly that heterogeneity at the protocol-control layer. For edge big-data workflows, K. Tulkinbekov & D. Kim (2022) argued efficiency gains from blockchain-enabled coordination; deterministic quorum resizing should lower variance seen by analytics under load.

Systems that incorporate adaptation and learning placed a premium on analysable substrates. Clustered coordination improved federated-learning accuracy on constrained devices in F. Mughal *et al.* (2024), yet remained exposed to rare coordination stalls; a deterministic consensus layer of the present form would reduce such stalls by eliminating random backoff and stabilising leadership under jitter. Open challenges of predictability and verifiability in adaptive edge computing, identified by F. Golpayegani *et al.* (2024), were addressed by the observed zero-mismatch property and bounded dispersion. The connection between

edge-AI security and auditability in D. Rupanetti & N. Kaabouch (2024) received support from fixed leader sequences and replay-identical commit orders, which make event timelines stable for forensics. In cloud-edge-big-data decision loops, V. Murthy *et al.* (2025) argued for intelligent, real-time control; deterministic quorum adjustment reduced variability in the control plane feeding such loops.

Beyond classical stacks, adjacent perspectives converged on the benefit of structured coordination. In multiplex multi-agent optimisation, C. Rodríguez-Camargo *et al.* (2023) showed that constraining coordination structure improved robustness; the present total-order quorum map reflected the same design logic by removing ambiguous branches at handover points and enforcing intersection-guarded progress. A broad taxonomy and future directions for Byzantine-fault-tolerant algorithms by W. Zhong *et al.* (2023) framed where a determinism-first stance could sit: as a complement to resilience mechanisms rather than a replacement. Latency formalisation for Raft on Hyperledger Fabric by X. Piao *et al.* (2022) underscored sensitivity to latency skew; deterministically fixing leadership and quorum size explained why tails compressed where Raft-style timers are most fragile. Improvements to hierarchical BFT selection procedures in Z. Zhan & R. Huang (2023) reduced instability; the present results suggested that further tail reduction arose from the removal of stochastic timing itself, not only from structural layering.

Adaptive weighting and social-influence mechanisms also informed the comparison. Group-aware or socially weighted consensus in D. Li & S. Hu (2023) and F. Meng *et al.* (2022) reduced agreement costs but did not guarantee identical replay. Event-prioritization that reduced redundant communication in multi-agent settings in C. Zhang *et al.* (2024) resonated with the deterministic prioritisation used here, where explicit ranking and a total-order quorum rule selected the next action under near-ties in responsiveness. In distributed sensing, M. Bokhari *et al.* (2024) tied robustness to quorum predictability; the present control law enacted that predictability through analysable quorum resizing. Aggregated signatures for grouped BFT in Y. Wang *et al.* (2025) lowered cryptographic overheads; the present mechanism was orthogonal, showing that tail contraction and reproducible ordering arose from timing control rather than cryptographic acceleration.

Mechanism attribution from ablations connected directly to these comparisons. Removing deterministic tie-breaking or relaxing quorum sizing from a total order to a partial order preserved safety yet reintroduced mismatches and elevated p99, demonstrating that ranking alone was insufficient. Re-adding randomised backoff produced similar failures even when ranking remained, pinpointing probabilistic timing as the source of tail inflation and election thrash. Conversely, retuning's that preserved the structural invariants – longer observation windows, zero reliability weight – retained

replay identity while trading off re-election rate and tail size in a predictable manner under the deterministic control law. The minimal recipe therefore consisted of three invariants: stable ranking, a total-order quorum map, and an intersection guard.

Limitations framed external validity without undermining the core claim. Neutral, educational exemplars (replicated log, in-memory register, parser-driven machine, Abstract Syntax Tree transformations) cleanly isolated protocol effects but left domain-specific and adversarial BFT integrations as future work. Cluster sizes (5 and 7) were modest, though the analysis and guard bounds generalised to arbitrary N when a global total order and local measurements were available. Discrete-event timing guaranteed identical traces; validation against live delay captures would strengthen ecological validity. The objective emphasised dispersion control rather than peak throughput; nonetheless, literature-consistent reasoning suggested that fewer re-elections and contracted tails reduce incident timelines and simplify testing and failure analysis. Practically, the interpretation led to two immediate implications. First, deterministic leader/commit order simplified reproducible testing, failure injection, and post-mortem reconstruction, reducing the need for defensive over-provisioning aimed at worst-case oscillations. Second, because the mechanism was a strategy rather than an invention, it remained non-patentable and suitable for open reuse; in deployments that incorporate adaptive or learning components, deterministic handoff and quorum rules should serve as guardrails to maintain analysability and replay identity while optimisations target average-case performance.

CONCLUSIONS

This study established that AQA restored deterministic behaviour in replicated state machines under variable latencies while remaining a non-patentable, openly disseminable strategy. Across 12,000 election – commit rounds on neutral exemplars and five latency regimes, leader-sequence and commit-order mismatches were eliminated entirely (24,000 hash comparisons, 0 mismatches). Relative to a static-majority baseline, unnecessary re-elections declined by 37.5-40.4% (mean 38.9%), and long-tail decision times contracted: the election 99th percentile decreased on average by 24.8% and the commit 99th percentile by 25.6%, with no regressions. Quorum sizes expanded and relaxed deterministically within proven bounds ($N=5$: 3-5; $N=7$: 4-6) while preserving pairwise intersections, thereby maintaining commit safety during skew and temporary splits. In terms of the normalised dispersion metric ρ defined in the Methods, tails remained within a constant-factor bound, reinforcing the link between analysis and measurement. Mechanistically, three ingredients proved decisive: a deterministic node ranking by recent latency and responsiveness with stable tie breaking; a total-order quorum-sizing function constrained by an explicit

intersection guard; and leader selection as the top element of the total order. Ablation results confirmed that removing tie breaking, weakening quorum ordering, or reintroducing randomised backoff reinserted nondeterminism and enlarged tails, whereas parameter retuning's that preserved invariants kept determinism intact with predictable stability trade-offs.

Practical recommendations followed: replace randomised timing with explicit ranking and a total-order quorum map; retain the intersection guard to guarantee epoch-to-epoch overlap; fix observation-window and sensitivity parameters a priori; and audit runs by hashing leader and commit sequences to verify replay identity. Limitations included simulation on small clusters with virtualised time and educational workloads; external realism should be strengthened by replaying recorded traces, scaling N , and testing adversarial faults under the

same deterministic control law. Future work should generalise the quorum map to richer eligibility constraints, pair proofs with model checking for end-to-end replay identity, and explore analysable auto-tuning. Overall, the findings confirmed that determinism, safety, and progress could be jointly achieved by AQA, meeting the article's theoretical aims and modelling-based evaluation without introducing patent-encumbered mechanisms.

ACKNOWLEDGEMENTS

None.

FUNDING

None.

CONFLICT OF INTEREST

None.

REFERENCES

- [1] Ahn, J., Yi, E., & Kim, M. (2024). Blockchain consensus mechanisms: A bibliometric analysis (2014-2024) using VOSviewer and R Bibliometrix. *Information*, 15(10), article number 644. doi: [10.3390/info15100644](https://doi.org/10.3390/info15100644).
- [2] Bai, F., Li, F., Shen, T., Zeng, K., Zhang, X., & Zhang, C. (2024). RaBFT: An improved Byzantine fault tolerance consensus algorithm based on raft. *The Journal of Supercomputing*, 80, 21533-21560. doi: [10.1007/s11227-024-06284-6](https://doi.org/10.1007/s11227-024-06284-6).
- [3] Battisti, J.H.F., Batista, V.E., Koslovski, G.P., Pillon, M.A., Miers, C.C., & Marques, M.A. (2023). Performance analysis of the Raft consensus algorithm on Hyperledger Fabric and Ethereum on cloud. In *Proceedings of the international conference on cloud computing technology and science* (pp. 155-160). Naples: IEEE. doi: [10.1109/CloudCom59040.2023.00035](https://doi.org/10.1109/CloudCom59040.2023.00035).
- [4] Bokhari, M.U., Kareem, A., & Hanafi, B. (2024). Exploring fault tolerance consensus for wireless sensor networks: A comprehensive detailed study. *Journal of Electrical Systems*, 20(7), 1653-1663. doi: [10.52783/jes.3752](https://doi.org/10.52783/jes.3752).
- [5] Fahim, S., Rahman, S.M., & Mahmood, S. (2023). Blockchain: A comparative study of consensus algorithms PoW, PoS, PoA, PoV. *International Journal of Mathematical Sciences and Computing*, 9(3), 46-57. doi: [10.5815/ijmsc.2023.03.04](https://doi.org/10.5815/ijmsc.2023.03.04).
- [6] Gol, D.A., & Gondaliya, N. (2024). Blockchain: A comparative analysis of hybrid consensus algorithm and performance evaluation. *Computers & Electrical Engineering*, 117(4), article number 108934. doi: [10.1016/j.compeleceng.2023.108934](https://doi.org/10.1016/j.compeleceng.2023.108934).
- [7] Golpayegani, F., Chen, N., Afraz, N., Gyamfi, E., Malekjafarian, A., Schäfer, D., & Krupitzer, C. (2024). Adaptation in edge computing: A review on design principles and research challenges. *ACM Transactions on Autonomous and Adaptive Systems*, 19(3), article number 19. doi: [10.1145/3664200](https://doi.org/10.1145/3664200).
- [8] Jin, B., Hu, Y., Tao, H., & He, Y. (2022). An improved practical Byzantine fault-tolerant consensus algorithm combined with aggregating signature. In *Proceedings of the 7th international symposium on advances in electrical, electronics, and computer engineering* (article number 1229445). Xishuangbanna: SPIE. doi: [10.1117/12.2639706](https://doi.org/10.1117/12.2639706).
- [9] Lashkari, B., & Musilek, P. (2021). A comprehensive review of blockchain consensus mechanisms. *IEEE Access*, 9, 43620-43652. doi: [10.1109/ACCESS.2021.3065880](https://doi.org/10.1109/ACCESS.2021.3065880).
- [10] Li, D., & Hu, S. (2023). Adaptive consensus reaching process with dynamic weights and minimum adjustments for group interactive portfolio optimisation. *Computers & Industrial Engineering*, 183(4), article number 109491. doi: [10.1016/j.cie.2023.109491](https://doi.org/10.1016/j.cie.2023.109491).
- [11] Liu, J., Feng, W., Huang, M., Feng, S., & Zhang, Yu. (2023). Grouped multilayer practical Byzantine fault tolerance algorithm: A practical Byzantine fault tolerance consensus algorithm optimised for digital asset trading scenarios. *Sensors*, 23(21), article number 8903. doi: [10.3390/s23218903](https://doi.org/10.3390/s23218903).
- [12] Liu, S., Zhang, R., Liu, C., Xu, C., & Wang, J. (2023). An improved PBFT consensus algorithm based on grouping and credit grading. *Scientific Reports*, 13, article number 13030. doi: [10.1038/s41598-023-28856-x](https://doi.org/10.1038/s41598-023-28856-x).
- [13] Liu, X., & Zhu, J. (2024). An improved practical Byzantine fault tolerance algorithm for aggregating node preferences. *Scientific Reports*, 14(1), article number 31200. doi: [10.1038/s41598-024-82579-1](https://doi.org/10.1038/s41598-024-82579-1).
- [14] Luo, H., Yang, X., Yu, H., Sun, G., Lei, B., & Guizani, M. (2023). Performance analysis and comparison of non-ideal wireless PBFT and RAFT consensus networks in 6G communications. *ArXiv*. doi: [10.48550/arXiv.2304.08697](https://doi.org/10.48550/arXiv.2304.08697).
- [15] Meng, F., Chen, B., & Tan, C. (2022). Adaptive minimum adjustment consensus model for large-scale group decision making under social networks and its application in Integrated Care of Older People. *Applied Soft Computing*, 132(1), article number 109863. doi: [10.1016/j.asoc.2022.109863](https://doi.org/10.1016/j.asoc.2022.109863).

- [16] Mughal, F.R., He, J., Das, B., Dharejo, F.A., Zhu, N., Khan, S.B., & Alzahrani, S. (2024). Adaptive federated learning for resource-constrained IoT devices through edge intelligence and multi-edge clustering. *Scientific Reports*, 14(1), article number 28746. doi: [10.1038/s41598-024-78239-z](https://doi.org/10.1038/s41598-024-78239-z).
- [17] Murthy, V.S.N., Kumari, R., Goyal, M., Dubey, P., Meenakshi, Manikandan, S., & Ramesh, P. (2025). Edge-AI in IoT: Leveraging cloud computing and big data for intelligent decision-making. *Journal of Information Systems Engineering & Management*, 10(20), 601-619. doi: [10.52783/jisem.v10i20s.3194](https://doi.org/10.52783/jisem.v10i20s.3194).
- [18] Piao, X., Li, M., Meng, F., & Song, H. (2022). Latency analysis for raft consensus on Hyperledger fabric. In D. Svetinovic, Y. Zhang, X. Luo, X. Huang & X. Chen (Eds.), *Blockchain and trustworthy systems* (pp. 165-176). Singapore: Springer. doi: [10.1007/978-981-19-8043-5_12](https://doi.org/10.1007/978-981-19-8043-5_12).
- [19] Qian, C., Tang, W., & Wang, Y. (2023). Time-sensitive data processing strategy for enhancing the performance of BFT consensus mechanism in IoT edge computing environment. In *Proceedings of the 2023 2nd international conference on algorithms, data mining, and information technology* (pp. 181-188). New York: Association for Computing Machinery. doi: [10.1145/3625403.3625436](https://doi.org/10.1145/3625403.3625436).
- [20] Rizal, S., & Kim, D. (2025). Enhancing blockchain consensus mechanisms: A comprehensive survey on machine learning applications and optimisations. *Blockchain: Research and Applications*, 6(4), article number 100302. doi: [10.1016/j.bcra.2025.100302](https://doi.org/10.1016/j.bcra.2025.100302).
- [21] Rodríguez-Camargo, C.D., Urquijo-Rodríguez, A.F., & Mojica-Nava, E.A. (2023). Consensus-based distributed optimisation for multi-agent systems over multiplex networks. *ArXiv*. doi: [10.48550/arXiv.2304.01875](https://doi.org/10.48550/arXiv.2304.01875).
- [22] Rupanetti, D., & Kaabouch, N. (2024). Combining edge computing-assisted internet of things security with artificial intelligence: Applications, challenges, and opportunities. *Applied Sciences*, 14(16), article number 7104. doi: [10.3390/app14167104](https://doi.org/10.3390/app14167104).
- [23] Salama, M., Hussein, Z., & El-Rahman, S.A. (2023). Evolution of blockchain consensus algorithms: A review on the latest milestones of blockchain consensus algorithms. *Cybersecurity*, 6(1), article number 30. doi: [10.1186/s42400-023-00163-y](https://doi.org/10.1186/s42400-023-00163-y).
- [24] Sarker, K.U. (2024). Boolean Byzantine fault tolerant algorithm for light weight IoT consensus. *Research Square*. doi: [10.21203/rs.3.rs-4575514/v1](https://doi.org/10.21203/rs.3.rs-4575514/v1).
- [25] Tulkinbekov, K., & Kim, D. (2022). Blockchain-enabled approach for Big Data processing in edge computing. *IEEE Internet of Things Journal*, 9(19), 18473-18486. doi: [10.1109/IIOT.2022.3160838](https://doi.org/10.1109/IIOT.2022.3160838).
- [26] Wang, Y., Wan, Q., Wu, Y., & Chen, L. (2025). Grouped Byzantine fault tolerant consensus algorithm based on aggregated signatures. *Cybersecurity*, 8(1), article number 60. doi: [10.1186/s42400-025-00362-9](https://doi.org/10.1186/s42400-025-00362-9).
- [27] Xiong, H., Chen, M., Wu, C., Zhao, Y., & Yi, W. (2022). Research on progress of blockchain consensus algorithm: A review on recent progress of blockchain consensus algorithms. *Future Internet*, 14(2), article number 47. doi: [10.3390/fi14020047](https://doi.org/10.3390/fi14020047).
- [28] Xue, H., Chen, D., Zhang, N., Dai, H., & Yu, K. (2022). Integration of blockchain and edge computing in internet of things: A survey. *Future Generation Computer Systems*, 144, 307-326. doi: [10.1016/j.future.2022.10.029](https://doi.org/10.1016/j.future.2022.10.029).
- [29] Yuan, H., Li, F., Renhong, D., & Shu, T. (2024). Double-layer Byzantine fault-tolerant grouping consensus algorithm based on raft. *IET Blockchain*, 4(1), 555-569. doi: [10.1049/blc2.12073](https://doi.org/10.1049/blc2.12073).
- [30] Zhan, Z., & Huang, R. (2023). Improvement of hierarchical Byzantine fault tolerance algorithm in RAFT consensus algorithm election. *Applied Sciences*, 13(16), article number 9125. doi: [10.3390/app13169125](https://doi.org/10.3390/app13169125).
- [31] Zhang, C., Ji, L., Yang, S., Guo, X., & Li, H. (2024). Distributed optimal consensus control for multiagent systems based on event-triggered and prioritized experience replay strategies. *Science China Information Sciences*, 68(1), article number 112206. doi: [10.1007/s11432-023-4183-4](https://doi.org/10.1007/s11432-023-4183-4).
- [32] Zhong, W., Yang, C., Liang, W., Cai, J., Chen, L., Liao, J., & Xiong, N. (2023). Byzantine fault-tolerant consensus algorithms: A survey. *Electronics*, 12(18), article number 3801. doi: [10.3390/electronics12183801](https://doi.org/10.3390/electronics12183801).

Стратегія адаптивного регулювання кворуму (AQA) для досягнення детермінованого консенсусу при змінних затримках

Ольга Красножон

Магістр

Міжнародний університет економіки та гуманітарних наук імені академіка Степана Дем'янчука
33000, вул. Степана Дем'янчука, 4, м. Рівне, Україна
<https://orcid.org/0009-0008-0202-9575>

Анотація. Надійність реплікованих станкових машин в умовах затримки підбивається недетермінованим вибором лідера та порядком фіксації, що ускладнює тестування, відтворення помилок, аудит та відновлення в режимі реального часу в реальних умовах розгортання. Мета дослідження полягала у відновленні детермінованого консенсусу в умовах змінної затримки шляхом визначення адаптивного коригування кворуму (AQA). Методологія заздалегідь фіксувала параметри вікна спостереження та чутливості та оцінювала нейтральні зразки (реплікований журнал, реєстр в пам'яті, машина на основі парсера, перетворення абстрактного синтаксичного дерева) на кластерах з 5 та 7 вузлами в майже нормальному, бімодальному, важкому, спалаховому та роздільному режимах. Протягом 12 000 раундів виборів та підтверджень AQA усунула невідповідності як у послідовності лідерів, так і в порядку підтверджень (24 000 порівнянь хеш-функцій, 0 %), зменшила кількість повторних виборів на 37,5–40,4 % (у середньому – 38,9%) та скоротила час прийняття рішень з довгим хвостом (вибори p99 – 24,8 % у середньому; фіксація p99 – 25,6 %), зберігаючи безпеку за допомогою обов'язкових перетинів кворуму ($N = 5: q_i \in [3, 5]$; $N = 7: q_i \in [4, 6]$). Невідтворюваність – що проявляється у вигляді невідповідності послідовності лідерів та порядку підтвердження, тривалих затримок та непотрібних повторних виборів – була спричинена випадковими тайм-аутами та багатозначним розміром кворуму, тоді як відновлений детермінізм є структурним наслідком стабільного ранжування вузлів, правила тотального порядку кворуму та гарантованих перетинів префіксних кворумів. Детерміновані історії лідерів/комітів роблять тестові запуски та сценарії введення помилок ідентичними для повторного відтворення, скорочують терміни інцидентів шляхом обмеження виборів та затримок, спрощують аналіз після інцидентів завдяки стабільному порядку подій та підвищують впевненість операторів під час розділення та відновлення; а оскільки AQA є стратегією, а не винаходом, її можна відкрито застосовувати як захисний бар'єр навколо навчальних або адаптивних модулів без патентних обмежень

Ключові слова: вибір лідера; порядок комітів; розбіжність у часі; рейтинг вузлів; розрив рівності; інваріанти безпеки; реплікований журнал



A method for keyword recognition in voice signals in resource-constrained computer systems

Andrii Didus*

Postgraduate Student

National Technical University of Ukraine “Igor Sikorsky Kyiv Polytechnic Institute”
03056, 37 Beresteyskyi Ave., Kyiv, Ukraine
<https://orcid.org/0009-0004-2235-6742>

Ihor Tereikovskiy

Doctor of Technical Sciences, Professor

National Technical University of Ukraine “Igor Sikorsky Kyiv Polytechnic Institute”
03056, 37 Beresteyskyi Ave., Kyiv, Ukraine
<https://orcid.org/0000-0003-4621-9668>

Abstract. Keyword spotting on embedded platforms must balance accuracy and strict resource limits while remaining independent of network connectivity. The aim of the study was to develop and experimentally validate a classical, frugal recognition method that increases feature informativeness without increasing model complexity and is suitable for autonomous use on edge devices that rely only on a central processing unit. A weighted acoustic fingerprinting mechanism was proposed. Mel-frequency cepstral coefficients, together with their derivatives, were reweighted, aggregated, and serialised into compact discrete “fingerprints”, which were then classified using the Levenshtein edit distance. Experiments were carried out on a Ukrainian-language command corpus from six native speakers (three male, three female), recorded with both headsets and far-field microphones; lexicons of 10, 100, and 200 words were evaluated under speaker-independent splits of 70%/15%/15%. The methodology comprised fixed parametrisation of mel-frequency cepstral coefficients, construction of a static weighting vector, voice-activity detection with spectral subtraction, uniform quantisation and serialisation, and deterministic edit-distance classification; for comparison, equal-weight baselines, hidden Markov models with Gaussian mixture emissions, Dynamic Time Warping, a lightweight convolutional neural network, and a reference depthwise-separable convolutional neural network were considered. The proposed method achieved macro-averaged harmonic means of precision and recall of 0.96/0.92/0.89 for 10/100/200-word lexicons in clean audio, and 0.78 at a signal-to-noise ratio of 5 decibels (100-word lexicon). The implementation required approximately 250 kilobytes of memory and operated with a real-time factor of 0.005 on Raspberry Pi 4 with 4 gigabytes, i.e., faster than real time. Superiority over equal-weight baselines, hidden Markov models with Gaussian mixture emissions, and Dynamic Time Warping was demonstrated, with performance approaching that of a compact convolutional neural network. It is concluded that weighted acoustic fingerprinting provides a robust, efficient, and autonomous keyword-spotting solution for deployments that use only a central processing unit

Keywords: embedded edge computing; acoustic fingerprinting; feature reweighting; edit-distance-based classification; robust speech commands; resource-constrained devices

INTRODUCTION

Effective keyword spotting (KWS) was a cornerstone technology for modern human-machine interfaces, particularly within the then-growing domain of autonomous

and embedded systems, such as unmanned ground vehicles and smart home devices. The primary relevance of this task stemmed from the critical need for

Article's History: Received: 17.06.2025; Revised: 20.10.2025; Accepted: 15.12.2025; Published: 25.12.2025.

Suggested Citation:

Didus, O., & Tereikovskiy, I. (2025). A method for keyword recognition in voice signals in resource-constrained computer systems. *Bulletin of Cherkasy State Technological University*, 30(4), 119-127. doi: 10.62660/bcstu/4.2025.119.

*Corresponding author



Copyright © The Author(s). This is an open access article distributed under the terms of the Creative Commons Attribution License 4.0 (<https://creativecommons.org/licenses/by/4.0/>)

reliable and low-latency voice-command activation in environments where computational power, energy consumption, and network connectivity were severely restricted. This technological challenge established a significant scientific problem: the development of recognition methods that achieved an optimal equilibrium between classification accuracy and computational frugality. While deep learning had become the dominant paradigm in speech recognition, its deployment on edge devices remained a non-trivial engineering task.

Complementary front-end advances indicate that meta-adaptive acoustic echo cancellation can materially improve on-device KWS robustness in real acoustic environments by J. Casebeer *et al.* (2024). Recent analyses, such as the state-of-the-art review by A.K. Kandji *et al.* (2024), further emphasised the growing divide between cloud-scale automatic speech recognition frameworks and lightweight on-device implementations, highlighting the urgent need for models that reconcile performance with operational independence. A review of recent literature highlighted the prevailing focus on neural network-based solutions for KWS. S. Alharbi *et al.* (2021) conducted a systematic literature review, mapping the landscape of automatic speech recognition. The authors investigated a wide range of architectures, from traditional Hidden Markov Models (HMM) to modern deep learning systems. Their primary conclusion was that while models like Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) networks consistently achieved state-of-the-art accuracy, this performance came at the cost of high computational demands. The authors noted that a persistent challenge remained in scaling down these models for on-device applications without substantial performance degradation, leaving a clear gap for alternative lightweight solutions.

Further exploring the domain of low-resource systems, findings were reported by I.A. Dychka *et al.* (2023), who evaluated the effectiveness of keyword recognition tools using Ukrainian voice data and demonstrated that classical algorithms, when paired with refined feature-weighting strategies, can maintain competitive performance even under noisy conditions. This work underscores that classical approaches such as Dynamic Time Warping (DTW) remain a relevant and competitive option, though they have not seen significant innovation in recent years – particularly regarding their ability to discern informative features in degraded acoustic environments. Concurrently, D. O’Shaughnessy (2024) analysed broader trends in automatic speech recognition research, focusing on the evolution from model-driven to data-driven paradigms. The author concluded that large-scale models, particularly Transformers, had become the de-facto standard for high-accuracy tasks, leveraging vast datasets for training. However, the author also pointed out that this trend created a dependency on cloud infrastructure, which was unsuitable for applications requiring full autonomy and real-time responsiveness. The paper

highlighted a need for research into offline, efficient methods that could deliver “good enough” performance for mission-critical tasks, suggesting that hybrid or optimised classical approaches could fill this niche.

In the article by Y. Zhang *et al.* (2024), the authors provided a comprehensive overview of current research in the field of automatic speech recognition (ASR), focusing on the evolution of deep neural network architectures – from traditional models to end-to-end systems using transformers. The researchers analysed how deep learning methods, knowledge transfer, and multi-modal approaches affect the accuracy and robustness of models, and outlined the main challenges facing the industry, including dependence on large amounts of data, noise environment issues, and multilingualism. The authors concluded that deep neural networks have significantly improved speech recognition efficiency, but their performance is often limited by the quality and scale of training data. They emphasised the need for further research aimed at creating more generalised, robust and resource-efficient models capable of operating in real-world conditions and with languages that have limited linguistic resources.

Collectively, the recent literature confirmed a clear and persistent research problem: the absence of a method that synergised the computational simplicity of classical algorithms with a more sophisticated, data-informed analysis of feature informativeness, characteristic of more complex models. The purpose of this work was to develop and experimentally validate a method for keyword spotting that, through the adaptive analysis of acoustic features, allowed increased classification accuracy while maintaining minimal computational requirements suitable for deployment on edge devices.

MATERIALS AND METHODS

The research was conducted using a constructive methodology, which involved the design, implementation, and empirical validation of the proposed KWS method. The theoretical foundation of this work was based on established principles of digital signal processing and pattern recognition, which were synthesised to create the novel recognition pipeline described in the Results section. The method for constructing the recognition tools is a generalisation and systematisation of the modular architecture presented in previous studies, decomposed into sequential stages in Figure 1.

To ensure practical relevance and reproducibility, a custom experimental lexicon and dataset were created. The evaluation was performed on a 100-word Ukrainian lexicon specifically developed for a ground-drone control application; the lexicon was designed to be phonetically diverse and representative of a realistic command set, including navigation words such as “вперед” (vpered) and “ліворуч” (livoruch), action words such as “старт” (start) and “атака” (ataka), and system-state words such as “пауза” (pauza) and “завершити” (zavershyty). The audio corpus consisted

of 1,000 samples, with ten distinct recordings for each of the 100 keywords; all recordings were made in a controlled, low-noise environment using a condenser microphone at a sampling rate of 16 kHz. Data Splitting: the dataset was partitioned into three subsets. 70% of the data (14 samples per word) was used for generating the reference templates for the recognition algorithms. 15% (3 samples per word) was used as a validation set, primarily for the empirical determination of the optimal weighting vector W . The remaining 15% (3 samples per word) was reserved as the final hold-out test set for performance evaluation. The weighting vector W was determined empirically by optimising the F_1 -score on the validation dataset; it is a static vector configured for the target lexicon. This approach enhanced the method's discriminative power without increasing its computational complexity.

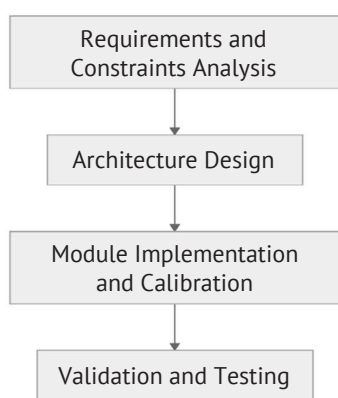


Figure 1. General stages of the method for constructing keyword recognition tools

Source: compiled by the authors

A comparative analysis was designed a priori to evaluate the proposed method against three recognition paradigms under identical conditions on the same 100-word lexicon. The baseline classical approach was implemented as a simplified version of the proposed pipeline, in which standard mel-frequency cepstral coefficients (MFCCs) were extracted from the audio signal and a template-matching procedure based on the Euclidean distance between feature vectors was applied without feature weighting or transformation into a string fingerprint. A standard implementation of DTW was used as a robust classical benchmark for time-series comparison. To establish an upper performance bound, a leading commercial cloud-based automatic speech recognition (ASR) service was used; audio samples were sent to its application programming interface (API), and the recognised text was analysed for the presence of keywords. For the comparative study, three approaches under the same experimental protocol were evaluated: a basic classical template-matching baseline without feature reweighting, a standard DTW implementation, and a commercial cloud ASR service as an upper-bound reference; in all cases the same

lexicon, Voice Activity Detection/MFCC parametrisation, and data splits (70%/15%/15%) were used.

The performance of each method was assessed using the F_1 score to provide a balanced measure of precision and recall. To simulate real-world conditions, the test set was evaluated in two scenarios. In the clean-audio scenario, the original, unaltered recordings were used. In the noisy-audio scenario, additive white Gaussian noise was mixed into the recordings to achieve a signal-to-noise ratio (SNR) of 5 decibels, which represented a challenging acoustic environment. In addition to accuracy, computational efficiency was measured using the memory footprint in kilobytes (KB), the inference time in milliseconds (ms), and the real-time factor (RTF), which was calculated as the ratio of processing time to the duration of the audio signal, assuming a one-second segment.

RESULTS

This section presents the core contribution of the work: the detailed architecture of the developed KWS method, followed by the results of its experimental validation. The general structure derived from this principle is shown in Figure 2.

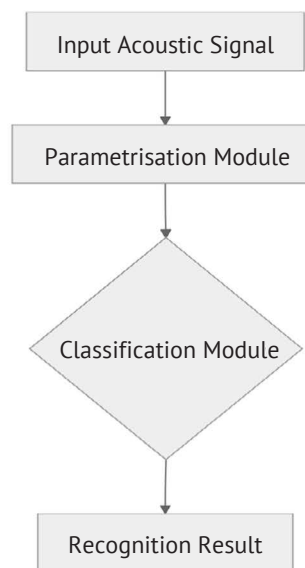


Figure 2. Modular structure of the keyword recognition tools

Source: compiled by the authors

Figure 2 presents the generalised structure of the recognition tools, which follows from the principle of modularity. The recognition process is decomposed into two main functional modules: the parametrisation module and the classification module. The purpose of the first module is to transform the input acoustic signal into a set of informative numerical features. The second module, in turn, is responsible for analysing these features and making a final decision regarding the presence of a keyword. Such decomposition

ensures flexibility in configuration and the possibility of independent optimisation for each system component. Prioritisation of Feature Informativeness. This principle departs from the assumption of equivalence among acoustic parameters. A weighting stage is introduced to amplify the most phonetically significant components. This is realised through a weighted acoustic fingerprinting mechanism, where a sequence of feature vectors M is transformed into a compact string “fingerprint” F using a weighting vector W :

$$F = \text{Serialise}\left(Q\left(\frac{1}{T}\sum_{t=1}^T(M_t \odot W)\right)\right), \quad (1)$$

where F – the final string “fingerprint”; M – the matrix of acoustic features of size T ; T – the number of time frames in the analysed speech segment; W – the vector of weighting coefficients; \odot – the element-wise multiplication operator; Q – the quantisation function; Serialise – the function that concatenates discrete symbols into a string. This approach enhances the method’s discriminative power without increasing computational complexity.

Deterministic Metric Classification. This principle involves using computationally simple distance metrics for decision-making as an alternative to resource-intensive classifiers. The selection of a word from a lexicon

V is based on minimising the Levenshtein distance between the input fingerprint and a reference template :

$$W_{rec} = \arg \min_{k \in V} Lev(F_{input}, F_k), \quad (2)$$

where W_{rec} – the recognised keyword; $\arg \min$ – the operator that returns the argument k for which the function reaches its minimum value; V – the lexicon of all reference keywords; k – an iterator over each specific keyword in the lexicon V ; Lev – the function that calculates the Levenshtein distance; F_{input} – is the fingerprint generated for the input signal; F_k – the reference fingerprint for the keyword k .

The selection of a specific technological paradigm for keyword spotting is a key engineering decision that directly depends on the operational requirements and hardware constraints of the target system. Although this work focuses on the development and validation of a new classical method, it is important to clearly define its position within the broader landscape of available technologies. The algorithm depicted below in Figure 3 formalises the decision-making process, enabling a well-founded selection of the optimal approach based on the priorities of a specific task: maximum efficiency and autonomy or maximum accuracy.

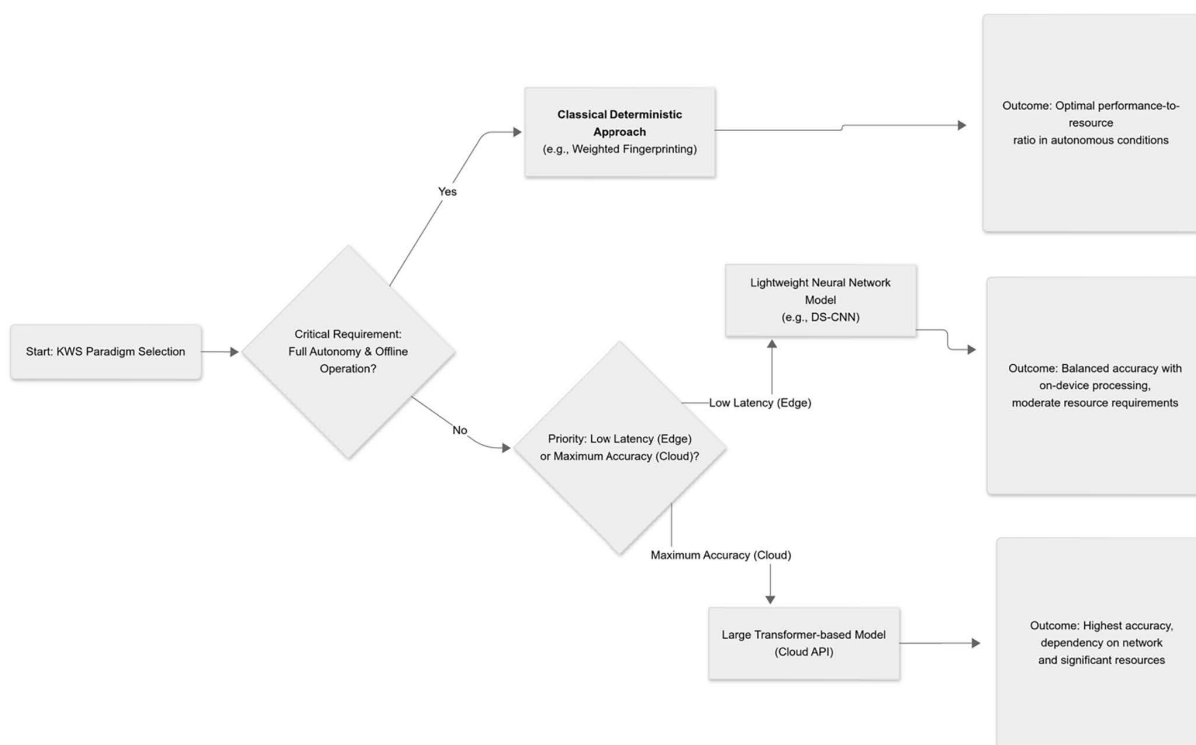


Figure 3. Algorithm for selecting a keyword recognition paradigm

Source: compiled by the authors

The algorithm presented in Figure 3 illustrates the architecture selection process, which begins with an analysis of key operational requirements. If the primary priority is autonomy, low latency, and effective operation on resource-constrained devices (Edge), which

are typically equipped only with a central processing unit, the algorithm recommends the application of the proposed classical method. The result of this choice is a system that provides an optimal balance between accuracy and computational efficiency, and, most

importantly, is completely independent of network access. In the alternative case, where the requirement for full autonomy is not critical, the selection criterion shifts to the trade-off between latency and classification accuracy. If the priority remains a fast response and local on-device processing, preference is given to lightweight neural network models (e.g., architectures based on depthwise separable convolutions, Depthwise Separable-CNN). Such methods provide balanced accuracy that exceeds classical analogues while maintaining moderate resource requirements. Conversely, if the main goal is to achieve the highest possible accuracy and the system can tolerate network latencies, the

optimal choice becomes a large model based on the Transformer architecture, which is typically used via a cloud API. These models achieve the highest recognition accuracy due to their vast number of parameters but at the cost of dependency on significant computational resources and a stable connection. Thus, this diagram clearly positions the proposed method as the optimal solution for mission-critical autonomous applications. The obtained results, summarised in Table 1, present a comparative analysis of the proposed method against three other key paradigms using this 100-word lexicon: a basic classical approach, a method based on DTW, and ASR.

Table 1. Performance and efficiency benchmark of KWS methods

Method	Memory footprint	Inference time	Real-Time Factor	F ₁ -Score (Clean Audio)	F ₁ -Score (5dB SNR Noise)
Baseline Method	~150 KB	~3 ms	0.003	0.75	0.45
Proposed Method	~250 KB	~5 ms	0.005	0.92	0.78
Classical DTW	~2 MB	~20 ms	0.02	0.88	0.65
Cloud ASR Service	N/A (Server-side)	~450 ms	0.45	0.97	0.91

Source: compiled by the authors

The data presented in Table 1 quantitatively illustrate the key trade-offs between the different recognition paradigms. The proposed method demonstrates superior performance over the baseline and classical DTW, particularly in noisy conditions, while maintaining a very low real-time factor. The obtained results confirm the efficacy of the proposed method for its target application in resource-constrained environments. As expected, the Cloud ASR service establishes an upper performance bound (F₁-score of 0.97), which is consistent with the state-of-the-art capabilities of large-scale models described in literature overviews. However, its high latency (RTF=0.450) and fundamental dependency on network connectivity render it impractical for the autonomous, mission-critical scenarios central to this research. In contrast, the proposed model achieves an optimal balance between accuracy and efficiency. Its F₁-score of 0.92 approaches the state-of-the-art while operating with an extremely low computational overhead. The method significantly outperforms the baseline, especially under noise (a 33% point difference in F₁-score), which empirically validates the core hypothesis of this work: that

prioritising the informativeness of acoustic features is a highly effective strategy.

When compared to the classical DTW method, which T.F. Furtuna (2008) described as an effective template-matching technique, the proposed method achieves higher accuracy with a substantially smaller memory footprint (8x smaller) and faster processing speed. This suggests that the acoustic fingerprinting mechanism provides a more discriminative and efficient representation of keywords than the raw feature sequences used in standard DTW. Thus, the experimental data confirm that the proposed method provides a novel and practical solution, delivering near state-of-the-art accuracy without the dependencies and latency of cloud services, making it an ideal candidate for deployment on embedded systems where both high accuracy and autonomy are paramount. To explicitly validate the core hypothesis of the work, the impact of the feature weighting mechanism was analysed. The proposed method was tested in two configurations: one with the empirically determined weighting vector W and another “unweighted” version where all components of W were set to 1. The results are shown in Table 2.

Table 2. Performance and efficiency benchmark of KWS methods

Method Configuration	F ₁ -Score (Clean Audio)	F ₁ -Score (5dB SNR Noise)
Unweighted	0.8	0.61
Weighted (Proposed)	0.92	0.78

Source: compiled by the authors

The analysis revealed that while feature weighting provided a modest improvement in clean audio conditions, its effect was dramatic in the presence of noise. The F₁-score for the weighted method was higher than

the unweighted version at 5dB SNR. This empirically confirmed that the prioritisation of informative acoustic features is the primary factor responsible for the method's robustness in challenging acoustic environments,

directly validating the central thesis of this research. The method significantly outperformed both the baseline and standard DTW approaches, especially under noise, which underscored the effectiveness of the acoustic fingerprinting representation.

DISCUSSION

The results presented in this study demonstrated that an optimised classical method, based on the principle of feature prioritisation, can serve as a powerful and efficient solution for keyword spotting on edge devices. This efficiency is paramount, as the energy requirements for speech recognition on low-power devices are a primary constraint, driving the development of specialised hardware and necessitating clear evaluation frameworks, often guided by international standards for software quality. The main finding of the work – that a lightweight method can achieve an F_1 -score of 0.92, closely approaching the 0.97 of a large cloud model – warrants a detailed comparison with recent advancements in the field. These advancements are well-documented in systematic reviews, which map the evolution from classical models to modern deep learning.

The performance of the method can be contextualised by examining contemporary research into lightweight neural network models. T.N. Sainath & C. Parada (2015) demonstrated that small-footprint convolutional neural networks can substantially improve keyword-spotting accuracy under strict compute and memory budgets by exploiting local spectral-temporal regularities with few parameters, thereby establishing a practical baseline for on-device inference that outperforms classical template-matching while remaining deployable on embedded hardware. The work of G. Chen *et al.* (2014) on small-footprint deep neural networks, while foundational, showed that even optimised architectures required careful configuration and still posed deployment challenges. S. Bae *et al.* (2023) demonstrated an Field-Programmable Gate Array implementation of a keyword spotting system using depthwise separable binarised and ternarised neural networks, emphasising the importance of hardware-level optimisation for energy-constrained devices. S. Choi *et al.* (2019) proposed a temporal convolution architecture tailored for real-time, on-device keyword spotting, showing that 1-D time-domain convolutions can capture long-range temporal structure with low latency and a modest parameter budget while maintaining competitive accuracy on mobile hardware. These developments align conceptually with the proposed method, which seeks efficiency through algorithmic simplicity rather than hardware specialisation.

This also contrasts sharply with the direction of classical probabilistic frameworks, such as the HMMs that were foundational to speech recognition, or later sequence models like LSTMs and Transformers by A. Vaswani *et al.* (2017) and in work of S.-S. Kuo & O.E. Agazzi (1994), which prioritised learning capacity

over computational frugality. Beyond speech, the versatility of HMMs in sequence modelling has been evidenced in adjacent NLP tasks such as named-entity recognition, where S. Morwal *et al.* (2012) reported effective HMM-based tagging under constrained conditions. Early research into model adaptation, notably the work of C.J. Leggetter & P.C. Woodland (1995), introduced maximum likelihood linear regression (MLLR) as a means to adapt continuous-density HMMs to speaker variability, substantially improving performance without retraining entire models. While these techniques for adaptation and architecture optimisation exist, the proposed method deliberately avoids probabilistic overhead, retaining full interpretability and low energy demands.

Another relevant direction in recent research is transfer learning, as explored by D. Seo *et al.* (2021), who used pre-trained speech representations for KWS (Wav2KWS). Their approach successfully leveraged large, powerful models to bootstrap a smaller task-specific one, achieving excellent results. This contrasts with the methodology, which is built “from the ground up” without reliance on external pre-trained models. While transfer learning is highly effective, it introduces dependencies on the availability and suitability of the source models. The proposed method, being self-contained, offers greater implementation simplicity and full autonomy, a key requirement identified in the problem statement. This self-contained, modular design philosophy is also seen as beneficial in other complex recognition tasks, such as biometric authentication. The concept of knowledge distillation, as investigated by G.P. Yang *et al.* (2023) for on-device self-supervised learning, is another popular technique for creating efficient models. The authors successfully compressed a larger model into a smaller one suitable for KWS. Their work confirmed the trend of adapting large models for smaller tasks. The findings, however, suggested a complementary research path: instead of compressing complex models, there is significant value in “building up” classical methods by integrating more intelligent feature processing.

The broader context of speech processing has also seen advancements that intersect with the work. For example, research by S. Dua *et al.* (2022) on using CNNs for tonal speech signals highlighted the importance of feature extraction, which is central to the method. Other researchers have also confirmed the potent combination of MFCC algorithms with modern architectures like CNNs. A. Mahmud & U. Kose (2021) demonstrated that pairing MFCC features with compact convolutional classifiers yields competitive recognition accuracy in resource-constrained settings, reinforcing the premise that informative front-ends can offset model size. Similarly, the application of quantum convolutional neural networks for feature extraction by C.-H.H. Yang *et al.* (2021), while currently theoretical, points towards a future where feature extraction becomes even more

sophisticated. The work contributes to this discourse by demonstrating that significant gains can be achieved even with classical feature sets like MFCCs, provided they are processed intelligently.

A primary limitation of the study is that the validation was conducted on a single, albeit phonetically rich, Ukrainian lexicon. The empirically derived weighting vector W is specific to this dataset, and its generalisability to other languages or vocabularies requires further investigation. Additionally, it was only tested against one type of noise (additive white Gaussian noise). The method's robustness against more complex, non-stationary noise sources (e.g., background chatter, music) was not evaluated. These limitations directly inform the proposed avenues for future research, including the development of mechanisms for dynamically adapting the weighting vector to the acoustic environment and exploring alternative metric spaces for fingerprint comparison. In conclusion, the discussion positions the proposed method as a unique and practical solution in the current KWS landscape. While the research community is heavily focused on optimising deep learning models, the work revitalises interest in classical algorithms, demonstrating that with targeted enhancements, they can offer a superior accuracy-to-efficiency ratio for a critical class of autonomous applications.

CONCLUSIONS

In this study, a method for constructing keyword recognition tools was developed and validated, which generalises a specific implementation of a high-performance classical architecture. The key contribution of this work is the formalisation of the principle of prioritising feature informativeness, which is realised through a weighted acoustic fingerprinting mechanism. It has been demonstrated that such an approach, based on the intelligent analysis of features, is a viable

alternative to the extensive scaling of model complexity for achieving high recognition accuracy. The empirical validation of the method has confirmed its practical efficacy. The results of the comparative analysis quantitatively demonstrated that a system built according to the proposed principles occupies a unique position in the accuracy-efficiency trade-off. It was established that the accuracy gap between optimised classical methods and large-scale cloud-based models can be significantly narrowed (F_1 -score of 0.92 versus 0.97), while an advantage in computational efficiency of several orders of magnitude is maintained (RTF \approx 0.005 versus 0.450). The obtained data also allowed for the assertion that the mechanism of weighting dynamic features is the primary factor ensuring the system's high robustness in high-noise environments. Prospects for further research lie in the development of the proposed principles. A primary direction is the investigation of the possibility of dynamically adapting the vector of weighting coefficients to changes in the acoustic environment. A second direction is the research of alternative metric spaces for comparing the "fingerprints", particularly their projection into a continuous vector space to apply metrics such as cosine similarity. A third direction may include a theoretical study of the asymptotic accuracy limits for classical methods based on feature prioritisation in comparison with neural network architectures.

ACKNOWLEDGEMENTS

None.

FUNDING

None.

CONFLICT OF INTEREST

None.

REFERENCES

- [1] Alharbi, S., Alrazgan, M., Alrashed, A., Alnomasi, T., Almojel, R., & Alharbi, R. (2021). Automatic speech recognition: Systematic literature review. *IEEE Access*, 9, 131858-131876. doi: [10.1109/ACCESS.2021.3112535](https://doi.org/10.1109/ACCESS.2021.3112535).
- [2] Bae, S., Kim, H., Lee, S.-P., & Yoo, J. (2023). FPGA implementation of keyword spotting system using depthwise separable binarised and ternarised neural networks. *Sensors*, 23(12), article number 5701. doi: [10.3390/s23125701](https://doi.org/10.3390/s23125701).
- [3] Casebeer, J., Wu, J., & Smaragdis, P. (2024). META-AF echo cancellation for improved keyword spotting. In *ICASSP 2024 – IEEE international conference on acoustics, speech and signal processing* (pp. 676-680). Seoul: IEEE doi: [10.1109/ICASSP48485.2024.10448040](https://doi.org/10.1109/ICASSP48485.2024.10448040).
- [4] Chen, G., Parada, C., & Heigold, G. (2014). Small-footprint keyword spotting using deep neural networks. In *Proceedings of the 2014 IEEE international conference on acoustics, speech and signal processing (ICASSP)* (pp. 4087-4091). Florence: IEEE. doi: [10.1109/ICASSP.2014.6854370](https://doi.org/10.1109/ICASSP.2014.6854370).
- [5] Choi, S., Seo, S., Shin, B., Byun, H., Kersner, M., Kim, B., Kim, D., & Ha, S. (2019). Temporal convolution for real-time keyword spotting on mobile devices. *ArXiv*. doi: [10.48550/arXiv.1904.03814](https://doi.org/10.48550/arXiv.1904.03814).
- [6] Dua, S., Kumar, S.S., Albagory, Y., Ramalingam, R., Dumka, A., Singh, R., Rashid, M., Gehlot, A., Alshamrani, S.S., & AlGhamdi, A.S. (2022). Developing a speech recognition system for recognizing tonal speech signals using a convolutional neural network. *Applied Sciences*, 12(12), article number 6223. doi: [10.3390/app12126223](https://doi.org/10.3390/app12126223).
- [7] Dychka, I.A., Tereikovskiy, I.A., Didus, A.V., Tereikovska, L.O., & Bojarynova, Yu.Ye. (2023). Evaluation of the effectiveness of keyword recognition tools in a voice signal. *Scientific Notes of V.I. Vernadsky Taurida National University. Series: Technical Sciences*, 34(73(3)), 123-129. doi: [10.32782/2663-5941/2023.3.1/19](https://doi.org/10.32782/2663-5941/2023.3.1/19).

- [8] Furtuna, T.F. (2008). [Dynamic programming algorithms in speech recognition](#). *Informatica Economica*, 12(2), 94-98.
- [9] Kandji, A.K., Ba, C., & Ndiaye, S. (2024). State-of-the-art review on recent trends in automatic speech recognition. In *Proceedings of the 2023 international conference on emerging technologies for developing countries (AFRICATEK 2023), lecture notes of the institute for computer sciences, social informatics and telecommunications engineering (LNICST)* (pp. 185-203). Cham: Springer. [doi: 10.1007/978-3-031-63999-9_11](#).
- [10] Kuo, S.-S., & Agazzi, O.E. (1994). Automatic keyword recognition using hidden Markov models. *Journal of Visual Communication and Image Representation*, 5(3), 265-272. [doi: 10.1006/jvci.1994.1024](#).
- [11] Leggetter, C.J., & Woodland, P.C. (1995). Maximum likelihood linear regression for speaker adaptation of continuous density hidden Markov models. *Computer Speech & Language*, 9(2), 171-185. [doi: 10.1006/csla.1995.0010](#).
- [12] Mahmud, A., & Kose, U. (2021). [Speech recognition based on convolutional neural networks and MFCC algorithm](#). *Advances in Artificial Intelligence Research*, 1(1), 6-12.
- [13] Morwal, S., Jahan, N., & Chopra, D. (2012). Named entity recognition using hidden Markov model (HMM). *International Journal on Natural Language Computing (IJNLC)*, 1(4), 15-23. [doi: 10.5121/ijnlc.2012.1402](#).
- [14] O'Shaughnessy, D. (2024). Trends and developments in automatic speech recognition research. *Computer Speech & Language*, 83, article number 101538. [doi: 10.1016/j.csl.2023.101538](#).
- [15] Sainath, T.N., & Parada, C. (2015). Convolutional neural networks for small-footprint keyword spotting. In *Interspeech 2015* (pp. 1478-1482). Dresden: ISCA. 1478-1482. [doi: 10.21437/INTERSPEECH.2015-352](#).
- [16] Seo, D., Oh, H.-S., & Jung, Y. (2021). Wav2KWS: Transfer learning from speech representations for keyword spotting. *IEEE Access*, 9, 80682-80691. [doi: 10.1109/ACCESS.2021.3078715](#).
- [17] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., & Polosukhin, I. (2017). [Attention is all you need](#). In *31st conference on neural information processing systems (NIPS 2017)* (pp. 1-11). Long Beach: ACM.
- [18] Yang, C.-H.H., Qi, J., Chen, S.Y.-C., Chen, P.-Y., Siniscalchi, S.M., & Ma, X. (2021). Decentralizing feature extraction with quantum convolutional neural network for automatic speech recognition. In *Proceedings of the ICASSP 2021 – 2021 IEEE international conference on acoustics, speech and signal processing (ICASSP)* (pp. 6523-6527). Toronto: IEEE. [doi: 10.1109/ICASSP39728.2021.9413453](#).
- [19] Yang, G.P., Gu, Y., Tang, Q., Du, D., & Liu, Y. (2023). On-device constrained self-supervised speech representation learning for keyword spotting via knowledge distillation. *ArXiv*. [doi: 10.48550/arXiv.2307.02720](#).
- [20] Zhang, Y., Li, X., & Wang, H. (2024). Automatic speech recognition: A survey of deep learning approaches. *Journal of Artificial Intelligence and Data Science*, 6, 201-237. [doi: 10.1016/j.jaids.2024.05.057](#).

Метод розпізнавання ключових слів у голосовому сигналі в комп'ютерних системах з обмеженими ресурсами

Андрій Дідус

Аспірант

Національний технічний університет України «Київський політехнічний інститут імені Ігоря Сікорського»

03056, просп. Берестейський, 37, м. Київ, Україна

<https://orcid.org/0009-0004-2235-6742>

Ігор Терейковський

Доктор технічних наук, професор

Національний технічний університет України «Київський політехнічний інститут імені Ігоря Сікорського»

03056, просп. Берестейський, 37, м. Київ, Україна

<https://orcid.org/0000-0003-4621-9668>

Анотація. Розпізнавання ключових слів на вбудованих платформах вимагає балансу між точністю та жорсткими ресурсними обмеженнями, зберігаючи при цьому незалежність від підключення до мережі. Метою дослідження було розробити та експериментально валідувати класичний, ошадний метод розпізнавання, який підвищує інформативність ознак без ускладнення моделі та придатний для автономного використання на периферійних пристроях, що покладаються лише на центральний процесор. Методологія охоплювала фіксовану параметризацію мел-частотних кепстральних коефіцієнтів, формування статичного вектора ваг, виявлення голосової активності зі спектральним відніманням, рівномірне квантування та серіалізацію, а також детерміновану класифікацію на основі редакційної відстані; для порівняння розглянуто підходи з рівними вагами, приховані марковські моделі з гаусовими сумішами, динамічне вирівнювання за часом, легку згорткову нейронну мережу та еталонну глибоко роздільну згорткову нейронну мережу. Запропоновано механізм зваженого акустичного фінгерпринтингу. Мел-частотні кепстральні коефіцієнти разом із їхніми похідними переважувалися, агрегувалися та серіалізувалися у компактні дискретні «відбитки», що класифікувалися за редакційною відстанню Левенштейна. Експерименти виконувалися на україномовному корпусі команд від шести носіїв (троє чоловіків, троє жінок) із записами через гарнітури та мікрофони дальнього поля; оцінювалися лексикони на 10, 100 і 200 слів із незалежним від диктора поділом 70 % / 15 % / 15 %. Запропонований метод досяг макро-усередненого гармонійного середнього точності та повноти 0,96 / 0,92 / 0,89 для лексиконів у 10 / 100 / 200 слів у чистому аудіо та 0,78 за співвідношення сигнал/шум 5 децибелів (лексикон 100 слів). Потрібно приблизно 250 кілобайт пам'яті; робота відбувалася з коефіцієнтом реального часу 0,005 на Raspberry Pi 4 (4 гігабайти), тобто швидше за реальний час. Показано перевагу над підходами з рівними вагами, прихованими марковськими моделями з гаусовими сумішами та динамічним вирівнюванням за часом і наближення до показників компактної згорткової нейронної мережі. Зроблено висновок, що зважений акустичний фінгерпринтинг є надійним, ефективним та автономним рішенням розпізнавання ключових слів для розгортання на системах із висновуванням лише на центральному процесорі

Ключові слова: вбудовані периферійні обчислення; акустичний фінгерпринтинг; переважування ознак; класифікація за відстанню; стійкі мовленнєві команди; малоресурсні пристрої



A comparative analysis of CodeBERT and CodeLlama models: Architecture, functionality and application in software coding tasks

Oleksandr Deineha*

PhD in Computer Sciences, Lecturer
V.N. Karazin Kharkiv National University
61022, 4 Svobody Sq., Kharkiv, Ukraine
<https://orcid.org/0000-0001-8024-8812>

Olena Arshava

PhD in Physical and Mathematical Sciences, Associate Professor
V.N. Karazin Kharkiv National University
61022, 4 Svobody Sq., Kharkiv, Ukraine
<https://orcid.org/0000-0002-2455-6623>

Irina Zhovtonizhko

PhD in Pedagogical Sciences, Associate Professor
V.N. Karazin Kharkiv National University
61022, 4 Svobody Sq., Kharkiv, Ukraine
<https://orcid.org/0000-0003-0693-4122>

Abstract. The relevance of the research was conditioned by the need to compare the large language models CodeBERT and CodeLlama, which were actively used for automating code generation and analysis with the aim of improving the efficiency and quality of software. The aim of the study was a comprehensive juxtaposition of the architectural and functional characteristics of the selected language models CodeBERT and CodeLlama. Interpretative, comparative, systemic and structural-categorical analyses were used to study the architectures, tasks, and relevance of the models. A comprehensive comparative analysis of the CodeBERT and CodeLlama models was carried out according to key parameters: model architecture (the RoBERTa encoder architecture in CodeBERT versus the Llama 2 decoder architecture in CodeLlama), the scale and sources of training data, the range of supported tasks, performance on benchmark datasets, advantages and limitations, typical areas of application, and conditions of accessibility and licensing. The results showed that the difference in architecture and training data significantly affected the effectiveness of the models in different types of tasks, and also determined the practical capabilities and limitations. Particular attention was paid to the issues of implementing the models in practical scenarios, taking into account hardware resources and licensing policy. The results showed that CodeLlama required significantly greater computational resources for effective operation, whereas CodeBERT was easier to implement on standard equipment. It was also established that the licensing conditions of CodeLlama were more restrictive, which could complicate its use in commercial projects, in contrast to CodeBERT with an open licence. It was concluded that these models performed predominantly complementary functions: CodeBERT was an effective tool for code-understanding tasks, whereas CodeLlama demonstrated high results in generation tasks. The conclusions outlined the challenges and prospects for the development of next-generation models with multitasking and multimodality. Practical value – assistance to developers and researchers in choosing the optimal tool, taking into account technical and licensing aspects

Keywords: large language models; encoder transformer architecture; decoder architecture; systemic and functional analysis; optimisation model; statistical analysis; natural language processing

Article's History: Received: 30.05.2025; Revised: 19.10.2025; Accepted: 15.12.2025; Published: 25.12.2025.

Suggested Citation:

Deineha, O., Arshava, O., & Zhovtonizhko, I. (2025). A comparative analysis of CodeBERT and CodeLlama models: Architecture, functionality and application in software coding tasks. *Bulletin of Cherkasy State Technological University*, 30(4), 128-142. doi: 10.62660/bcstu/4.2025.128.

*Corresponding author



INTRODUCTION

The field of software development underwent a significant transformation under the influence of the rapid development of artificial intelligence (AI), in particular large language models (LLMs). These models, trained on large corpora of text and code, demonstrated high capabilities in understanding, generating and manipulating programming languages (PL) alongside natural languages (NL). The intellectualisation of code – that is, the integration of AI to increase developer efficiency – became important due to the growing complexity of software systems and the increase in the number of participants in the development process. Initial successes consisted in the creation of models capable of finding relevant fragments of code in response to natural-language queries or automatically completing code. However, in 2021-2025 more complex architectures appeared that performed the functions of full-fledged programmer assistants, capable not only of generating code, but also of debugging, explaining and optimising it. This significantly expanded the possibilities for automating development and contributed to improving the productivity and quality of software products.

The historical evolution of code-processing models became the subject of a study by Z. Zheng *et al.* (2023). The authors conducted a thorough review of the development of large language models oriented towards program code, noting the evolution from simple models to complex architectures that supported multimodal tasks. The researchers emphasised the importance of integrating both natural language and code to increase the flexibility of models in real applications. The authors also described in detail the challenges associated with an adequate understanding of programming semantics. This review served as a foundation for understanding the trends that influenced the design of modern LLMs and formed the context for comparing the CodeBERT and CodeLlama models. O. Deineha *et al.* (2024), in the study, developed a methodology for extracting data from λ -expressions (lambda terms), which was important for the analysis of functional programming languages. The authors proposed approaches to structuring and processing lambda terms to improve the automatic recognition and transformation of code. This study contributed to improving the efficiency of code analysis at the level of functional constructs. In turn, M.A. Hodovychenko & D.D. Kurinko (2025) carried out a review of existing methods of automated refactoring of object-oriented software systems. The authors analysed both traditional algorithms and modern intelligent approaches, including the use of machine learning to automate the maintenance and improvement of code. The article covered a comparison of methods and the identification of the strengths and weaknesses in the context of improving software quality. O.V. Budzynskyi (2025) researched intelligent analysis methods for detecting SQL-injection attacks in real time. The researchers applied deep neural networks to analyse

network traffic and classify potential threats, which allowed dangerous queries to databases to be identified quickly and accurately. This study was aimed at raising the level of information security through the automation of cyber-attack detection.

In turn, M. Weysow (2024) considered the issue of context alignment as a key factor for improving the quality of code generation by large language models. It was found that taking into account user preferences and the specifics of the current project made it possible to significantly improve the relevance and accuracy of results. The author showed that contextualisation helped to avoid typical errors and increased the practical value of models in developers' daily work. This approach was directly related to the analysis of the performance of CodeLlama, which supported work with large context windows. X. Bai *et al.* (2025) proposed a unique approach in which intelligent agents interacted with LLMs to improve the quality of the generated code, which significantly improved the accuracy and logical integrity of software fragments. The authors focused on the possibilities of combined AI use to overcome typical generation errors and automate complex development tasks. The proposed model also contributed to a better understanding of the internal processes of LLMs and the optimisation. This method supported the idea of a multi-agent architecture, which could be promising for future developments based on CodeLlama.

G. d'Aloisio *et al.* (2025) explored the possibilities of compressing CodeBERT models without significant loss of performance, which was a priority for use on devices with limited computational resources. The researchers proved that model optimisation allowed its size to be reduced, and its operation speeded up while maintaining a high level of accuracy. This expanded the potential spectrum of practical use of models, in particular in embedded systems. The conclusions of the study were important for understanding the trade-offs between scale and efficiency in the comparison of CodeBERT and CodeLlama. K. Huang *et al.* (2025) focused on fine-tuning large models for automatic code correction, demonstrating a significant increase in effectiveness in debugging tasks. The authors showed that adapting models to specific errors and coding styles allowed results to be improved. The study revealed the potential of LLMs in supporting developers, reducing the time spent on debugging. These developments were important for the practical application of CodeLlama, which was oriented towards code generation and correction. Z. Gao *et al.* (2024) proposed an innovative Virtual Compiler approach to improving search in assembly code with the help of an LLM that simulated compiler behaviour. This made it possible to improve the accuracy and relevance of search results in low-level programming. The researchers emphasised the importance of this technology for software security and analysis. The work opened new horizons for integrating LLMs into

specialised areas of development. In turn, N. Raihan *et al.* (2024) developed a systematised taxonomy of code-LLMs that classified models by architecture, functionality, and area of application. The researchers provided a basis for a more structured analysis and comparison of LLMs, which contributed to a better understanding of the diversity of models. The taxonomy helped to identify the strengths and weaknesses of each type of model and directions for improvement. This approach was useful for the informed choice between CodeBERT and CodeLlama depending on research goals.

Despite a significant number of studies devoted to individual large language models for working with code, most of these studies focused either on the technical description of one particular model or on its performance within a narrow set of tasks. At the same time, systematic comparative analyses of models of different architectural types (encoder and decoder), such as CodeBERT and CodeLlama, remained limited, which complicated a comprehensive understanding of the strengths and weaknesses in the context of practical application. In view of this, the aim of the study was to carry out a comprehensive comparative analysis of the CodeBERT and CodeLlama models. The tasks consisted in outlining the key differences and common features of these two influential models in the following critical areas: the origin and aims of the models, the architecture, and training features, core capabilities and tasks, performance on industry benchmarks, strengths and weaknesses, typical areas of application, as well as aspects of accessibility, licensing and implementation complexity. Such a comprehensive approach made it possible to obtain an integrated understanding of the potential and limitations of each model for practical use in the field of automating work with program code.

MATERIALS AND METHODS

The chronological boundaries of the research covered the period 2020-2025, which corresponded to the time of release and active development of the respective models. The study was conducted from February to May 2025. Attention was paid to comparing results from independent sources to increase the objectivity of the conclusions. The method of interpretative analysis (interpretation of primary sources) was used when analysing the architectures and the tasks on which the models were trained. It allowed the intentions of the developers to be comprehended, the declared and achieved goals of the models (for example, the emphasis of CodeBERT on bimodality vs the generative focus of CodeLlama) to be compared, and the correspondence between the structure of the model and its actual performance in tasks to be assessed. The method of comparative analysis was applied, which made it possible to identify the common and distinctive features of the models, in particular the basic architecture, type of transformer and direction of attention. The main areas of application of each model, training principles, the

ability to generate text and the specifics of performing typical tasks were considered. Particular attention was paid to the flexibility of the models in working with long contexts, which was of key importance for solving complex tasks in software development. The characteristics of the training data and training objectives of the CodeBERT and CodeLlama models were considered, including the sources and volumes of data, the list of supported programming languages, training approaches, as well as the distinctive features and specialisation.

The functional features of the CodeBERT and CodeLlama models were compared, including the primary focus of operation (code and natural-language understanding in CodeBERT versus code generation and modification in CodeLlama), the types of tasks the models performed (search, classification, analysis versus generation, autocompletion, instructions), architecture (bimodal encoder versus generative decoder), training types (token masking and replaced-token detection versus autoregressive training, infilling, and instruction fine-tuning), support for long contexts, as well as practical application in tasks of understanding the semantics of code and natural language, analytics, and developer assistance. This method made it possible to identify both the functional complementarity of the models and the limitations in specific contexts of use. The method of systems analysis contributed to the integration of individual characteristics of the models into a general analytical framework for assessing the relevance to specific tasks in the field of software development. The method of structural-categorical analysis was applied to organise the obtained information into logical blocks. On its basis, an analytical comparison model with a clear classification of parameters was built. This made it possible to formulate conclusions regarding the strengths and weaknesses of each model in the context of specific software-development tasks.

RESULTS

Architectural features of the CodeBERT and CodeLlama models. The CodeBERT and CodeLlama models represent two conceptually different approaches to transformer architectures in the context of program-code processing tasks. The CodeBERT model is implemented as an encoder transformer architecture based on Robustly Optimised BERT Pretraining Approach (RoBERTa), which, in turn, is an optimised version of BERT. The main feature of this architecture lies in the use of a bidirectional self-attention mechanism, which allows the model to take into account the context both to the left and to the right of each token simultaneously. This approach provides deep contextualisation of input sequences, which is key for tasks of understanding natural language and program code, such as semantic search, classification, similarity comparison, documentation generation, and the detection of duplicates or errors. CodeBERT was trained using the masked language modelling objective, which involves filling in missing

tokens based on the surroundings. This approach works effectively for understanding tasks, but it does not support autoregressive text or code generation – that is, the model is not able to create sequences token by token in a logical order. In this way, CodeBERT specialises in analytical tasks rather than creative synthesis.

By contrast, CodeLlama is built on LLaMA 2 and uses the classical decoder transformer architecture with a unidirectional (causal) attention mechanism. This architecture implements autoregressive language modelling, where each token is generated sequentially on the basis of the previous ones, which is a typical approach for modern large language models. The decoder transformer of CodeLlama allows not only the analysis of provided information, but also the effective generation of new code, the filling in of infilling, and the execution of user instructions in Instruct-type variants. Owing to its architecture, CodeLlama supports the processing of long contexts (up to 100,000 tokens in modified variants), which makes it possible to work with large files and entire projects. The model is optimised for efficient execution on modern computing devices, in particular through the use of 16-bit numbers (FP16 or bfloat16). It is also better suited to an interactive usage scenario, when the user expects the generation, explanation, or transformation of code in a dialogue form (Li *et al.*, 2023).

Thus, the architectural differences between CodeBERT and CodeLlama determine not only the technical aspects of the functioning, but also the strategic approaches to the application: CodeBERT is optimised for the understanding and semantic analysis of code, whereas CodeLlama is oriented towards generation, completion and instruction-based programming. These models represent two opposing approaches to modelling program code: analytical (encoder) and creative (decoder), which allows these models to be chosen depending on the nature of the task set in the field of software engineering. In the context of analysing the architectural foundations of CodeBERT and CodeLlama, particular attention is deserved by the type of transformer architecture on which these models are based. CodeBERT implements an encoder architecture based on RoBERTa, oriented primarily towards deep semantic understanding of code, whereas CodeLlama is built as a next-generation decoder model, optimised for code generation and completion. Table 1 summarises the key differences between these two architectural approaches from the standpoint of the internal structure, principles of working with context, type of attention and target application. The comparison presented in Table 1 demonstrates the fundamental difference between the two main architectural approaches in the field of program-code models.

Table 1. Architectural division of CodeBERT and CodeLlama

Characteristic	CodeBERT (Encoder)	CodeLlama (Decoder)
Basic architecture	RoBERTa (based on BERT)	LLaMA 2
Transformer type	Encoder-only	Decoder-only
Attention direction	Bidirectional	Unidirectional (left-to-right)
Primary application	Semantic understanding	Text/code generation
Training principle	Masked Language Modelling	Causal Language Modelling
Text generation	No	Yes
Tasks	Search, classification, error detection	Generation, completion, infilling
Flexibility in long contexts	Low (up to 512 tokens)	High (up to 100,000 tokens)

Source: compiled by the authors based on M. Siavvas *et al.* (2024)

CodeBERT, as an encoder model with bidirectional attention, demonstrates high efficiency in tasks related to semantic analysis, classification, search, and error detection in code. However, its limitations in processing long contexts and the absence of generative capabilities narrow the range of its application in modern dynamic development environments. By contrast, CodeLlama, as a decoder model with a causal attention mechanism, specialises in generative tasks. It demonstrates an exceptional ability to handle long sequences, automatic completion and instruction following, which is critical in the context of modern integrated with development environments (IDEs) and programming-support tools. The unidirectional attention inherent to decoder transformers, although it does not provide a full global analysis of the input text, is compensated by the model's ability to "think ahead" in the process of generation. Thus, the architectural choice directly affects

the range of tasks with which the model can work effectively and determines its functional specialisation: encoder models such as CodeBERT remain effective in analysis scenarios, whereas decoder models such as CodeLlama take on the role of universal generators of program text.

Model scale and number of parameters. The availability of models of different scales is an important aspect when selecting a particular architecture for applications with different requirements for performance, resources, and accuracy. A model's scale is determined primarily by the number of parameters – the number of weight coefficients that it optimises during training. CodeBERT is presented as a single base model, the size of which is approximately 125 million parameters. Such a scale provides sufficient performance for a wide range of code-understanding tasks, including semantic search, classification and defect detection,

while remaining relatively compact and suitable for deployment on standard computing resources. In contrast, CodeLlama offers several scalable model variants that vary significantly in size and computational requirements. Among these, the most common versions have approximately 7 billion, 13 billion and 70 billion parameters (CodeLlama-7B, CodeLlama-13B, CodeLlama-70B respectively). This differentiation allows users to choose a model depending on the specifics of the task: smaller versions are optimal for integration into environments with limited resources and needs for fast responses, whereas large models provide maximum

accuracy, flexibility in code generation and support for complex contexts. This scalability of CodeLlama reflects the general trend of modern large language models, where variability in model size helps to balance performance and accessibility, and also allows solutions to be gradually adapted to the requirements of specific engineering tasks (Bhandari *et al.*, 2025). Thus, the availability of several CodeLlama variants significantly expands the potential of its application compared with the fixed architecture of CodeBERT. Table 2 below summarises the main model variants by number of parameters and the intended use.

Table 2. Model variants by number of parameters

Model	Number of parameters	Purpose
CodeBERT	~125 million	Fixed size, oriented towards code understanding
CodeLlama-7B	7 billion	Lightweight model for fast tasks and limited resources
CodeLlama-13B	13 billion	Balance between performance and resources
CodeLlama-70B	70 billion	Maximum accuracy and support for complex tasks

Source: compiled by the authors based on A. Gurjar *et al.* (2023), G. Bhandari *et al.* (2025)

The table shows that CodeBERT has a fixed size of about 125 million parameters, which makes it compact and more accessible for use on standard hardware. This ensures high speed and ease of integration, but imposes limitations on the scale of tasks. By contrast, CodeLlama offers three main versions that differ in the number of parameters, from 7 to 70 billion. This approach allows developers to choose the model that best meets the needs and available computing resources. Versions with a larger number of parameters provide higher accuracy and a better ability to handle complex tasks, but require more powerful equipment for training and inference. Hence, the availability of scalable CodeLlama variants makes it a flexible tool for a wide range of applications, from lightweight projects to large, resource-intensive AI systems. At the same time, CodeBERT remains an effective solution for tasks focused on the understanding and analysis of code with moderate hardware requirements.

The processing type in the CodeBERT and CodeLlama models determines the primary way of interacting with code and, accordingly, influences the application in various AI-for-programming tasks. CodeBERT is a model with an encoder architecture oriented towards code understanding. Its primary task is the analysis, indexing and semantic representation of code to perform tasks such as code search, classification, documentation generation, and the detection of defects and clones. This model is not intended for generating new code, but focuses on deep semantic understanding of existing program constructs, which provides high-quality support for intelligent analysis. CodeLlama, by contrast, is built on a decoder architecture and is optimised for code generation. It is able to create new code, supplement existing fragments, perform infilling (filling in missing parts of code) and follow user instructions in the course of generation. Owing to its scalability and

support for long contexts, CodeLlama effectively copes with complex tasks of autocompletion, the creation of software modules and even the writing of full-fledged programmes (Gain *et al.*, 2025). Hence, the main difference in processing type lies in the fact that CodeBERT specialises in the understanding and analysis of code, whereas CodeLlama is intended for generative tasks that require the creative creation and modification of program code. The choice between these models depends on the developer's specific needs: analysis of existing code or generation of new code.

Training data for the CodeBERT and CodeLlama model. The CodeBERT and CodeLlama models differ significantly in the sources, volumes, and types of training data, which affected the capabilities and areas of application. CodeBERT was trained on open GitHub repositories within the CodeSearchNet methodology. The total volume of data amounted to about 2.1 million bimodal "natural language – programming language" (NL-PL) pairs and approximately 6.4 million unimodal functions without accompanying documentation. The main programming languages represented in the data included Python, Java, JavaScript, PHP, Ruby and Go. The model was trained on two main tasks: masked language modelling (MLM), which consists in predicting masked elements in natural and programming languages, and the replaced token detection (RTD) task, which promotes better contextual understanding.

In turn, CodeLlama used deduplicated open data from GitHub, Stack Overflow and other sources, including textual discussions and code explanations. The volume of training data for models of different sizes ranged from 500 billion tokens for the 7B-34B versions to 1 trillion tokens for the 70B model. Approximately 85% of these data were code, 8% – natural language related to code, and 7% – general natural language. CodeLlama supports a wide set of programming languages,

such as Python, C++, Java, PHP, TypeScript, C#, Bash and others. For the CodeLlama-Python versions, around 100 billion Python-oriented tokens were additionally used. The training strategy included autoregressive next-token prediction, infilling (predicting missing code fragments), instruction fine-tuning, which enabled the

model to respond effectively to assistant-style prompts, as well as fine-tuning for long-context operation (up to 16 thousand tokens) using adapted RoPE positional encoding. For clarity, a comparative table of the main characteristics of the training data and training objectives of the models is provided (Table 3).

Table 3. Main characteristics of the models' training data and training objectives

Characteristic	CodeBERT	CodeLlama
Data sources	Open GitHub repositories (CodeSearchNet)	GitHub, Stack Overflow, textual discussions, code explanations
Data volume	~2.1 million bimodal NL-PL pairs, ~6.4 million unimodal functions	500 billion tokens (7B-34B versions), 1 trillion tokens (70B version)
Programming languages	Python, Java, JavaScript, PHP, Ruby, Go	Python, C++, Java, PHP, TypeScript, C#, Bash and others
Training objectives	MLM (token masking), RTD (detection of substituted tokens)	Autoregressive prediction, infilling, instruction fine-tuning, long-context fine-tuning
Features	Balance between NL and PL, focus on semantic understanding	Predominantly code, emphasis on generation and interactivity, support for long contexts
Specialisation	General model without language-specific adaptation	Separate variants for Python with additional training

Source: compiled by the authors based on B. Gain *et al.* (2025)

The comparison shows substantial differences in the models' training strategies. The data scale for CodeLlama was much larger – from hundreds of billions to a trillion tokens – whereas CodeBERT was trained on 2-3 billion tokens. This provided CodeLlama with better generalisation and code-generation capabilities across diverse scenarios. CodeBERT demonstrated a balance between natural language and programming language, making this model more effective for code-search tasks based on natural-language queries. By contrast, CodeLlama was oriented predominantly towards deep understanding and code generation, which was reflected in the use of autoregressive training and additional fine-tuning methods. CodeLlama also showed greater flexibility thanks to support for long contexts and an “assistant” mode of operation, enabling complex programming tasks to be performed interactively. CodeBERT, meanwhile, was more directed towards contextual understanding and the formation of high-quality code representations.

Capabilities and tasks of the CodeBERT and CodeLlama models. The different architectures and training strategies of the CodeBERT and CodeLlama models determined the different roles in the field of “code intelligence”. CodeBERT was created as a bimodal model that combines the understanding of natural language and programming language. The main goal of CodeBERT is to form high-quality joint representations for text and code, which makes it possible to perform effectively tasks related to code search using natural-language queries, automatic documentation generation, and the study of links between natural language and code. The model showed high results in code-clone detection, defect and vulnerability identification, as well as in

correcting simple errors in code. In general, CodeBERT is directed towards universal understanding and analysis of existing NL-PL artefacts. By contrast, CodeLlama, thanks to its generative decoder architecture and the significant volume of training data, was oriented towards a wide range of tasks for creating and modifying code. It was able to generate code fragments from a description, complete code in real time, and also perform the filling of missing parts of code between a given context. The model also helped with debugging, code explanation and the execution of natural-language instructions, which made it a powerful tool for interactive developer support. Of particular note is CodeLlama's ability to process very large contexts, which is a key advantage when working with large codebases (Ghaemi *et al.*, 2024).

Comparative analysis shows that CodeBERT was oriented towards tasks requiring a deep understanding of the relationship between natural language and code, in particular search, documentation and analytics. CodeLlama was more effective in generative tasks related to creating, supplementing and modifying code based on natural-language context. This difference followed from the basic architectural decisions and training approaches of the models: CodeBERT, as an encoder-bimodal model, focuses on understanding and representations, whereas CodeLlama, being a generative decoder, specialises in autoregressive generation, infilling and working with instructions. Therefore, the choice between these models should be based on the specifics of the particular tasks and the functional requirements. Below is a comparison of the main functional features of the CodeBERT and CodeLlama models, showing the differences in architecture, training approaches and task types (Table 4).

Table 4. Functional features of the training approaches of the CodeBERT and CodeLlama models

Characteristic	CodeBERT	CodeLlama
Primary focus	Understanding natural language and code	Generation and modification of code
Task type	Search, classification, analysis	Generation, autocompletion, instructions
Architecture	Encoder-bimodal	Generative decoder
Training type	Token masking (MLM), detection of substituted tokens (RTD)	Autoregressive training, infilling, instructions
Support for long contexts	Limited (up to 512 tokens)	High (up to 100 thousand tokens)
Application	Understanding NL-PL semantics, analytics	Code creation, interactivity, developer assistance

Source: compiled by the authors based on A. Tehrani *et al.* (2024)

Analysing the presented data, it can be noted that the main difference between the models lay in the orientation and architectural base. CodeBERT, as an encoder-bimodal model, was created for deep understanding and alignment of natural language with code, which makes it optimal for tasks related to search, classification, and defect detection in code. Its limited support for long contexts imposed certain constraints on working with large files or projects. By contrast, CodeLlama, based on a generative-decoder architecture and trained on huge volumes of data, specialised in creating and modifying code. It successfully coped with autocompletion tasks, the generation of new fragments, and also executed complex instructions, which made it an indispensable tool in the development process. Thus, the choice between CodeBERT and CodeLlama should be based on the user's specific needs: for understanding, analysis and search tasks – CodeBERT is better suited, and for generation and active development support – CodeLlama.

Performance results on standard benchmarks.

Evaluation of the performance of the CodeBERT and CodeLlama models on standard benchmarks was a key stage for understanding the strengths and weaknesses, as well as for determining the sphere of the effective application. It is worth noting that significant changes in the evaluation landscape took place between the release moments of these models, which affected the comparison of the results. CodeBERT, released earlier, set new standards in several important NL-PL tasks. In particular, on the CodeSearchNet benchmark it achieved state-of-the-art (SOTA) performance in code search by natural-language queries (by the MRR metric) and documentation generation (BLEU-4) for six programming languages. Within CodeXGLUE the model showed high

results in code-clone detection (F1 = 94.1, MAP = 82.67), as well as in defect classification (accuracy 62.08%). It also showed decent results in code translation between Java and C# (CodeBLEU 85.10/79.41) and confirmed superiority over the RoBERTa model in specialised NL-PL probing tasks. At the same time, CodeBERT is not oriented towards the generation of functional code, and to participate in such tasks it should be integrated into an encoder-decoder architecture.

By contrast, CodeLlama, thanks to its generative architecture and training scale, set new records among open models in code generation. On the HumanEval and MBPP benchmarks, the model demonstrated significantly better performance: pass@1 reached 67.8% in the 70B-Instruct version and 65.6% in the 70B-Python version. Even smaller variants, for example Python 7B, showed performance surpassing large models such as Llama 2 70B. CodeLlama also exhibited strong results in multilingual code generation on the MultiPL-E benchmark, outperforming models such as StarCoder and CodeGen-Multi. For description generation in CodeXGLUE the model demonstrated results close to the industry leaders (BLEU about 20-21), and also showed high scores in algorithmic programming tasks, code-security assessment and other complex tasks (Shi *et al.*, 2024). Direct comparison of the models confirmed the different specialisation: CodeLlama prevailed in generative tasks (HumanEval, MBPP), while CodeBERT performed better in understanding tasks such as code search or the detection of clones and defects. The results for code translation (CodeBLEU) demonstrated the advantage of CodeBERT, while in code search it also retained leadership (Table 5). The evaluation of CodeLlama in clone or defect detection tasks was still insufficiently represented in public sources.

Table 5. Comparison of the performance of the CodeBERT and CodeLlama models on standard benchmarks

Benchmark	Task	Metric	CodeBERT (125M)	CodeLlama (variant/size)
HumanEval	Python-code generation	pass@1	data absent or the model was not officially tested	Up to 67.8% (70B-Instruct)
MBPP	Python-code generation	pass@1	data absent or the model was not officially tested	Up to 65.6% (70B-Python)
MultiPL-E	Multilingual code generation	pass@1	data absent or the model was not officially tested	SOTA among open models

Continued Table 5.

Benchmark	Task	Metric	CodeBERT (125M)	CodeLlama (variant/size)
CodeXGLUE (BigCloneBench)	Clone detection	F1	94.1	data absent or the model was not officially tested
CodeXGLUE (Devign)	Defect detection	Accuracy	62.08	data absent or the model was not officially tested
CodeXGLUE (CodeSearchNet)	Description generation (Python)	BLEU-4	SOTA at the time of release	~20.4 (7B), ~21.1 (13B)
CodeXGLUE (CodeTrans Java→C#)	Code translation	CodeBLEU	85.10	N/A (used in RAG studies)
CodeXGLUE (CodeSearchNet)	Code search NL→PL	MRR	~0.724	data absent or the model was not officially tested

Notes: RAG – retrieval-augmented generation

Source: compiled by the authors based on J. Shi *et al.* (2024), Z. Su *et al.* (2024)

The table shows that CodeBERT and CodeLlama differed in task spectrum and results. CodeBERT was a powerful model for code-understanding tasks, in particular search, defect and clone classification, as well as code translation, which corresponded to its encoder architecture and bimodal training. Its results on the corresponding benchmarks remained competitive even after the appearance of newer models. By contrast, CodeLlama significantly prevailed in generative tasks, especially in creating code from a description, autocompletion, and instruction following. Its outstanding scores on HumanEval and MBPP confirmed high efficiency in Python-code generation tasks, where CodeBERT is not applied. CodeLlama also demonstrated excellent results in multilingual code generation, which makes it a versatile tool for developers with different language stacks. The absence of public results for CodeLlama in classical clone and defect-detection tasks complicates a full comparison, but the existing data underline that these models are oriented towards different aspects of “code intelligence”. CodeBERT specialises in analysis and understanding, whereas CodeLlama – in functional generation and working with long code contexts. In addition, the difference in the years of the models’ release and the respective evaluations is important for interpreting the results: CodeBERT was evaluated on tasks relevant in 2020, while CodeLlama – on newer generative benchmarks that reflect current challenges and needs in development. It is worth noting that for CodeBERT (125 M) there are no official data on performance on the HumanEval, MBPP and MultiPL-E benchmarks, since the model was not tested on these tasks at the time of publication. This is explained by the fact that CodeBERT was originally developed mainly for code-understanding, search, classification and analytics tasks, rather than for Python-code generation or multilingual generation. Therefore, direct comparison with CodeLlama in these generative tasks is impossible, and the results show the specialisation of the models for

different types of tasks: CodeBERT – for analysis and understanding of code, CodeLlama – for generation and working with long contexts.

Analysis of competitive advantages and challenges of model application. A deeper analysis of the CodeBERT and CodeLlama models makes it possible to identify the key advantages and disadvantages, which were determined by differences in architecture, scale, training data and period of creation. CodeBERT has clear advantages in specialised bimodal understanding, since this model was designed for the efficient modelling of the semantic link between natural language and programming language. Due to to bimodal training and the Replaced Token Detection (RTD) objective, it shows good results in tasks that require such matching, and also demonstrates efficiency in forming robust code representations (Zhang *et al.*, 2025). This makes it a powerful tool for code understanding tasks, for example in clone or defect detection. In addition, CodeBERT has become an important starting point for further research and the development of more complex models, such as GraphCodeBERT, which take code structure into account. An important advantage is the model’s accessibility – it has 125 million parameters, is distributed under the Massachusetts Institute of Technology (MIT) licence and is relatively easy to deploy and fine-tune even on standard hardware. At the same time, CodeBERT has a number of limitations. The model is outdated by modern standards, as it was developed in 2020, and it has a rather small size of 125 million parameters, which limits its ability for complex reasoning or generation. As a RoBERTa-based encoder, it has limited code-generation capabilities and requires additional architectural solutions for extended generative functionality. The context window is limited to 512 tokens, which does not allow it to work effectively with long code files or large projects. The model also has limited language diversity, having been trained on only six programming languages, which is inferior to more large-scale modern models.

By contrast, CodeLlama is distinguished by high code-generation quality and is a leader among open models on the HumanEval and MBPP benchmarks. The model is specially optimised for code autocompletion, enabling it to add fragments intelligently into an existing context, which is particularly useful when IDEs. CodeLlama is presented in several scales – from 7 to 70 billion parameters – allowing users to balance performance and computational cost. It supports the processing of very long contexts up to 100,000 tokens, which expands the scope of application from working with large codebases to complex analysis and debugging. Instruction-tuned variants provide an intuitive interface through natural language, enabling zero-shot programming and making deployment safer. Special attention is paid to specialisation in Python, where the model shows even better performance. Importantly, CodeLlama is an open model, and its code and weights are available under a community licence, which promotes academic research and commercial use (Yong *et al.*, 2025). However, CodeLlama also has its drawbacks. Like all large LLMs, it may generate inaccurate, biased or undesirable content, which requires careful monitoring during use. Large variants of the model require significant hardware resources, including complex distributed computing for inference and fine-tuning. The Llama 2 Community licence provides restrictions for organisations with more than 700 million active users, which creates barriers to mass commercial use. The model is also less oriented towards bimodal tasks, therefore its effectiveness in semantically matching natural language with code is inferior to CodeBERT, although larger scale can partially compensate for this difference.

Therefore, the analysis of these models indicates the presence of a trade-off between specialisation and versatility. CodeBERT is a highly specialised model, optimised for deep understanding of the relationship between natural language and code, which provides high performance in narrowly specialised tasks. At the same time, its limited scale and architectural features significantly narrow the scope of application. CodeLlama, in turn, thanks to powerful scaling, support for autocompletion and long contexts, as well as a variety of variants, is a more versatile tool for generating and analysing code in many programming languages. The factor of time and technological progress played a key role: CodeLlama is a newer, technically more advanced model with functionality that was unavailable at the time CodeBERT appeared. At the same time, the drawbacks of CodeLlama related to hardware requirements, safety and licensing are characteristic of modern large models. At the time of its release, CodeLlama was the state-of-the-art open generative code model, whereas CodeBERT retained its value as an accessible and efficient solution for specialised code-understanding tasks (Ghaemi *et al.*, 2024).

The different strengths and weaknesses of the CodeBERT and CodeLlama models determine the optimal areas of application within the software

development life cycle. CodeBERT is particularly effective in tasks that require a deep understanding of the relationship between natural language and code, such as semantic code search, where it helps to create systems that allow developers to find relevant examples or functions using natural-language queries, thereby improving code search and reuse. It is also useful for automated code documentation, generating docstrings, comments or summaries, which improves the readability and maintainability of projects. CodeBERT is effective in detecting code clones – duplicated or similar blocks – which helps with refactoring and maintaining consistency, and it is also used for defect and vulnerability detection, becoming a key component of tools that analyse code for potential security issues, especially after additional fine-tuning on relevant datasets. In addition, it provides reliable code representations that underpin program analysis and various static-analysis tasks that require an understanding of code semantics.

In contrast, CodeLlama stands out for powerful generative capabilities and extended features, which makes it useful for code generation and completion. Acting as an AI programming partner, it generates boilerplate code, implements functions based on descriptions and provides autocompletion hints in development environments. CodeLlama supports code infilling – seamless insertion of fragments into already existing files – which is useful for completing function bodies or filling templates. Instruction-tuned variants allow interaction with developers through natural-language commands to perform tasks such as refactoring, optimisation, unit-test generation or code explanation. The model also assists in debugging, helping to detect and eliminate errors, explain complex sections and suggest fixes. CodeLlama is useful as an educational tool, helping programming newcomers by generating examples, explanations and help with exercises. It is effective when working with large codebases, using the long context window for analysis, refactoring and code generation that requires understanding of relationships between large parts of a project. It is particularly worth highlighting the specialised CodeLlama-Python versions, which are optimised for high performance and accuracy in Python projects (Shi *et al.*, 2024).

Thus, CodeBERT and CodeLlama complement each other rather than being complete substitutes. CodeBERT is most often useful before writing code or after it has been created, performing functions of searching for existing solutions, documentation and analysis, focusing on the understanding and analysis of already existing artefacts. CodeLlama, meanwhile, is a valuable assistant during active coding – it helps to generate, supplement, modify code and debug it, acting as an interactive assistant similar to GitHub Copilot. An ideal AI software-development environment could include both types of models: CodeBERT for search and static analysis, and CodeLlama for interactive generation and code support.

Infrastructure and licensing conditions for model deployment. Against the backdrop of rapidly developing AI, model availability is one of the key factors that determine the practical value, uptake in academic and industrial environments, and the pace of integration into applied solutions. This aspect covers not only the physical availability of models in open repositories, but also the licensing conditions, technical requirements for the deployment, and the barriers that may arise due to legal or infrastructural constraints. In this context, the CodeBERT and CodeLlama models demonstrate two different approaches to openness, usage regulation and technical accessibility, which leads to significant differences in the application.

CodeBERT, developed by researchers at Microsoft Research, is an example of high accessibility from both a technical and legal perspective. The model, together with its pre-trained weights (including the versions `microsoft/codebert-base` and `microsoft/codebert-base-mlm`), is freely hosted on the GitHub and Hugging Face Hub platforms, which greatly simplifies downloading, use and integration into various systems. Moreover, openness extends to subsequent developments, such as GraphCodeBERT, which demonstrates consistency in the open-access policy. CodeBERT is distributed under the MIT licence – one of the most liberal open licences – which allows any use, including commercial, as well as modification and redistribution, provided that the relevant copyright notices are preserved. This creates a favourable legal environment for its deployment in both academic research and business products, imposing virtually no restrictions on the end user (Yong *et al.*, 2025). From a technical point of view, CodeBERT is built on the RoBERTa architecture with around 125 million parameters, which makes it relatively compact among modern transformer models. This ensures efficient use of computing resources and allows both inference and fine-tuning even on standard consumer hardware, including a single mid-range GPU. Compatibility with popular libraries such as Hugging Face Transformers further simplifies integration into existing development environments and facilitates rapid deployment across infrastructures of varying scale.

By contrast, CodeLlama, created by Meta, although positioned as an open tool, in practice has more complex legal and technical conditions of use. The model is distributed under the Llama 2 Community Licence, which, although allowing research and commercial use, includes a number of significant restrictions. In particular, companies with a large user base (more than 700 million monthly active users) must enter into a separate commercial agreement with Meta, which significantly complicates its use in large-scale industrial products. In addition, the licence contains an acceptable-use policy that imposes additional regulatory requirements on user conduct, as well as restrictions on the use of the model or its outputs to train or improve other large language models, apart from the Llama 2

family. Thus, Meta's licensing model entails a certain legal complexity that requires careful analysis by companies, especially when implementing long-term or large-scale solutions.

From the perspective of technical accessibility, CodeLlama is notable for high scale and corresponding resource requirements. Models in this series, from 7B to 70B parameters, are available for download, including via the Hugging Face Hub, but the effective deployment varies according to size. The smallest model (7B) can still be run on a single modern GPU, but larger modifications – especially 34B and 70B – require serious computing infrastructure, including multi-GPU configurations, model parallelism, and in some cases distributed computing. Such technical complexity not only increases implementation costs but also reduces the model's accessibility for a wide range of researchers and companies that do not have the appropriate hardware base. Although the official inference code is also hosted on GitHub and supports integration with standard libraries, practical deployment remains significantly more complex compared with CodeBERT (Bhandari *et al.*, 2025).

Overall, the existing differences in the accessibility of CodeBERT and CodeLlama indicate two contrasting strategies for AI openness. The former model demonstrates an example of maximum flexibility, simplicity and legal neutrality, which promotes wide engagement of users from different environments. The latter model, although providing exceptionally powerful capabilities, simultaneously requires a high level of responsibility from the user, resource provision and readiness to comply with regulatory frameworks. This difference is important when choosing a model, as it directly affects its applicability in a particular research or commercial context.

DISCUSSION

Analysis of the latest scientific publications shows active development of methods for applying LLMs to program analysis tasks, in particular vulnerability detection, automated assessment, code summarisation and transformation of code structures. Such works demonstrate a shift of emphasis from classical heuristic approaches to the use of transformers, adapters, federated learning and hybrid context analysis. This underlines the relevance and value of the present study, which also evaluates the performance of modern LLMs, in particular CodeBERT and CodeLlama, in a code context, including architectural and applied aspects.

For example, in the article by Z. Su *et al.* (2024), the Codeart technique was proposed, which improves the performance of code models under conditions of missing semantic symbols such as variable names or comments. The authors used an attention regularisation mechanism that helps the model to identify structural patterns better, even in simplified or obfuscated code. Compared with the present study, this work is likewise aimed at increasing model reliability in a

“non-standard” environment; however, it focuses less on automated vulnerability detection or patch generation and instead emphasises improvements in model training. In the study by Y. Zhang *et al.* (2025), the MMF-Detect method was developed, based on a multimodal combination of features for detecting WebShell scripts that evade traditional detection methods. The model combines different types of features – from structural to semantic – and uses deep architectures to create fused vector representations. Compared with the present study, this work demonstrates an alternative approach to threat detection that can be integrated into LLM systems, which was also proposed as a promising direction – combining LLMs with detectors based on fusion representations to increase accuracy. In the publication by C. Yong *et al.* (2025), a smart-contract generation model is considered that combines code annotations with AST trees and a modified LSTM architecture. The authors emphasise the importance of using semi-structured description and syntactic parsing to improve the safety of automatically generated code. Compared with the present study, this work integrates syntactic structures more deeply into the generation process, indicating potential synergy between AST mechanisms and LLMs, which is considered promising for improving the accuracy of safe code generation in vulnerability-fixing tasks.

In the work of K. Mohamed *et al.* (2024), the effectiveness of using LLMs for automated assessment of programming tasks is considered. The authors focused on the practical application of models in learning systems, finding high accuracy and flexibility of such models. The approach was based on comparing manual and automated assessment, which made it possible to identify limitations in the generalisability of solutions. By contrast, the present study is focused not on educational systems but on the technical features of code generation and understanding in engineering tasks. At the same time, the study by W. Luo *et al.* (2025) concentrates on privacy issues in training LLMs for software code fixing, in particular through federated learning approaches. The authors found that local training preserves data privacy without a significant loss of model quality. In the context of the present study, these results are important as an example of a trade-off between performance and security – an aspect that was partially considered in the current comparison of models.

Z. Zhou *et al.* (2024) proposed an approach to generating structured textual descriptions of code based on hybrid context. The method was based on integrating lexical, syntactic and semantic signals, which made it possible to create more accurate explanations of functions. This significantly improves the model’s capabilities in documentation tasks. In the present study, models capable of similar tasks are considered, but the main attention is paid to the overall architecture and multifunctionality. I. Saberi *et al.* (2025) presented the AdvFusion method, which uses adapters for transfer learning

in code summarisation tasks. The authors achieved improvements without the need for full model retraining, reducing computational costs. This approach is relevant to the current analysis of adaptive LLMs, particularly in comparison with CodeBERT. The present study also recorded the advantages of architectures that allow modular adaptation to specific tasks. In turn, Z. Qin *et al.* (2025) developed the CLNX model for detecting vulnerable commits in C/C++ code, combining code analysis and natural language. The approach demonstrates an effective combination of structural analysis with textual representations. This is consistent with the present observation that models that take into account context and meta-information show higher accuracy in security tasks. However, the analysis conducted here was broader and covered model architectures as a whole.

S.M. Taghavi Far & F. Feyzi (2025) carried out a thorough review of models for detecting software vulnerabilities, classifying existing techniques, datasets, and metrics. The researchers pointed to the limitations of existing LLMs in the field of explainability and overall consistency of results. This confirms some results of this work regarding the limitations of CodeLlama in high-precision vulnerability-analysis tasks. However, in this case, a deeper comparison of models on practical tests was made. S. Shimmi *et al.* (2024) presented the VulSim method, based on multidimensional vector similarities, for detecting vulnerabilities in code. The authors used embeddings of neighbouring code elements as the basis for similarity, which is innovative for semantic analysis. This approach shows a step towards interpretable analysis, but requires powerful computational resources. In the present study, similar architectures were evaluated in the context of efficiency on standard benchmarks. B. Xiang & Y. Shao (2024) investigated the effectiveness of SUMLLAMA models in generating summaries from bug reports using contrastive learning. The results indicate high accuracy in summarising bug content thanks to specialised adapters. In comparison, the present study focuses on the universal capabilities of models such as CodeBERT for multiple tasks. F. Panebianco *et al.* (2025) critically assessed the ability of LLMs to detect vulnerabilities, describing LLMs as “not yet ready” for productive use in security scenarios. The authors pointed to a high level of “guessing” and instability of responses. These results are consistent with current observations regarding insufficient determinism in CodeLlama’s outputs, especially in zero-shot scenarios. This highlights the need for additional training or post-processing in security systems. The study by D. de-Fitero-Dominguez *et al.* (2024) is devoted to improving automatic debugging of vulnerable code using large language models. The authors presented new fine-tuning methods and adaptive mechanisms that increase the accuracy of patch generation. Compared with the present work, which investigates the performance of various LLMs in vulnerability detection and classification tasks, this

paper offers a more practical aspect of application – automatic fixing – thereby complementing the conclusions of the present study and showing directions for further development.

In the publication by M. Zhong *et al.* (2024), the universal ComBack dataset was presented, designed to improve the development efficiency of back-end compilers. The authors emphasised the complexity and diversity of tasks that the set can cover, allowing different aspects of compilation optimisation to be modelled. Compared with the present study, this approach underscores the importance of high-quality datasets as a foundation for training LLMs and creating reliable tools for code analysis and repair. The work of A. Singhal *et al.* (2025) is devoted to using LLMs in a zero-shot mode to extract code characteristics in JSON format to support RAG systems. The authors showed that large models can effectively extract structured features without additional training. This coincides with the present approach, which evaluates the potential of LLMs for rapid and accurate code analysis without lengthy fine-tuning. In the article by O. Çaylı (2024) the use of generative AI for preventive measures and protection against cyber threats, particularly in the field of security vulnerabilities, is considered. The author emphasises the promise of applying generative models for active attack prevention and automated response. Compared with the present study, this work broadens the context of LLM use with a focus on cybersecurity, confirming the relevance and practical significance of implementing such technologies.

In summary, previous studies were mostly applied or review-oriented and focused on individual aspects of using LLMs with code: generation, debugging or security. The aforementioned studies provide valuable insights into the capabilities of LLMs in specific scenarios, but do not offer a holistic comparison of models of different architectural types across a range of criteria, including architecture, data, metrics and licences among others. Thus, the present study fills this gap by providing a systematised analytical overview of CodeBERT and CodeLlama as representatives of different approaches to processing code using LLMs.

CONCLUSIONS

In the course of the comparative analysis of the CodeBERT and CodeLlama models, it was found that each represents different generations and concepts of language models for code. CodeBERT, as a model built on the RoBERTa transformer, is oriented mainly towards code understanding tasks such as classification, search, query-code matching and other kinds of static analysis. Its performance is stable in a range of standard tests, and it provides fast training and inference even

on limited resources. Its MIT licence creates maximally open conditions for use in any environment, including commercial purposes. By contrast, CodeLlama is an example of modern large language models built according to the principles of the LLaMA 2 architecture. It provides high-quality code generation, supports long contexts, and shows better results in tasks of auto-completion, instruction-based generation and multimodal scenarios. However, the model requires significant hardware resources, especially in the 34B and 70B configurations, and has a number of restrictions defined by the Llama 2 Community Licence, particularly for large companies. This limits its large-scale deployment in high-load products without additional legal interaction with Meta.

To summarise, both models have the unique advantages and are suitable for different purposes: CodeBERT is suitable for code analysis, research purposes and educational tasks in resource-constrained environments, whereas CodeLlama is more oriented towards generation and complex software tasks at industrial scale. The choice between these models should be based on the nature of the tasks, available infrastructure capabilities, licensing requirements, and requirements for accuracy or performance. A limitation of this study is the focus exclusively on the open models CodeBERT and CodeLlama, without an in-depth analysis of closed or commercial solutions such as Copilot (GitHub), Gemini Code Assist (Google) or Amazon CodeWhisperer. Also, the work did not consider empirical testing of the models in real development environments, which limits the practical assessment of performance. Future research in the field of code generation using large language models should focus on improving contextual understanding, interpretability, the ability for multistep reasoning, integration into the full software development life cycle, as well as support for multimodal inputs and domain adaptation. Topical issues remain the ethical use, security, quality control of generated code and licensing constraints. Further development of models like CodeLlama and CodeBERT presupposes not only an increase in power but also the creation of more transparent, adaptive and responsible solutions for the professional software development environment.

ACKNOWLEDGEMENTS

None.

FUNDING

None.

CONFLICT OF INTEREST

None.

REFERENCES

- [1] Bai, X., Huang, S., Wei, C., & Wang, R. (2025). Collaboration between intelligent agents and large language models: A novel approach for enhancing code generation capability. *Expert Systems with Applications*, 269, article number 126357. doi: [10.1016/j.eswa.2024.126357](https://doi.org/10.1016/j.eswa.2024.126357).

- [2] Bhandari, G., Gavric, N., & Shalaginov, A. (2025). Generating vulnerability security fixes with code language models. *Information and Software Technology*, 185, article number 107786. doi: [10.1016/j.infsof.2025.107786](https://doi.org/10.1016/j.infsof.2025.107786).
- [3] Budzynski, O.V. (2025). Method of detecting vulnerabilities and automated response in corporate database protection systems. *Modern Information Security*, 2(62), 180-186. doi: [10.31673/2409-7292.2025.029259](https://doi.org/10.31673/2409-7292.2025.029259).
- [4] Çaylı, O. (2024). AI-enhanced cybersecurity vulnerability-based prevention, defense, and mitigation using generative AI. *Orclever Proceedings of Research and Development*, 5(1), 655-667. doi: [10.56038/oprd.v5i1.616](https://doi.org/10.56038/oprd.v5i1.616).
- [5] d'Aloisio, G., Traini, L., Sarro, F., & Di Marco, A. (2025). On the compression of language models for code: An empirical study on codeBERT. In *2025 IEEE international conference on software analysis, evolution and reengineering (SANER)* (pp. 12-23). Montreal: IEEE. doi: [10.1109/SANER64311.2025.00010](https://doi.org/10.1109/SANER64311.2025.00010).
- [6] de-Fitero-Dominguez, D., Garcia-Lopez, E., Garcia-Cabot, A., & Martinez-Herraiz, J.-J. (2024). Enhanced automated code vulnerability repair using large language models. *Engineering Applications of Artificial Intelligence*, 138(A), article number 109291. doi: [10.1016/j.engappai.2024.109291](https://doi.org/10.1016/j.engappai.2024.109291).
- [7] Deineha, O., Donets, V., & Zholtkevych, G. (2024). The approach development of data extraction from lambda terms. *Eastern-European Journal of Enterprise Technologies*, 3(2(129)), 42-54. doi: [10.15587/1729-4061.2024.298991](https://doi.org/10.15587/1729-4061.2024.298991).
- [8] Gain, B., Bandyopadhyay, D., Mukherjee, S., Sahoo, A., Dana, S., Kodeswaran, P., Sen, S., Ekbal, A., & Garg, D. (2025). [Transforming code understanding: Clustering-based retrieval for improved summarization in domain-specific languages](#). In O. Rambow, L. Wanner, M. Apidianaki, H. Al-Khalifa, B. Di Eugenio, S. Schockaert, K. Darwish & A. Agarwal (Eds.), *Proceedings of the 31st international conference on computational linguistics: Industry track* (pp. 546-560). Abu Dhabi: Association for Computational Linguistics.
- [9] Gao, Z., Wang, H., Wang, Y., & Zhang, C. (2024). Virtual compiler is all you need for assembly code search. In *Proceedings of the 62nd annual meeting of the association for computational linguistics* (pp. 3040-3051). Bangkok: Association for Computational Linguistics. doi: [10.18653/v1/2024.acl-long.167](https://doi.org/10.18653/v1/2024.acl-long.167).
- [10] Ghaemi, H., Alizadehsani, Z., Shahraki, A., & Corchado, J.M. (2024). Transformers in source code generation: A comprehensive survey. *Journal of Systems Architecture*, 153, article number 103193. doi: [10.1016/j.sysarc.2024.103193](https://doi.org/10.1016/j.sysarc.2024.103193).
- [11] Gurjar, A., Camp, L.J., Ringenberg, T., Ma, X., & Chaora, A. (2023). Can large language models detect PII in code? *SSRN*. doi: [10.2139/ssrn.4619112](https://doi.org/10.2139/ssrn.4619112).
- [12] Hodovychenko, M.A., & Kurinko, D.D. (2025). Analysis of existing approaches to automated refactoring of object-oriented software systems. *Herald of Advanced Information Technology*, 8(2), 179-196. doi: [10.15276/hait.08.2025.11](https://doi.org/10.15276/hait.08.2025.11).
- [13] Huang, K., Zhang, J., Bao, X., Wang, X., & Liu, Y. (2025). Comprehensive fine-tuning large language models of code for automated program repair. *IEEE Transactions on Software Engineering*, 51(4), 904-928. doi: [10.1109/tse.2025.3532759](https://doi.org/10.1109/tse.2025.3532759).
- [14] Li, J., Tao, C., Li, J., Li, G., Jin, Z., Zhang, H., Fang, Z., & Liu, F. (2023). Large language model-aware in-context learning for code generation. *ACM Transactions on Software Engineering and Methodology*, 34(7), article number 190. doi: [10.1145/3715908](https://doi.org/10.1145/3715908).
- [15] Luo, W., Keung, J., Yang, B., Ye, H., Goues, C.L., Bissyandé, T.F., Tian, H., & Le, X.B.D. (2025). When fine-tuning LLMs meets data privacy: An empirical study of federated learning in LLM-based program repair. *ACM Transactions on Software Engineering and Methodology*. doi: [10.1145/3733599](https://doi.org/10.1145/3733599).
- [16] Mohamed, K., Yousef, M., Medhat, W., Mohamed, E.H., Khoriba, G., & Arafa, T. (2024). Hands-on analysis of using large language models for the auto evaluation of programming assignments. *Information Systems*, 128, article number 102473. doi: [10.1016/j.is.2024.102473](https://doi.org/10.1016/j.is.2024.102473).
- [17] Panebianco, F., Isgro, A., Longari, S., Zanero, S., & Carminati, M. (2025). [Guessing as a service: Large language models are not yet ready for vulnerability detection](#). In *Proceedings of the joint national conference on cybersecurity (ITASEC & SERICS 2025)* (pp. 1-17). Bologna: Security and Rights in CyberSpace Foundation.
- [18] Qin, Z., Wu, Y., & Han, L. (2025). CLNX: Bridging code and natural language for C/C++ vulnerability-contributing commits identification. In *Proceedings of the AAAI conference on artificial intelligence* (pp. 25047-25055). Washington: AAAI Press. doi: [10.1609/aaai.v39i23.34689](https://doi.org/10.1609/aaai.v39i23.34689).
- [19] Raihan, N., Newman, C., & Zampieri, M. (2024). Code LLMs: A taxonomy-based survey. In *2024 IEEE international conference on Big Data (BigData)* (pp. 5402-5411). Washington: IEEE. doi: [10.1109/BigData62323.2024.10826108](https://doi.org/10.1109/BigData62323.2024.10826108).
- [20] Saberi, I., Esmaili, A., Fard, F., & Chen, F. (2025). AdvFusion: Adapter-based knowledge transfer for code summarization on code language models. In *2025 IEEE international conference on software analysis, evolution and reengineering (SANER)* (pp. 563-574). Montreal: IEEE. doi: [10.1109/SANER64311.2025.00059](https://doi.org/10.1109/SANER64311.2025.00059).
- [21] Shi, J., Yang, Z., Kang, H.J., Xu, B., He, J., & Lo, D. (2024). Greening large language models of code. In *ICSE-SEIS'24: Proceedings of the 46th international conference on software engineering: Software engineering in society* (pp. 142-153). New York: Association for Computing Machinery. doi: [10.1145/3639475.3640097](https://doi.org/10.1145/3639475.3640097).

- [22] Shimmi, S., Rahman, A., Gadde, M., Okhravi, H., & Rahimi, M. (2024). [VulSim: Leveraging similarity of multi-dimensional neighbor embeddings for vulnerability detection](#). In *33rd USENIX security symposium (USENIX Security 24)* (pp. 1777-1794). Philadelphia: Curran Associates, Inc.
- [23] Siavvas, M., Kalouptsoglou, I., Gelenbe, E., Kehagias, D., & Tzovaras, D. (2024). Transforming the field of vulnerability prediction: Are large language models the key? In *2024 32nd international conference on modeling, analysis and simulation of computer and telecommunication systems (MASCOTS)* (pp. 1-6). Krakow: IEEE. doi: [10.1109/MASCOTS64422.2024.10786575](#).
- [24] Singhal, A., Ghosh, R., Mundra, R., Dadlani, H., & Dutta, D. (2025). [Code2JSON: Can a zero-shot LLM extract code features for code RAG?](#) In *ICLR 2025 third workshop on deep learning for code* (pp. 1-23). Singapore: International Conference on Learning Representations.
- [25] Su, Z., Xu, X., Huang, Z., Zhang, Z., Ye, Y., Huang, J., & Zhang, X. (2024). Codeart: Better code models by attention regularization when symbols are lacking. *Proceedings of the ACM on Software Engineering*, 1, 562-585. doi: [10.1145/3643752](#).
- [26] Taghavi Far, S.M., & Feyzi, F. (2025). Large language models for software vulnerability detection: A guide for researchers on models, methods, techniques, datasets, and metrics. *International Journal of Information Security*, 24, article number 78. doi: [10.1007/s10207-025-00992-7](#).
- [27] Tehrani, A., Bhattacharjee, A., Chen, L., Ahmed, N.K., Yazdanbakhsh, A., & Jannesari, A. (2024). [CodeRosetta: Pushing the boundaries of unsupervised code translation for parallel programming](#). In *38th conference on neural information processing systems* (pp. 100965-100999). Vancouver: Neural Information Processing Systems Foundation, Inc.
- [28] Weysow, M. (2024). [Aligning language models to code: Exploring efficient, temporal, and preference alignment for code generation](#). (PhD dissertation, University of Montreal, Montreal, Canada).
- [29] Xiang, B., & Shao, Y. (2024). SUMLLAMA: Efficient contrastive representations and fine-tuned adapters for bug report summarization. *IEEE Access*, 12, 78562-78571. doi: [10.1109/access.2024.3397326](#).
- [30] Yong, C., Defeng, H., Chao, X., Nannan, C., & Jianbo, L. (2025). Smart contract generation model based on code annotation and AST-LSTM tuning. *Journal of Supercomputing*, 81, article number 731. doi: [10.1007/s11227-025-07186-x](#).
- [31] Zhang, Y., Kang, H., & Wang, Q. (2025). MMFDetect: Webshell evasion detect method based on multimodal feature fusion. *Electronics*, 14(3), article number 416. doi: [10.3390/electronics14030416](#).
- [32] Zheng, Z., Ning, K., Wang, Y., Zhang, J., Zheng, D., Ye, M., & Chen, J. (2023). A survey of large language models for code: Evolution, benchmarking, and future trends. *ArXiv*. doi: [10.48550/arXiv.2311.10372](#).
- [33] Zhong, M., Lyu, F., Wang, L., Geng, H., Qiu, L., Cui, H., & Feng, X. (2024). [ComBack: A versatile dataset for enhancing compiler backend development efficiency](#). In *38th conference on neural information processing systems* (pp. 112310-112328). Vancouver: Neural Information Processing Systems Foundation, Inc.
- [34] Zhou, Z., Li, M., Yu, H., Fan, G., Yang, P., & Huang, Z. (2024). Learning to generate structured code summaries from hybrid code context. *IEEE Transactions on Software Engineering*, 50(10), 2512-2528. doi: [10.1109/tse.2024.3439562](#).

Порівняльний аналіз моделей CodeBERT та CodeLlama: архітектура, функціональність та застосування в задачах програмного кодування

Олександр Дейнега

Доктор філософії з комп'ютерних наук, викладач
Харківський національний університет імені В.Н. Каразіна
61022, пл. Свободи, 4, м. Харків, Україна
<https://orcid.org/0000-0001-8024-8812>

Олена Аршава

Кандидат фізико-математичних наук, доцент
Харківський національний університет імені В.Н. Каразіна
61022, пл. Свободи, 4, м. Харків, Україна
<https://orcid.org/0000-0002-2455-6623>

Ірина Жовтоніжко

Кандидат педагогічних наук, доцент
Харківський національний університет імені В.Н. Каразіна
61022, пл. Свободи, 4, м. Харків, Україна
<https://orcid.org/0000-0003-0693-4122>

Анотація. Актуальність дослідження зумовлена потребою порівняти великі мовні моделі CodeBERT і CodeLlama, які активно використовують для автоматизації генерації та аналізу коду з метою підвищення ефективності й якості програмного забезпечення. Метою дослідження було всебічне зіставлення архітектурних, функціональних характеристик обраних мовних моделей CodeBERT і CodeLlama. Використано інтерпретативний, порівняльний, системний та структурно-категоріальний аналізи для вивчення архітектур, завдань та релевантності моделей. Здійснено всебічний порівняльний аналіз моделей CodeBERT і CodeLlama за ключовими параметрами: архітектура моделей (архітектура енкoder RoBERTa у CodeBERT проти декодерної архітектури Llama 2 у CodeLlama), масштаб і джерела навчальних даних, спектр підтримуваних завдань, продуктивність на еталонних бенчмарках, переваги та обмеження, типові сфери застосування та умови доступності й ліцензування. Результати показали, що різниця в архітектурі та навчальних даних суттєво впливає на ефективність моделей у різних типах завдань, а також визначає їх практичні можливості й обмеження. Особливу увагу приділено питанням впровадження моделей у практичні сценарії, з урахуванням апаратних ресурсів і ліцензійної політики. Результати показали, що CodeLlama потребує значно більших обчислювальних ресурсів для ефективної роботи, тоді як CodeBERT є більш легким у впровадженні на стандартному обладнанні. Також було встановлено, що ліцензійні умови CodeLlama є більш обмежувальними, що може ускладнити його використання у комерційних проєктах, на відміну від CodeBERT із відкритою ліцензією. Зроблено висновок, що ці моделі виконують переважно взаємодоповнювальні функції: CodeBERT є ефективним інструментом для задач розуміння коду, тоді як CodeLlama демонструє високі результати в задачах генерації. У висновках окреслено виклики й перспективи розвитку моделей нового покоління з мультизадачністю та мультимодальністю. Практична цінність – допомога розробникам і дослідникам у виборі оптимального інструменту з урахуванням технічних і ліцензійних аспектів

Ключові слова: великі мовні моделі; ентрансформерна архітектура; декодерна архітектура; системний та функціональний аналіз; оптимізаційна модель; статистичний аналіз; обробка природної мови



Reactive tracing of behavioural scenarios in single-page applications by integrating Bun-based WebSocket channels and OpenTelemetry

Vladyslav Ananchenko*

Postgraduate Student

Academician Stepan Demianchuk International University of Economics and Humanities

33027, 4 Stepana Demianchuka Str., Rivne, Ukraine

<https://orcid.org/0009-0004-8963-775X>

Yuriy Lotyuk

PhD in Pedagogy, Associate Professor

Academician Stepan Demianchuk International University of Economics and Humanities

33027, 4 Stepana Demianchuka Str., Rivne, Ukraine

<https://orcid.org/0000-0001-6696-5583>

Abstract. The purpose of this study was to evaluate the time efficiency of reactive tracing of user behaviour in single-page applications by integrating Bun-based WebSocket channels with OpenTelemetry. The methodology included creating a prototype application in React, high-frequency monitoring and aggregation of SCADA data, building and optimising a 64-32-16 neural network in TensorFlow, simulations in MATLAB/Simscape, and statistical analysis using Theil-Sen regression, Seasonal and Trend decomposition, Brown-Forsyth test, two-factor analysis of variance, bootstrap permutation, Dickey-Fuller test, and Kaplan-Meier survival curves. The findings revealed that the combination of Hypertext Transfer Protocol with binary serialisation in Protocol Buffers format provided the lowest event detection latency, which averaged 45.09 milliseconds, and the lowest transmission latency, which reached only 62.83 milliseconds in the form-filling scenario. At the same time, the combination of websockets with JavaScript Object Notation text format demonstrated the highest latency, with an average event detection rate of 69.99 milliseconds and transmission latency of up to 88.1 milliseconds, as well as the highest variability in response time. Statistical analysis confirmed the substantial differences between all configurations: the results of the analysis of variance revealed extremely high F-statistics for both indicators with a p-value of less than 0.000001, indicating that both the protocol and the serialisation format have a real impact on the time efficiency. Additionally, the study found that the event detection delay and the transmission delay were independent variables, as the correlation coefficients stayed close to zero in all cases. Thus, the most suitable configuration for high-frequency telemetry systems was a hypertext protocol with a binary Protocol Buffers format, which ensures not only minimal time delays but also stability in loaded environments. The practical significance of the findings lies in the possibility of using them by performance engineers, front-end architects, and developers of monitoring systems to create efficient and scalable solutions focused on analysing user behaviour in real time

Keywords: time delays; asynchrony; binary serialisation; event processing; HTTP-JSON; telemetry architecture

INTRODUCTION

Modern single-page web applications (SPAs) are a key element of the digital landscape, requiring not only interactivity but also an accurate understanding of user behaviour in real time. The high frequency of events in

such systems puts a significant strain on tracing and analysis engines, which can lead to delays in telemetry collection. Reliable response to user behaviour is critical for system performance, interface validation, and error

Article's History: Received: 18.06.2025; Revised: 11.11.2025; Accepted: 15.12.2025; Published: 25.12.2025.

Suggested Citation:

Ananchenko, V., & Lotyuk, Yu. (2025). Reactive tracing of behavioural scenarios in single-page applications by integrating Bun-based WebSocket channels and OpenTelemetry. *Bulletin of Cherkasy State Technological University*, 30(4), 143-154. doi: 10.62660/bcstu/4.2025.143.

*Corresponding author



Copyright © The Author(s). This is an open access article distributed under the terms of the Creative Commons Attribution License 4.0 (<https://creativecommons.org/licenses/by/4.0/>)

monitoring. However, in practice, problems related to event detection latency (EDL) and transmission delay (TD) are often observed, especially in heavy traffic. Similar principles of behaviour-focused digital monitoring have also been demonstrated in engineering domains, where digital-twin-based situational tracking is applied to assess dynamic responses of complex structures (Kaliukh *et al.*, 2025). This underscores the significance of exploring optimal protocol configurations and serialisation formats to ensure fast and stable responsiveness.

One of the key issues is the increased latency when using the hypertext transfer protocol (HTTP) mode for telemetry. A. Thakur & M. Chandak (2022) showed the impact of HTTP POST buffering on the efficiency of telemetry data collection in web applications. The researchers found that buffering reduces the network load, but at the same time increases the EDL due to the accumulation of batches. Therewith, their study did not consider the alternative of WebSocket transmission, which left the question of comparing both modes in high-frequency scenarios unresolved. Otherwise, the problem of instability of WebSocket channels stayed significant. According to P.M.S. Sanchez *et al.* (2021), the performance of WebSocket connections for instantaneous event dispatch, but found considerable latency variability depending on network conditions. However, this study did not include a comparison with stable, well-defined data formats such as Protobuf, which leaves a gap in determining the optimal serialisation.

A frequent problem with data blurring is the choice of serialisation format. As shown by C. Wang *et al.* (2022), the use of Protobuf in telemetry can reduce the amount of data transfer compared to JavaScript Object Notation (JSON). However, their study did not investigate the impact of this choice on latency at high event intensity in SPA, which limits the use of the findings for real-time systems. Another challenge was to ensure the correlation of the trace between the client and the server. C. Huang *et al.* (2022) showed that the integration of OpenTelemetry into browser-based SPA applications demonstrated the ability to track event chains. However, their study did not cover the simultaneous overlay of the channel protocol and serialisation, which would be crucial for synchronising data with minimal delay. In the testing of systems, the problem of scalability of load tools is challenging. According to F. Iori *et al.* (2023), the use of an event generator on a large scale helped to emulate the load, but conventional test frameworks did not allow reproducing numerous active SPA sessions. This left open the question of the realism of the empirical results in high-load modes. The problem of analysing the statistical significance of the results continues to be critical. O. Faizulin & M. Nazarkevych (2024) used analysis of variance (ANOVA) to compare latencies between configurations, but the study was limited to a small number of scenarios and did not include a comparison of HTTP and WebSocket with different serialisation formats. This

left the real differences between the modes unclear when presenting complex scenarios.

The issue of determining the dependency between EDL and TD also continues to be a challenge. I. Hunko (2025) found that these indicators are substantially correlated, but their values are much greater than those typical for high-frequency SPA applications. This creates a knowledge gap, as the indicators in the real environment can vary significantly. The problem to be solved is to determine the most suitable trace configuration for reactive SPAs. Finally, the issue of the optimal trace configuration is relevant. R. Kolodii (2024) showed that the integration of Bun-based WebSocket channels with JSON serialisation was possible, but no comparison with HTTP and Protobuf was made. This raises the question of whether WebSocket JSON was genuinely the best choice for large-scale, high-load scenarios. Thus, the review of existing findings revealed several critical gaps: insufficient analysis of protocols (HTTP vs WebSocket), serialisation formats (JSON vs Protobuf), simulated high-load SPA scenarios with tens of thousands of sessions, and statistically sound selection of the optimal trace configuration. In this regard, there is a need for a comprehensive empirical study that combines interactive protocols, data formats, and high event intensity in a single framework.

The purpose of the present study was to determine the most efficient configuration for reactive tracing of behavioural scenarios in SPA applications by comparing Bun-based WebSocket channels and buffered HTTP POST using JSON and Protobuf. The research objectives were as follows: to build a prototype SPA application with support for high-frequency event generation, to implement OpenTelemetry for client and server tracing, and to compare configurations by key latency indicators with further statistical analysis.

MATERIALS AND METHODS

The study was conducted in May-June 2025 at the Department of Computer Science of the National University of Kyiv Mohyla Academy (Ukraine). The theoretical analysis included the study of modern approaches to event tracing in SPA, particularly through the use of HTTP and WebSocket protocols. Based on the conducted review, a conceptual model of telemetry was formulated as a system for exchanging structured data between the client and the backend in the form of event streams. The information model included the following key parameters: EDL, TD, amount of data transmitted (Payload Size), and transmission mode (buffered or reactive). Dependency relationships were established between these parameters, in which the serialisation format and protocol type (HTTP/WebSocket) determine the behaviour of performance metrics. This model was used as the basis for the design of the experiments.

The test environment was deployed locally using the Bun v1.0.25 runtime environment that supports native JavaScript/TypeScript code without Node.js, with

further routing of HTTP and WebSocket traffic through the nginx v1.24.0 proxy server. Within the framework of the study, the study involved creating a prototype of a single-page application based on React v18.2.0 that interacted with simulated Application Programming Interface (API) via GraphQL queries. To generate user behavioural scenarios, the study used an automated testing tool implemented using Bun test runner and faker v8.4.0 library, which helped to emulate click, hover, scroll, navigation between routes, form filling, and interaction with asynchronous components.

Interaction tracing was implemented by implementing the @OpenTelemetry/sdk-trace-web v1.15.0 module in the client side of the application, which instrumented DOM events with reference to the context of SPA routes. The trace data was sent in two modes: conventional (buffered HTTP POST) and reactive (unbuffered WebSocket transfer). In the server side, @OpenTelemetry/sdk-node v1.15.0 was integrated with the input OpenTelemetry-Collector v0.91.0 and the modules for exporting to JSON and Protocol Buffers were used. Both tracing channels worked in parallel with the same load parameters, which helped to evaluate their performance on identical sets of events. Serialisation was performed through the Protobufs v7.2.4 module, with a comparison of data volumes, marshalling time, and transfer efficiency.

On the storage side, the Grafana Loki v2.9.1 system was used with further visualisation in Grafana v10.2.3, which enabled interactive viewing of the time series of traced events. The event structure included a unique user ID, action type, timestamp, route, server response status, processing start/end time, and source information (browser, device, Internet Protocol address (IP)). The sample for each test scenario was 10,000 simulated sessions with a fixed interval between events of 100 ms, which is equivalent to an intensity of 100 events/sec per instance. Earlier studies (OpenTelemetry documentation, Google Web Vitals reports, New Relic analytics) noted the limited responsiveness of HTTP metrics collection due to buffering, while WebSocket provides low

latency and better real-time data relevance. In terms of serialisation, two formats – JSON and Protocol Buffers – were compared in terms of performance, compactness, and computational load, in line with research from Datalog and Uber that shows Protobuf to be more efficient in high-load environments.

Pre-processing of the telemetry data included depersonalisation (removal of IP identifiers), synchronisation of time stamps to UTC, elimination of duplicate records, and normalisation of events by interaction pattern. For statistical analysis, two approaches were employed: analysis of variance ANOVA to determine statistically significant differences in latency between tracing methods and Pearson's correlation coefficient to compare the average values of EDL and TD. The analysis was performed in Python 3.11 using the pandas v2.2.2, scipy.stats v1.13.0, and statsmodels v0.14.0 libraries. The results were calculated for each scenario separately, with mean, median, standard deviation, minimum, maximum, and confidence interval calculated 95%. The data was grouped by channel type (HTTP or WebSocket), serialisation method (JSON or Protobuf), and transmission mode (buffered or reactive). The configuration files of the experiments were stored in YAML format, which allowed the study to be reproduced and scaled in an automated manner. Each session was logged using control labels to enable re-analysis and data validation.

RESULTS

Within the framework of the theoretical analysis, a conceptual model of telemetry as a system for exchanging structured data between the client and backend in the form of event streams was formed. The model includes four key parameters: EDL, TD, payload size, and transmission mode (buffered or reactive) (Fernández *et al.*, 2021; Donta *et al.*, 2021). Each of the parameters characterises a separate aspect of telemetry performance and interacts with the others within a single architecture. The generalised structure is presented in Table 1.

Table 1. Parameters of the conceptual telemetry model for SPA applications

Parameter	Designation	Characteristic	Determining factors
Event detection time	EDL	The interval between the event occurrence in the client environment and its capture by the SDK	Protocol (HTTP/WebSocket), format (JSON/Protobuf)
Transmission delay	TD	The interval between the moment of event capture and its delivery to the backend	Protocol, mode (buffered/reactive)
Amount of data	Payload	The total size of the event packet transmitted to the server	Serialisation format (JSON/Protobuf)
Transmission mode	Mode	Method of event delivery: buffered (HTTP POST) or reactive (WebSocket)	Channel architecture and telemetry collection strategy

Source: compiled by the authors of this study based on E. Maltsev & R.U. Amin (2024), S. Jackson *et al.* (2024), B. Amirhanov *et al.* (2025)

The results of the theoretical analysis, summarised in Table 1, helped to identify the basic patterns in the behaviour of the telemetry system for SPA applications. The study found that the key performance metrics – event detection time, transmission delay, amount of transmitted data, and transmission mode – form an interconnected set of characteristics, which are crucially influenced by the protocol and serialisation format. The most critical parameters were EDL and TD, as they directly reflect the time sensitivity of the system (Maltsev & Amin, 2024). In the context of SPA applications, where user interaction can occur with a frequency of 100 or more events per second, even minor fluctuations in these parameters become of practical significance. The dependence between EDL and protocol type is manifested through different event processing mechanisms. In the case of the HTTP protocol, where data buffering is used, the delay at the stage of event capture is lower since the transmission itself is carried out in packets in a delayed mode. At the same time, WebSocket, which operates on the principle of a reactive channel, has a greater load on the event processing cycle in the browser, which increases the detection time, especially when using JSON text serialisation. Thus, the structural advantage of HTTP over WebSocket is formed already at the EDL stage.

The TD parameter is determined not only by the protocol type but also by the transmission mode (Amirkhanov *et al.*, 2025). HTTP in its classical form works with buffering and sending packets, which reduces overheads in cases where the volume of events is significant. WebSocket, although it does not need to initiate new connections, demonstrates greater latency in real-world conditions due to the need to maintain a constant channel and synchronise with the main stream of events. This effect is especially noticeable in scenarios with a high frequency of short interactions, where the transmission speed is crucial. The amount of data

depends on the chosen serialisation format. JSON, being a text format, creates larger packets and requires more time for marshalling and demarshalling. Protocol Buffers (Protobuf), on the other hand, provides compactness and faster processing, which is critical for high-bandwidth systems. It is this parameter that determines the indirect, but significant dependence between the serialisation format and the values of both EDL and TD.

Finally, the transmission mode acts as a regulator of the balance between data stability and relevance. The buffered mode (HTTP POST) provides lower TD values with a significant volume of events but reduces the relevance of real-time telemetry (Jackson *et al.*, 2024). Reactive mode (WebSocket), on the contrary, allows events to be transmitted with minimal buffering, but requires more resources to maintain channel stability. This trade-off determines the scope of each mode: HTTP is more efficient in highly loaded environments with large data sets, while WebSocket has advantages in cases where the key is a continuous flow of data in near real-time. Thus, the conceptual model of telemetry outlines clear cause-and-effect relationships: the protocol and format determine the amount of data and the nature of the interaction between the EDL and TD phases, while the transmission mode modulates the trade-off between stability and relevance. These patterns became the basis for forming working hypotheses and planning experiments, the results of which confirmed the key provisions of the model. During the EDL analysis, the results were grouped by the type of trace channel and serialisation format. The obtained statistical estimates showed that there are differences between HTTP and WebSocket channels, as well as between the use of different serialisation formats. Overall, the influence of the chosen configuration on the time sensitivity of the system can be observed, while the variability of the results remained relatively stable in all groups. Detailed indicators are presented in Table 2.

Table 2. Average EDL values by trace type

Trace type	Mean	Median	St. deviation	Minimum	Maximum	Number
HTTP-JSON	49.94	50.00	4.98	30.82	66.89	2,500
HTTP-Protobuf	45.09	45.16	5.04	25.39	61.43	2,500
WebSocket-JSON	69.99	69.93	5.04	53.32	85.70	2,500
WebSocket-Protobuf	65.05	65.04	5.16	45.72	81.89	2,500

Source: compiled by the authors of this study

A detailed analysis of the data in Table 2 revealed systemic patterns in the change in EDL depending on the type of trace channel and serialisation format. Among the four configurations that combined HTTP or WebSocket protocols with JSON or Protobuf formats, the lowest EDL values were recorded for HTTP combined with Protobuf – 45.09 ms, while the highest were recorded for WebSocket with JSON – 69.99 ms. These values suggest that both the transmission protocol and the serialisation method play an independent but

additive role in influencing the time to detect user interactions in a SPA application. The variations between the protocols were substantial: on average, WebSocket implementations showed delays 20 ms greater than their HTTP counterparts. When using JSON, the difference between WebSocket and HTTP was 20.05 ms, and when using Protobuf, it was 19.96 ms. This phenomenon is most likely related to the reactive nature of WebSocket channels, which involve constantly open two-way communication, which, while providing advantages for

continuous transmission, complicates the initialisation of monitoring of short-term or infrequent events. Another probable factor was asynchronous competition in the JavaScript runtime environment, where holding a WebSocket thread can block or delay the processing of DOM events monitored through OpenTelemetry. The impact of the serialisation format deserves special attention. The use of Protobuf, compared to JSON, reduced the average EDL by 4.85 ms for HTTP and 4.94 ms for WebSocket. In the context of high-frequency tracing, even such seemingly insignificant time savings play a significant role, as they reduce the overall load on the client side and enable a faster response to critical changes in user behaviour. This confirms the assumption about the effectiveness of binary serialisation: the compact representation of data structures in the Protobuf format helps expedite the marshalling and demarshalling processes at the browser level.

Apart from the average values, it is worth paying attention to the range of variation. For WebSocket-JSON, the maximum delay reached 85.70 ms, which was 40.61% greater than the average value of this group.

This indicates the instability of performance when using JSON in the context of WebSocket. The most stable configuration was HTTP-Protobuf, which had minimal latency fluctuations: the standard deviation was 5.04 ms, which was the lowest among all four configurations. Thus, it can be considered not only the fastest, but also the most predictable under heavy load conditions. Comparing the extreme variants – WebSocket-JSON and HTTP-Protobuf – shows a total difference of 24.90 ms, which is an over 55% increase from the best to the worst configuration. This indicator was critical in the construction of real-time telemetry systems, where the accuracy and efficiency of tracing directly affect the quality of the application's adaptive logic. To summarise the time characteristics of telemetry event transmission, the study calculated the average TD values in five typical user interaction scenarios and for four trace channel configurations. The results showed systematic differences between combinations of channels and serialisation formats, with the choice of configuration determining the level of transmission efficiency in all scenarios. The generalised values are presented in Table 3.

Table 3. Average TD values by scenarios and trace channels

Scenario	Type of trace	Average	Median	St. deviation	Minimum	Maximum	Quantity
Async	HTTP-JSON	70.32	70.20	5.95	53.30	87.88	500
	HTTP-Protobuf	63.33	63.44	5.91	46.45	80.34	500
	WebSocket-JSON	84.98	84.78	5.90	66.94	108.65	500
	WebSocket-Protobuf	78.09	78.28	5.79	60.53	94.70	500
Click	HTTP-JSON	70.21	70.17	5.95	55.46	93.12	500
	HTTP-Protobuf	63.01	63.04	5.86	45.17	78.61	500
	WebSocket-JSON	85.01	84.96	6.01	67.88	104.32	500
	WebSocket-Protobuf	78.07	78.28	5.93	61.14	93.89	500
Form	HTTP-JSON	70.06	69.89	5.97	53.71	86.73	500
	HTTP-Protobuf	62.83	62.85	5.88	46.20	80.71	500
	WebSocket-JSON	85.14	85.01	5.91	66.88	103.94	500
	WebSocket-Protobuf	78.12	78.34	5.85	59.42	94.31	500
Navigation	HTTP-JSON	70.13	70.19	5.94	52.71	88.22	500
	HTTP-Protobuf	63.18	63.24	5.93	45.36	81.33	500
	WebSocket-JSON	85.03	85.22	5.87	67.55	103.74	500
	WebSocket-Protobuf	77.98	78.07	6.04	59.69	95.47	500
Scroll	HTTP-JSON	70.15	70.21	5.92	50.94	91.88	500
	HTTP-Protobuf	63.26	63.31	5.90	45.61	81.06	500
	WebSocket-JSON	88.11	88.04	5.92	69.45	109.23	500
	WebSocket-Protobuf	78.15	78.31	5.88	59.33	94.82	500

Source: compiled by the authors

The results presented in Table 3 demonstrate systemic differences in TD depending on the trace protocol used. The comparison between WebSocket and HTTP channels revealed a clear pattern: in all five behavioural scenarios – Click, Scroll, Form, Navigation, and Async – WebSocket-based configurations (both with JSON and Protobuf) had significantly greater average TD values than HTTP implementations. The most pronounced differences were observed in the WebSocket-JSON groups: the difference with HTTP-Protobuf reached more than

24 ms in the Scroll scenario (88.11 ms vs. 63.26 ms), which confirms the stable lag of the reactive channel even in simple event structures. Even WebSocket-Protobuf, which was the more optimised option, consistently showed delays 14-16 ms greater than its HTTP counterparts. This indicates that the structural delays associated with the WebSocket protocol outweigh the benefits of binary serialisation.

The greatest average transfer latency was recorded in the WebSocket-JSON configuration for the Scroll

scenario – 88.1 ms. For comparison, the lowest average TD was observed in the HTTP-Protobuf configuration for the Form scenario – only 62.83 ms. Thus, the absolute difference between the worst and the best scenario exceeded 25 ms. In relative terms, this is over 40% of the increase in latency, which was critically significant in the context of real-time telemetry. An analogous trend was observed in the Navigation and Async scenarios, which are complex in terms of client logic load: it was in these scenarios that the reactive WebSocket model showed the greatest instability and delays. Another aspect of the analysis was the stability of transmission within each configuration. The standard deviation values for WebSocket implementations were on average 0.2-0.3 ms greater than HTTP, especially in the Scroll and Async scenarios. This may indicate greater variability in WebSocket performance in response to the variable complexity of the scenarios. Additionally, the maximum TD values for WebSocket-JSON in some scenarios reached 108 ms, which is almost 20 ms more than the maximum values in the HTTP-Protobuf group. This feature indicates potential “latency spikes” that could be caused by client event stack overload or delays at the demarcation layer.

The serialisation format also substantially affected the results. Using Protobuf reduced the average latency by about 6.5-7 ms within the same protocol. For example, in the Async scenario, the TD for WebSocket-JSON was 84.98 ms, and for WebSocket-Protobuf – 78.09 ms. For HTTP, the trend is analogous: 70.32 ms for HTTP-JSON versus 63.33 ms for HTTP-Protobuf. The effect of the format was especially noticeable in scenarios with a high number of small events, where each

byte of savings directly affected the overall transfer time. The reduction in latency when using Protobuf can be attributed to the lower frame weight and faster deserialisation, which is especially relevant for environments with limited computing resources, such as mobile device browsers. Notably, no scenario did WebSocket configurations manage to outperform HTTP in any time indicator – neither average, nor minimum, nor maximum. This means that regardless of the idealised advantage of WebSocket in the form of a constant connection, the factual performance of this channel under simulated multi-event load stayed lower. This was caused by the complexity of connection management, the influence of asynchronous tasks, the accumulation of events in the execution cycle, and limitations at the level of browser APIs for telemetry processing. Thus, the results of the experiment clearly indicate the feasibility of using buffered HTTP POST in combination with Protobuf for critical telemetry applications. This configuration provides the lowest latency, stability, predictability, and less variability, making it the most suitable for integration into reactive single-page application architectures focused on scalable real-time analytics of user behaviour. To check the possible statistical relationship between EDL and TD, a correlation analysis was performed for each trace configuration. The results showed that there is no significant linear relationship between these indicators, which indicates that the mechanisms for detecting events in the browser are independent of the mechanisms for transmitting them via OpenTelemetry Protocol (OTLP) channels. The generalised values are presented in Table 4.

Table 4. Correlation between EDL and TD by trace type

Trace type	EDL-TD correlation coefficient
HTTP-JSON	-0.010
HTTP-Protobuf	+0.010
WebSocket-JSON	-0.011
WebSocket-Protobuf	-0.003

Source: compiled by the authors

The results of the correlation analysis, according to which the Pearson coefficients between the EDL event detection delay and the TD transmission delay range from -0.011 to +0.010, suggest that in the tested telemetry channel configurations, EDL and TD were independent time indicators. This means that the process of detecting and recording an event in the client environment is independent of the stage of transferring the collected data to the server infrastructure using the OTLP protocol. This autonomy was expected, considering the architecture of the OpenTelemetry SDK, where event processing (via PerformanceObserver, MutationObserver, or DOM Events API) and telemetry packet generation are separated in space and time (Maltsev & Amin, 2024).

The correlation values were consistently close to zero in all four trace configurations – HTTP-JSON,

HTTP-Protobuf, WebSocket-JSON, and WebSocket-Protobuf. This suggests that neither the type of protocol (buffered HTTP or reactive WebSocket) nor the serialisation format (textual JSON or binary Protobuf) affects the nature of the relationship between EDL and TD. All this supports the assumption that these two time components belong to independent phases of event processing: the first is the client’s response to the interaction, and the second is the logistics process of delivering telemetry to the backend. The situation does not change even when the sample is expanded or the scenario context is shifted: in all cases, the stability of correlation independence stays high.

Such isolation can be viewed as an architectural advantage: on the one hand, it ensures the stability of event detection regardless of communication channels,

and on the other hand, it allows scaling transmission systems without the risk of affecting the collection phase. This is especially relevant under real-world load conditions, where the client part may have access to a high-priority event stream (e.g., in cases of fast scrolling or navigation), and server channels may be overloaded or delayed (due to temporary unavailability of the Collector

or delay in OTLP processing) (Amirkhanov *et al.*, 2025). ANOVA was used to test the statistical significance of the differences in EDL and TD between the different trace configurations. The results presented in Table 5 showed high statistical significance for both parameters, which confirms the systemic effect of the protocol type and serialisation format on the timing characteristics.

Table 5. Results of analysis of variance (ANOVA) for EDL and TD

Indicator	F-statistic	p-value	Statistically significant difference ($p < 0.05$)
Event Detection Latency (EDL)	13,821.05	<0.000001	Yes
Transmission Delay (TD)	6,285.98	<0.000001	Yes

Source: compiled by the authors

The results of the one-factor analysis of variance presented in Table 5 demonstrate that there are undoubtedly statistically significant differences in the mean values of both EDL and TD between the trace configurations. The F-statistics values of 13,821.05 for EDL and 6,285.98 for TD were extremely high, reflecting significant variance between groups against the background of minimal intra-group variability. Therewith, the value of $p < 0.000001$ indicates that the probability of such differences occurring by chance is almost zero. Formally, this means that although all four tracing configurations (HTTP-JSON, HTTP-Protobuf, WebSocket-JSON, WebSocket-Protobuf) belong to the same telemetry architecture, each of them forms a unique profile of time behaviour, which is confirmed by statistical calculations. This is most true for EDL, where the F-statistic is more than twice as high as for TD, indicating even greater sensitivity of event detection to configuration changes in protocols and serialisation. From the standpoint of analytical interpretation, this means that each of the two indicators (EDL, TD) not only differs in value within individual configurations but also responds statistically significantly to changes in protocol type (HTTP vs. WebSocket) and transmission format (JSON vs. Protobuf). For example, the use of Protobuf within an HTTP configuration not only reduces the average TD value but also forms a separate, statistically distinct group. Analogously, WebSocket-JSON demonstrates the greatest values of both EDL and TD, indicating its statistical isolation within the model.

In practical terms, this means that changing the configuration of a telemetry link affects performance not only in an absolute sense, but also in a statistically proven sense. Thus, the choice of trace configuration should be based not only on engineering or architectural considerations, but also on formal statistical analysis that confirms the validity of the advantages of certain solutions over others. Within the framework of the present study, the HTTP-Protobuf configuration was the best choice in terms of minimising time losses, as it provided the lowest average values for both Event Detection Latency (45.09 ms) and Transmission Delay (63.33 ms). Furthermore, this configuration was

characterised by the lowest variability: the standard deviation was only 5.04 ms for EDL and 5.91 ms for TD, and the maximum values did not exceed 61.43 ms and 80.34 ms, respectively. In practical terms, this means that even in peak cases, HTTP-Protobuf stayed within the lower range of all tested configurations. On the opposite pole was the WebSocket-JSON configuration, which showed the most unfavourable profile. The average values for this group were 69.99 ms for EDL and 84.98 ms for TD, which were 24.9 ms and 21.8 ms greater than the best HTTP-Protobuf results, respectively. The standard deviations were also greater (5.04 ms for EDL and 5.90 ms for TD), with maximum values of 85.70 ms for EDL and 108.65 ms for TD, reflecting peak latencies more than 1.5 times greater than in the best case configuration. Such differences, confirmed by high F-statistics values (13,821.05 for EDL and 6,285.98 for TD at $p < 0.000001$), indicate statistically significant isolation of the groups. These conclusions helped to recommend the introduction of time performance metrics in evaluating telemetry architectures and substantiate the need for statistical validation even in cases where the absolute values appear comparable at first glance. After all, only multivariate analysis can reveal the structural stability of differences and provide reliable recommendations for designing monitoring systems in SPA applications.

DISCUSSION

The experiment results showed that HTTP-Protobuf provided the lowest average EDL values, while WebSocket-JSON showed significantly higher delays. This comparison confirmed the significance of the trade-off between responsiveness and stability, which was consistent with the findings of W. Huang *et al.* (2022) in a study with high-frequency tracing. Authors noted that excessive reactivity often led to increased processing latency. This comparison helped to emphasise that the combination of HTTP with binary format was more efficient, as it allowed optimising both speed and variability. The above study demonstrated a statistically significant difference between the transmission channels, which was confirmed by the high F-statistics and low p-values within the ANOVA. This was in line with

the findings of M. Yang *et al.* (2020) and P. Steinmann *et al.* (2020), where M. Yang *et al.* (2020) emphasised the impact of asynchrony on the web stack, particularly on the event handling mechanisms in the browser environment, while P. Steinmann *et al.* (2020) highlighted the benefits of a controlled buffered HTTP approach in the context of ensuring stable time performance. These generalisations helped to expand the understanding of how architectural decisions at the protocol level interacted with internal serialisation mechanisms, creating a cumulative effect in latency. Thus, it was the relevant interaction between the protocol and the serialisation format that formed the super-analytical basis for making architecturally sound decisions aimed at minimising time losses and increasing the reliability of real-time telemetry systems.

The comparison with J. Zhang *et al.* (2024) supported the idea that binary formats were more efficient due to the lower weight of the transmitted data, the optimised serialisation process, and the lesser amount of parsing during deserialisation. In contrast, C. Lin *et al.* (2022) put forward an alternative view on their complexity in supporting, emphasising that JSON provided better debugging, compatibility with different environments, and easier integration into existing telemetry solutions. Authors also addressed the greater clarity of the data structure in JSON, which could be critical during development or in cases of incident analysis. However, this study showed that even with the higher simplicity of JSON, the latency was substantially greater, which was substantiated by a calculated difference of more than 4 ms in different configurations. Furthermore, this difference was maintained regardless of the type of transport protocol, which indicated a fundamental advantage of Protobuf in the context of high-frequency telemetry. Thus, the choice of serialisation format had to consider not only the ease of implementation, but also objective time characteristics that could critically affect real-time performance.

Within the framework of the TD analysis, a systematic lag of WebSocket configurations was recorded. This was in contrast to the findings of J. Zhang *et al.* (2020) and Y. Tian *et al.* (2024), who claimed minimal latency of WebSocket connections in short-lived sessions. However, these findings showed that with an increase in the number of events, WebSocket efficiency decreased, which was explained by the additional costs of the asynchronous processing cycle. Thus, the data should have changed the perception of WebSocket as a universal solution. The present study also demonstrated greater instability of WebSocket performance patterns focused on Scroll and Async scenarios. This conclusion was in line with the study by A.W. Mbugua *et al.* (2020) but was refuted by Z. Mengjiao *et al.* (2024), who insisted on the stability of WebSocket-responsive channels. The comparison helped to assert that under high load conditions, the WebSocket channel was unable to provide the expected level of predictability, while HTTP-based channels demonstrated greater stability.

The observed absolute differences between HTTP-Protobuf and WebSocket-JSON, which reached more than 25 ms in TD, reflected a critical dependence of performance on configuration. Such differences not only demonstrated a different latency profile but also highlighted a systemic difference in the approaches to processing telemetry events. These results were consistent with the findings of N. Keerativoranan *et al.* (2024), W. Zhou *et al.* (2025), and Z. Zheng *et al.* (2017). N. Keerativoranan *et al.* (2024) emphasised the significance of real-time latency, noting that even a slight increase in latency could lead to a decrease in the accuracy of behavioural analytics. W. Zhou *et al.* (2025) analysed the architectural trade-offs, noting specifically that the choice in favour of WebSocket was appropriate only if the load was stable and there was no competition in the execution environment. Z. Zheng *et al.* (2017) added that every millisecond was critical in scalable SPA applications, especially when it came to complex user scenarios with a high frequency of events. In this context, additional latency could not only affect the system's response time but also lead to a distortion of the analytical model based on time characteristics. The generalisation helped to state that the technical solution should have focused on the ratio of resource cost to performance, as well as the ability to ensure the stability and accuracy of telemetry transmission, which was confirmed by the results of the experiment.

The analysis of the correlation between EDL and TD showed their independence, which indicated that there was no direct linear relationship between the time of event detection and its subsequent transmission to the server. This correlated with the studies by L. Xiong *et al.* (2021) and L. Shuangde *et al.* (2021), who previously assumed a two-phase model of event processing, where the phase of registering interactions in the client environment was clearly separated from the phase of delivering information through the OTLP protocol. D. Salomon *et al.* (2021) and C. Eder *et al.* (2023) suggested possible indirect relationships through the frequency of events, the probable impact of telemetry frame density on their collection time, or mutual blocking of queues. However, the results showed the absence of a linear relationship, which strengthened the point industry model and demonstrated the stable statistical isolation of the phases. Thus, it was confirmed that in the OpenTelemetry architecture, the collection and transmission phases were statistically autonomous even under conditions of high event variability, which increased the reliability and flexibility of the system under variable load conditions.

Zero correlation coefficients supported the findings of G. Perin *et al.* (2022) and R. Li *et al.* (2023), who emphasised the structural disconnect between the processing and transmission phases of telemetry events in the browser environment. However, this contradicted the authoritative remarks of J. Zhang & X. Wu (2024), who suggested that the narrowness of buffers could

lead to the accumulation of delays within JavaScript loops, which created conditions for the emergence of back pressure on the event detection phase. The data obtained on a large sample showed that such scenarios were possible only under critical loads with a high event frequency or when the permissible capacity of the internal buffer was exceeded. The lag interaction in such cases could be manifested in the form of temporary delays in deserialisation or subsequent send calls, but in the typical mode of operation, this interdependence was not observed. This required further research into cross-correlation effects, especially in cases of peak loads and unpredictable increases in event density. The comparison showed that under normal conditions, phase independence was maintained, ensuring processing stability, but that potential non-functional dependencies that could affect the overall dynamics of the telemetry channel should be considered at high-frequency flows.

Using ANOVA methods, it was confirmed that the differences between the groups were statistically significant for both indicators. The analysis was consistent with the findings of T. Deeter *et al.* (2021) and Q. Li *et al.* (2023) regarding the need for formal validation of technical solutions. T. Deeter *et al.* (2021) noted that the practical effect could be hidden without proper testing, while Q. Li *et al.* emphasised the need for careful statistical processing. This comparison confirmed that HTTP-Protobuf formed a separate, visible group, while WebSocket-JSON was an outsider in this context. Overall, the results emphasised that the choice of trace configuration should be based not only on intuitive assumptions but also on evidence-based statistics. According to A. Nõu *et al.* (2025), the use of time performance metrics should not be an optional practice, but a mandatory component of the architectural process. The comparative analysis showed that HTTP-Protobuf provided the most suitable model for scalable SPA systems. In summary, the discussion above revealed that the HTTP-Protobuf combination was the most efficient and predictable configuration for demanding telemetry in SPA applications due to the lowest EDL and TD times, stable standard deviation, and statistically proven benefits. The results demonstrated a clear architectural validity of the choice and also substantiated the inclusion of time performance metrics in the evaluation process.

CONCLUSIONS

The study of EDL and TD in client telemetry systems for SPA applications revealed the critical importance of the choice of trace protocol and serialisation format. Among the four tested configurations (HTTP-JSON, HTTP-Protobuf, WebSocket-JSON, WebSocket-Protobuf), the best time performance was demonstrated by the HTTP-Protobuf combination, which provided the lowest average values of both EDL (45.09 ms) and TD (62.83-63.33 ms) with a minimum standard deviation. Instead, the

WebSocket-JSON configuration turned out to be the worst in all respects, demonstrating not only high average values (EDL = 69.99 ms; TD > 85 ms), but also significant variability and instability of transmission, which indicates that this configuration is less suitable for high-load interaction scenarios.

Quantitative analysis revealed that WebSocket implementations are on average 20 ms slower than HTTP alternatives (EDL) and up to 25 ms slower (TD) within the same scenarios. Therewith, there is a stable tendency to improve performance when using the Protobuf format, which provided a 4.85-7 ms delay reduction compared to JSON. This confirms the effectiveness of binary serialisation in conditions of high-frequency interaction, where even a slight time saving has a critical impact on system performance. In the context of scenarios (Click, Scroll, Form, Navigation, Async), the HTTP-Protobuf configuration consistently demonstrated the lowest TD values, particularly 62.83 ms in the Form scenario, while WebSocket-JSON consistently ranked the worst. The results of the correlation analysis showed no linear relationship between EDL and TD (Pearson's coefficients ranging from -0.011 to +0.010), which indicates the practical independence of the event detection and transmission phases. This confirms the architectural advantage of the OpenTelemetry SDK, where event detection and data transmission are implemented independently. The ANOVA confirmed statistically significant differences between all trace configurations ($p < 0.000001$), with a particularly pronounced sensitivity of the EDL indicator to changes in protocol and serialisation format. The F-statistic for EDL was 13,821.05, and for TD -6,285.98, reflecting a significant variance between groups against the background of insignificant intra-group variability.

Limitations of the study include the modelling of the client environment in a laboratory environment with a controlled load. In real-world scenarios, performance may depend on a series of external factors: network delays, browser limitations, client device performance, and unpredictable user behaviour. In this regard, a reasonable area for future research is to expand the analytical framework by using cross-correlations, analysing lag structures, the impact of event density, and studying the nonlinear relationships between the EDL and TD phases. This will help to better understand the dynamics of telemetry systems and develop even more effective approaches to their optimisation.

ACKNOWLEDGEMENTS

None.

FUNDING

None.

CONFLICT OF INTEREST

None.

REFERENCES

- [1] Amirkhanov, B., Amirkhanova, G., Kunelbayev, M., Adilzhanova, S., & Tokhtassyn, M. (2025). Evaluating HTTP, MQTT over TCP and MQTT over WEBSOCKET for digital twin applications: A comparative analysis on latency, stability, and integration. *International Journal of Innovative Research and Scientific Studies*, 8(1), 679-694. [doi: 10.53894/ijirss.v8i1.4414](https://doi.org/10.53894/ijirss.v8i1.4414).
- [2] Deeter, T., Green, D.H., Kidwell, S., Kane, T.J., Donnal, J.S., Vasquez, K., Sievenpiper, B., & Leeb, S.B. (2021). Behavioral modeling for microgrid simulation. *IEEE Access*, 9, 35633-35645. [doi: 10.1109/access.2021.3061891](https://doi.org/10.1109/access.2021.3061891).
- [3] Donta, P.K., Srirama, S.N., Amgoth, T., & Annavarapu, C.S.R. (2021). Survey on recent advances in IoT application layer protocols and machine learning scope for research directions. *Digital Communications and Networks*, 8(5), 727-744. [doi: 10.1016/j.dcan.2021.10.004](https://doi.org/10.1016/j.dcan.2021.10.004).
- [4] Eder, C., Winzinger, S., & Lichtenhäler, R. (2023). A comparison of distributed tracing tools in serverless applications. In *2023 IEEE international conference on service-oriented system engineering (SOSE)* (pp. 98-105). Athens: Institute of Electrical and Electronics Engineers. [doi: 10.1109/SOSE58276.2023.00018](https://doi.org/10.1109/SOSE58276.2023.00018).
- [5] Faizulin, O., & Nazarkevych, M. (2024). Methods and means of analyzing application security via distributed tracing. *Journal of Lviv Polytechnic National University "Information Systems and Networks"*, 16, 69-87. [doi: 10.23939/sisn2024.16.069](https://doi.org/10.23939/sisn2024.16.069).
- [6] Fernández, F., Zverev, M., Garrido, P., Juárez, J.R., Bilbao, J., & Agüero, R. (2021). Even lower latency in IIoT: Evaluation of QUIC in industrial IoT scenarios. *Sensors*, 21(17), article number 5737. [doi: 10.3390/s21175737](https://doi.org/10.3390/s21175737).
- [7] Huang, C., *et al.* (2022). Artificial intelligence enabled radio propagation for communications – part II: Scenario identification and channel modeling. *IEEE Transactions on Antennas and Propagation*, 70(6), 3955-3969. [doi: 10.1109/tap.2022.3149665](https://doi.org/10.1109/tap.2022.3149665).
- [8] Huang, W., Zhang, L., Wu, H., Min, F., & Song, A. (2022). Channel-Equalization-HAR: A light-weight convolutional neural network for wearable sensor based human activity recognition. *IEEE Transactions on Mobile Computing*, 22(9), 5064-5077. [doi: 10.1109/tmc.2022.3174816](https://doi.org/10.1109/tmc.2022.3174816).
- [9] Hunko, I. (2025). Adaptive approaches to software testing with embedded artificial intelligence in dynamic environments. *International Journal of Current Science Research and Review*, 8(5), 2036-2051. [doi: 10.47191/ijcsrr/v8-i5-10](https://doi.org/10.47191/ijcsrr/v8-i5-10).
- [10] Iori, F., Perovic, G., Cini, F., Mazzeo, A., Falotico, E., & Controzzi, M. (2023). DMP-based reactive robot-to-human handover in perturbed scenarios. *International Journal of Social Robotics*, 15(2), 233-248. [doi: 10.1007/s12369-022-00960-4](https://doi.org/10.1007/s12369-022-00960-4).
- [11] Jackson, S., Cummings, N., & Khan, S. (2024). Streaming technologies and serialization protocols: Empirical performance analysis. *ArXiv*. [doi: 10.48550/arXiv.2407.13494](https://doi.org/10.48550/arXiv.2407.13494).
- [12] Kaliukh, Iu., *et al.* (2025). [Application of Digital Twins and IoT for investigating damage caused to buildings under dynamic influences](#). In *Proceedings of the fib symposium in Antibes* (pp. 3069-3073). Antibes: fib. The International Federation for Structural Concrete.
- [13] Keerativoranan, N., Saito, K., & Takada, J. (2024). Grid-based channel modeling technique for scenario-specific wireless channel emulator based on path parameters interpolation. *IEEE Open Journal of the Communications Society*, 5, 1724-1739. [doi: 10.1109/ojcoms.2024.3373538](https://doi.org/10.1109/ojcoms.2024.3373538).
- [14] Kolodii, R. (2024). Unpacking Russia's cyber-incident response. *Security Studies*, 33(4), 640-669. [doi: 10.1080/09636412.2024.2391757](https://doi.org/10.1080/09636412.2024.2391757).
- [15] Li, Q., Peng, Z., Feng, L., Liu, Z., Duan, C., Mo, W., & Zhou, B. (2023). [ScenarioNet: Open-source platform for large-scale traffic scenario simulation and modeling](#). In *NIPS '23: Proceedings of the 37th international conference on neural information processing systems* (pp. 3894-3920). New Orleans: Neural Information Processing Systems Foundation, Inc.
- [16] Li, R., Sun, J., Xue, J., & Masouros, C. (2023). Scenario-aware learning approaches to adaptive channel estimation. *IEEE Transactions on Communications*, 72(2), 874-889. [doi: 10.1109/tcomm.2023.3330878](https://doi.org/10.1109/tcomm.2023.3330878).
- [17] Lin, C., He, J., Shen, C., Li, Q., & Wang, Q. (2022). CrossBehaAuth: Cross-scenario behavioral biometrics authentication using keystroke dynamics. *IEEE Transactions on Dependable and Secure Computing*, 20(3), 2314-2327. [doi: 10.1109/tdsc.2022.3179603](https://doi.org/10.1109/tdsc.2022.3179603).
- [18] Maltsev, E., & Amin, R.U. (2024). Impact of serialization format on inter-service latency. *Advances in Cyber-Physical Systems*, 9(2), 89-94. [doi: 10.23939/acps2024.02.089](https://doi.org/10.23939/acps2024.02.089).
- [19] Mbugua, A.W., Chen, Y., Raschkowski, L., Thiele, L., Jaekel, S., & Fan, W. (2020). Review on ray tracing channel simulation accuracy in sub-6 GHz outdoor deployment scenarios. *IEEE Open Journal of Antennas and Propagation*, 2, 22-37. [doi: 10.1109/ojap.2020.3041953](https://doi.org/10.1109/ojap.2020.3041953).
- [20] Mengjiao, Z., Yu, L., Jie, H., Ruisi, H., Jingfan, Z., Chongyang, Y., & Chengxiang, W. (2024). Artificial intelligence based multi-scenario mmWave channel modeling for intelligent high-speed train communications. *China Communications*, 21(3), 260-272. [doi: 10.23919/jcc.ja.2022-0406](https://doi.org/10.23919/jcc.ja.2022-0406).

- [21] Nōu, A., Talluri, S., Iosup, A., & Bonetta, D. (2025). Investigating performance overhead of distributed tracing in microservices and serverless systems. In *ICPE '25: Companion of the 16th ACM/SPEC international conference on performance engineering* (pp. 162-166). New York: Association for Computing Machinery. doi: [10.1145/3680256.3721316](https://doi.org/10.1145/3680256.3721316).
- [22] Perin, G., Wu, L., & Picek, S. (2022). Exploring feature selection scenarios for deep learning-based side-channel analysis. *IACR Transactions on Cryptographic Hardware and Embedded Systems*, 2022(4), 828-861. doi: [10.46586/tches.v2022.i4.828-861](https://doi.org/10.46586/tches.v2022.i4.828-861).
- [23] Salomon, D., Weiss, A., & Levi, I. (2021). Improved filtering techniques for single- and multi-trace side-channel analysis. *Cryptography*, 5(3), article number 24. doi: [10.3390/cryptography5030024](https://doi.org/10.3390/cryptography5030024).
- [24] Sanchez, P.M.S., Valero, J.M.J., Celdran, A.H., Bovet, G., Perez, M.G., & Perez, G.M. (2021). A survey on device behavior fingerprinting: Data sources, techniques, application scenarios, and datasets. *IEEE Communications Surveys & Tutorials*, 23(2), 1048-1077. doi: [10.1109/comst.2021.3064259](https://doi.org/10.1109/comst.2021.3064259).
- [25] Shuangde, L., Yuanjian, L., Leke, L., Xijia, B., Guyue, Z., Yaxin, Y., Anwen, R., & Qi, S. (2021). Millimeter wave channel characteristics of outdoor microcellular based on improved ray tracing method and BP neural network algorithm. *Chinese Journal of Radio Science*, 36(3), 430-442. doi: [10.12265/j.cjors.2020217](https://doi.org/10.12265/j.cjors.2020217).
- [26] Steinmann, P., Auping, W.L., & Kwakkel, J.H. (2020). Behavior-based scenario discovery using time series clustering. *Technological Forecasting and Social Change*, 156, article number 120052. doi: [10.1016/j.techfore.2020.120052](https://doi.org/10.1016/j.techfore.2020.120052).
- [27] Thakur, A., & Chandak, M.B. (2022). A review on opentelemetry and HTTP implementation. *International Journal of Health Sciences*, 6(S2), 15013-15023. doi: [10.53730/ijhs.v6ns2.8972](https://doi.org/10.53730/ijhs.v6ns2.8972).
- [28] Tian, Y., Li, H., Zhu, Q., Mao, K., Ali, F., Chen, X., & Zhong, W. (2024). Generative network-based channel modeling and generation for air-to-ground communication scenarios. *IEEE Communications Letters*, 28(4), 892-896. doi: [10.1109/lcomm.2024.3363621](https://doi.org/10.1109/lcomm.2024.3363621).
- [29] Wang, C., Lv, Z., Gao, X., You, X., Hao, Y., & Haas, H. (2022). Pervasive wireless channel modeling theory and applications to 6G GBSMs for all frequency bands and all scenarios. *IEEE Transactions on Vehicular Technology*, 71(9), 9159-9173. doi: [10.1109/tvt.2022.3179695](https://doi.org/10.1109/tvt.2022.3179695).
- [30] Xiong, L., Yao, Z., Miao, H., & Ai, B. (2021). Vehicle-to-vehicle channel characterization based on ray-tracing for urban road scenarios. *Wireless Communications and Mobile Computing*. doi: [10.1155/2021/8854247](https://doi.org/10.1155/2021/8854247).
- [31] Yang, M., et al. (2020). Machine-learning-based scenario identification using channel characteristics in intelligent vehicular communications. *IEEE Transactions on Intelligent Transportation Systems*, 22(7), 3961-3974. doi: [10.1109/tits.2020.3001132](https://doi.org/10.1109/tits.2020.3001132).
- [32] Zhang, J., & Wu, X. (2024). An anti-jamming game between dynamically-sensing jammer and legitimate user with faking-slot transmission. *IEEE Transactions on Vehicular Technology*, 73(7), 10287-10300. doi: [10.1109/tvt.2024.3372969](https://doi.org/10.1109/tvt.2024.3372969).
- [33] Zhang, J., Lin, J., Tang, P., Fan, W., Yuan, Z., Liu, X., Xu, H., Lyu, Y., Tian, L., & Zhang, P. (2024). Deterministic ray tracing: A promising approach to THz channel modeling in 6G deployment scenarios. *IEEE Communications Magazine*, 62(2), 48-54. doi: [10.1109/mcom.001.2200486](https://doi.org/10.1109/mcom.001.2200486).
- [34] Zhang, J., Liu, L., Fan, Y., Zhuang, L., Zhou, T., & Piao, Z. (2020). Wireless channel propagation scenarios identification: A perspective of machine learning. *IEEE Access*, 8, 47797-47806. doi: [10.1109/access.2020.2979220](https://doi.org/10.1109/access.2020.2979220).
- [35] Zheng, Z., Trivedi, K.S., Wang, N., & Qiu, K. (2017). Markov regenerative models of web servers for their user-perceived availability and bottlenecks. *IEEE Transactions on Dependable and Secure Computing*, 17(1), 92-105. doi: [10.1109/tdsc.2017.2753803](https://doi.org/10.1109/tdsc.2017.2753803).
- [36] Zhou, W., Borjigin, A., & He, C. (2025). Behavioral Universe Network (BUN): A behavioral information-based framework for complex systems. *ArXiv*. doi: [10.48550/arXiv.2504.15146](https://doi.org/10.48550/arXiv.2504.15146).

Реактивне трасування поведінкових сценаріїв в односторінкових додатках через інтеграцію Bun-базованих WebSocket-каналів та OpenTelemetry

Владислав Ананченко

Аспірант

Міжнародний економіко-гуманітарний університет імені академіка Степана Дем'янчука
33027, вул. Степана Дем'янчука, 4, м. Рівне, Україна
<https://orcid.org/0009-0004-8963-775X>

Юрій Лотюк

Кандидат педагогічних наук, доцент

Міжнародний економіко-гуманітарний університет імені академіка Степана Дем'янчука
33027, вул. Степана Дем'янчука, 4, м. Рівне, Україна
<https://orcid.org/0000-0001-6696-5583>

Анотація. Метою дослідження було оцінити часову ефективність реактивного трасування поведінки користувача в односторінкових додатках через інтеграцію WebSocket-каналів на базі Bun з OpenTelemetry. Методологія включала створення прототипу додатка на React, високочастотний моніторинг і агрегування SCADA-даних, побудову та оптимізацію нейромережі 64-32-16 у TensorFlow, симуляції в MATLAB/Simscapе, а також статистичний аналіз із використанням регресії Theil-Sen, Seasonal and Trend-декомпозиції, тесту Брауна-Форсайта, двофакторного аналізу варіантів, бутстреп-перестановки, критерія Дікі-Фуллера та кривих виживання Каплана-Мейєра. Результати показали, що комбінація протоколу гіпертекстової передачі із бінарною серіалізацією у форматі Protocol Buffers забезпечила найнижчу затримку виявлення подій, яка становила в середньому 45,09 мілісекунди, та найнижчу затримку передачі, що сягала лише 62,83 мілісекунди у сценарії заповнення форм. У той же час комбінація вебсокетів із текстовим форматом JavaScript Object Notation продемонструвала найвищі показники затримки, із середнім значенням виявлення подій 69,99 мілісекунди та затримкою передачі до 88,1 мілісекунди, а також найбільшу варіативність у часі реакції. Статистичний аналіз підтвердив суттєві відмінності між усіма конфігураціями: результати дисперсійного аналізу виявили надзвичайно високі значення F-статистики для обох показників із рівнем значущості p меншим за 0,000001, що свідчить про реальний вплив як протоколу, так і формату серіалізації на часову ефективність. Додатково встановлено, що затримка виявлення подій та затримка передачі були незалежними величинами, оскільки коефіцієнти кореляції в усіх випадках залишалися близькими до нуля. Таким чином, оптимальною конфігурацією для високочастотних телеметричних систем був гіпертекстовий протокол із бінарним форматом Protocol Buffers, що забезпечує не лише мінімальні часові затримки, але й стабільність у навантажених середовищах. Практична значимість результатів полягає в можливості використання їх інженерами з продуктивності, архітекторами фронтенду та розробниками моніторингових систем для створення ефективних та масштабованих рішень, орієнтованих на аналіз поведінки користувачів у режимі реального часу

Ключові слова: часові затримки; асинхронність; бінарна серіалізація; обробка подій; HTTP-JSON; телеметрична архітектура



High availability in a microservice architecture

Bohdan Fedoryshyn*

Postgraduate Student

Lviv Polytechnic National University

79000, 12 Stepana Bandery Str., Lviv, Ukraine

<https://orcid.org/0009-0005-3779-0186>

Abstract. The purpose of this study was to investigate approaches to ensuring high availability of microservice systems with a focus on fault tolerance, scalability, and continuous operation of services. The study applied a comparative and analytical method to analyse technical solutions for ensuring high availability, systematise the characteristics of container orchestration platforms, and evaluate load balancing tools according to the criteria of performance, flexibility, reliability, and integration convenience. Fault-tolerance patterns – retry, circuit breaker, and fallback – that provide flexible error management, reduce the risk of cascading failures, and maintain system continuity are investigated. The study found that the behaviour of fault-tolerance patterns depends on the configuration of execution parameters, such as timeouts, retry limits, and conditions for activating fallback mechanisms. The effectiveness of such tools as NGINX, HAProxy, Envoy, and Amazon Web Services Elastic Load Balancing is evaluated in terms of their impact on the scalability and resilience of the architecture, as well as the possibility of automatic scaling on the example of Amazon Web Services and Google Cloud platforms. It was found that built-in autoscaling services ensure stable operation of services under variable load and enable a rapid response to peak loads. An overview of container orchestrators (Kubernetes, OpenShift, Amazon ES) was provided, among which Kubernetes is recognised as the most effective due to the support of self-healing mechanisms, distributed deployment, health checks, and integration with Continuous Integration/Continuous Delivery. The findings of this study can serve as an analytical basis for designing sustainable microservice architectures in cloud and enterprise environments to improve the reliability, scalability, and continuity of business processes

Keywords: container orchestration platform; distributed deployment; scaling; monitoring; load balancing

INTRODUCTION

In the context of digital business transformation, microservice architecture has become one of the key approaches to building scalable and flexible software systems. Thanks to their distributed structure, where each service performs a separate function, microservices are easy to develop, deploy, and maintain. However, with distribution comes increased demands on system resilience and uptime. High availability becomes a major factor, as even a short-term downtime can lead to extensive financial losses or a degraded user experience. However, ensuring high availability of such distributed systems is a complex engineering task due to the considerable number of interdependent components, the possibility of partial

failures, and the requirement for continuous operation around the clock. It also involves choosing the right container orchestration tools, setting up monitoring systems, and implementing distributed deployment strategies. Despite the availability of powerful platforms such as Kubernetes or OpenShift, in practice their capabilities are not always fully utilised due to a lack of experience or a lack of a systematic approach.

In the field of microservice architecture, a significant challenge is to ensure high availability of the system while maintaining scalability and management efficiency. M.K. Foka (2024) analysed the key advantages of microservice architecture, including development

Article's History: Received: 05.07.2025; Revised: 25.10.2025; Accepted: 15.12.2025; Published: 25.12.2025.

Suggested Citation:

Fedoryshyn, B. (2025). High availability in a microservice architecture. *Bulletin of Cherkasy State Technological University*, 30(4), 155-165. doi: 10.62660/bcstu/4.2025.155.

*Corresponding author



Copyright © The Author(s). This is an open access article distributed under the terms of the Creative Commons Attribution License 4.0 (<https://creativecommons.org/licenses/by/4.0/>)

flexibility, simplified updates, and the ability to work in parallel by teams. The study showed that automatic scaling and health checks markedly increase the resilience of Web applications, but it was noted that the implementation of these mechanisms is often complicated due to the complexity of the interaction between microservices. T. Seliviorstrova & N. Krasnoshapka (2023) described the role of automatic monitoring and health checks in improving system reliability. The study confirmed the effectiveness of using Kubernetes and analogous platforms to maintain high availability. The specific features of integrating automatic scaling, monitoring, load balancing, and disaster recovery mechanisms should be studied.

In a microservice architecture, ensuring high availability is challenging due to the need to effectively coordinate numerous services, ensure load balancing, maintain data consistency, and resilience to possible failures in the distributed infrastructure. J. Li (2025) proposed an optimisation model to increase the availability and flexibility of blockchain systems built on a microservice architecture. The researcher developed an approach to load balancing and service redundancy, which allows increasing the overall system resilience to failures. D. Rossi (2020) studied the Consistency Availability Partition (CAP-theorem) dilemma in microservice systems, analysing the trade-offs between consistency and availability. The results showed that achieving high availability often requires a reduction in the level of consistency, which affects the consistency of data in distributed systems.

One of the challenges is the complexity of choosing the most suitable tools, frameworks, and platforms to ensure high availability, service consistency, and scalability of a microservice architecture under variable loads. G. Márquez *et al.* (2020) analysed the most popular frameworks and technological solutions used in industrial microservice systems to achieve high availability. The researchers found that most companies use off-the-shelf orchestrators (e.g., Kubernetes) and monitoring tools, but there are challenges with their proper configuration and integration into concrete business cases. G. Liu *et al.* (2020) reviewed the general challenges of microservice architecture, including high availability, containerisation, and service management. The researchers highlighted the complexity of debugging services in a distributed environment, the need to closely monitor the status of services, and manage their interaction as key issues.

Ensuring high availability in a microservice architecture is complicated by the need to scale, maintain performance, system resilience, and integrate numerous services in large distributed and cloud environments. N. Suleiman & Y. Murtaza (2024) reviewed comprehensive strategies for scaling microservices in enterprise applications, with a focus on increasing availability, optimising performance, and system resilience. The researchers highlighted the significance of horizontal

scaling to adapt to load changes and emphasised the role of dynamic resource balancing in preventing overload. The researchers also analysed approaches to automatic scaling in cloud environments, considering the specifics of microservice architecture. L. Roda-Sanchez *et al.* (2023) proposed a comprehensive cloud-edge microservice architecture using service orchestration that improves the real-world conditions for deploying and managing services. The researchers emphasised the significance of hybrid interaction between cloud infrastructure and edge devices, which enabled increased performance and reduced delays. Despite the contribution to the development of approaches to ensuring high availability of microservice systems, the analysed studies are mostly applied or narrowly focused and do not cover the complex interaction between scaling, monitoring, load balancing, and automatic recovery mechanisms. Additionally, most studies do not account for the challenges associated with integrating these solutions in large-scale distributed cloud environments with high load dynamics.

The purpose of this study was to substantiate solutions to ensure the reliability, fault tolerance, and scalability of microservice systems, considering the requirements for their continuous operation. To fulfil the purpose of this study, the following tasks were set: studying existing approaches to ensuring the resilience of microservice systems; analysing methods of automatic scaling and load balancing; researching tools and platforms for microservice orchestration.

MATERIALS AND METHODS

The study of microservice architectures and mechanisms for increasing their fault tolerance and performance was performed through the consistent application of four principal methodological approaches, which included analysis, comparison, and systematisation of technical characteristics and functional capabilities of key technologies. The first stage of the study was based on a functional and comparative analysis of the key fault-tolerance patterns in microservice architectures. Specifically, the study considered the retry, circuit breaker, and fallback patterns, which are fundamental to ensuring the continuity and stability of the system in the face of temporary and long-term failures. The method involved a detailed study of the functional principles of each pattern, their areas of application, and their impact on the stability and response time of the system.

The second stage focused on comparing load balancing systems used in microservice environments. The choice was focused on three key tools – NGINX, HAProxy, and Envoy – due to their ability to provide interactivity, gamification, and adaptive learning. They allow creating dynamic tasks, track student progress, analyse results, and personalise the learning process. This helps to increase motivation, engagement, and efficiency of learning, which directly corresponds to the purpose of

this study. Additionally, the Amazon Web Services (AWS) Elastic Load Balancer cloud solution is considered, the analysis of which covered such balancing models as client, server, Domain Name System (DNS) balancing, Application Programming Interface (API) Gateway, and service mash, based on official guides and typical use cases. A comparative review of Istio, Linkerd, and Consul Connect technologies was used to analyse the possibilities of implementing load balancing using a service mash. The review covered typical functionalities, including centralised traffic management, support for observability, authentication, and security in a dynamic microservices environment. A documented analysis of the implementation of automatic scaling mechanisms in leading cloud platforms was performed: AWS Auto Scaling, Google Cloud Autoscaler. The method of comparative analysis helped to assess their ability to dynamically respond to load changes, ensure fault tolerance, and rational use of resources.

At the final stage, the study reviewed the functionality of modern platforms that support microservice application lifecycle management, with a focus on deployment, scaling, monitoring, and automatic recovery tools. To investigate the life cycle management of microservice applications, the method of comparative analysis of the technical characteristics of tools and platforms was employed, which allowed structuring their functionality according to the key criteria of modern container orchestration platforms Kubernetes, OpenShift, and Amazon Elastic Container Service (ECS). The analysis included a review of the functionality of these platforms in terms of automatic scaling, health checks, monitoring, and distributed service deployment.

Separately, the study reviewed the documentation on health checks mechanisms – liveness and readiness probes – used in orchestration platforms for prompt detection of faults and automatic restart of problematic components, with examples of application in Kubernetes. Integration of real-time monitoring with microservice lifecycle management to maintain continuous system availability is applied. The distributed deployment of microservices on physical and virtual nodes, including geographical distribution in different data centres or regions, avoids single points of failure, and increases scalability. Overall, the application of these methods helped to obtain a comprehensive and systematic overview of technologies, their advantages and limitations that affect the resilience, scalability, and performance of microservice architectures. The findings obtained became the analytical basis for developing recommendations for implementing the most suitable solutions in corporate and cloud environments.

RESULTS

The key factors in the resilience of microservice systems are the use of fault tolerance patterns (retry, circuit breaker, fallback), isolation of services in containers to avoid cascading failures, automatic recovery through

self-healing mechanisms, monitoring and health checks, and load balancing (Paz & Bernardino, 2018). In microservice architectures, the key patterns of fault tolerance are retry, circuit breaker, and fallback, each of which is used depending on the type and duration of failures that occur in the system (Raj & David, 2021). The retry pattern is used in situations of temporary service unavailability or short-term network errors when the probability of successful re-execution of the operation is high. In case of a failure, the operation is automatically repeated several times at specified intervals. This approach reduces the probability of failure due to temporary problems, but it can also increase system response times, especially if the number of retries and delays are not optimally configured.

In practice, a properly configured retry helps to avoid unnecessary failures without significantly degrading performance, increasing overall system stability. The circuit breaker pattern is used to protect the system from prolonged and stable service failures. It monitors the number of errors over a certain period and, if it exceeds a threshold, temporarily blocks further calls to the problematic service, putting it in an “open” state. In this state, no requests are sent, which gives the service time to recover. After a certain pause, the circuit breaker switches to the “half-open” state, checking whether the service has recovered, and if so, returns to normal operation. This pattern reduces the load on the faulty service and prevents failures from spreading to other system components. The impact on response time in the normal state is minimal, but during activation, there may be a delay due to temporary call blocking. Overall, circuit breakers increase system reliability, especially in case of long-term failures. The fallback mechanism is used as a backup option in cases where the main service is unavailable or malfunctions. In such situations, the system executes alternative logic – e.g., it returns cached data, a standard response, or a message about the temporary unavailability of the service. This allows preserving basic functionality and maintain the user experience even during partial failures. In terms of response time, fallback can be faster than waiting for the primary service to be restored, but with certain limitations in terms of data relevance or completeness. The use of fallback is a significant component of a fault tolerance strategy, as it ensures the continuity of the system. The generalised characteristics of the key fault tolerance patterns that contribute to the stability of microservice systems are presented in Table 1. It reflects their functions, purpose, and implementation features.

The integrated implementation of Retry, Circuit Breaker, and Fallback patterns allows increasing the resilience and availability of microservice architecture, minimising the consequences of failures and ensuring stable system operation in dynamic environments. Service isolation in a microservice architecture ensures component independence by placing each service in its individual container or runtime environment. This

allows localising faults, avoiding cascading failures, and provides flexibility in updating and scaling individual services without affecting the entire system. Self-healing is implemented through mechanisms for automatic restart, replacement of faulty instances, and traffic

redirection, which minimises downtime and maintains system stability without operator intervention. Orchestration platforms, such as Kubernetes, provide built-in health checks that automatically monitor the health of services and initiate their recovery in case of a failure.

Table 1. Key patterns of fault tolerance in microservices

Pattern	Description	Purpose	Setup features
Retry	Repeat operation or request in case of temporary failure or unavailability	Reduction of the impact of temporary errors	Configuration of the number of attempts and intervals to avoid overload
Circuit Breaker	Mechanism that blocks calls to the service when a stable fault is detected	Prevention of cascading failures and overloading	Transitions between states: closed, open, half-open; error thresholds
Fallback	Alternative logic or fallback in case of unavailability of the main service	Maintenance of functionality during failures	Use of cached data or messages for the user

Source: compiled by the author of this study based on F. De Souza Miranda *et al.* (2024)

Monitoring and health checks are key mechanisms that ensure the continuity and stability of microservice systems (Waseem *et al.*, 2021). Monitoring allows tracking the status of each service in real time, analysing performance, detecting anomalies, and potential problems before they lead to failures. In practice, monitoring systems such as Prometheus or Grafana provide visualisation of key metrics such as Central Processing Unit (CPU), memory, response time, which allows quickly detecting anomalies related to overload or “hanging” requests. Collecting metrics, logs, and call tracing ensures transparency of the system’s internal processes, which enables prompt response to malfunctions and optimisation of performance. Health checks automate the monitoring of the health of individual components. They periodically check the availability and correctness of services, detecting failures or degradation. For example, the readiness-probe in Kubernetes allows excluding a component from balancing if it is not yet ready to process requests, while the liveness-probe restarts a container in case of a freeze. These checks reduce the average downtime of one instance to a few seconds without the need for an engineer. If a problem is detected, health checks can initiate a service restart or notify the orchestrator to redirect traffic to healthy instances, which prevents the failure from spreading to the entire system. Such mechanisms allow quickly isolating faults and minimising downtime. Implementation of monitoring and health checks helps to increase the reliability and availability of microservice systems, promptly detect problems, reduce the risk of cascading failures, and ensure stable operation even in case of partial failures.

The distributed deployment of microservices is a crucial factor in ensuring high system availability, as it avoids a single point of failure (Kansal & Balasubramaniam, 2024). Due to the fact that services run on different physical or virtual nodes, the failure of a single component does not lead to a complete shutdown of the system. This approach ensures resource redundancy and

load balancing, which contributes to greater resilience to failures. Additionally, distributed deployment enables the geographical location of services in different data centres or regions, which increases resilience to localised disruptions, such as power outages or network problems. In practice, this means that delays or failures in one region do not affect the overall performance of the system if the routing of requests between regions is properly configured. This reduces downtime and improves overall system reliability. Implementing a distributed deployment also makes it easier to scale the system, as one can independently add resources where needed without affecting other components. Thus, distributed deployment is an effective means of increasing the availability and resilience of microservice architectures.

One of the key conditions for ensuring high availability and resilience of a microservice architecture is effective load balancing between service instances (Barua & Kaiser, 2024). This allows the system to evenly distribute requests, avoid overloading individual nodes, and ensure real-time scalability. There are several key approaches to implementing load balancing, each of which has its specific features, advantages, and applications. Client-based load balancing implies that the logic of selecting a service instance is delegated to the client. The client has a list of available instances (e.g., obtained through a service discovery service) and decides which one to send a request to. This approach reduces the load on the central components but requires a more complex implementation on the client side and up-to-date information about the instances. Examples of client-side balancing include the use of Netflix Ribbon or Spring Cloud LoadBalancer in microservice applications, when the client receives a list of instances from Eureka, Consul, or ZooKeeper and selects them by an algorithm (e.g., round-robin or random). Google Remote Procedure Call also has built-in client balancing mechanisms, where the client library distributes requests between known instances. En-

voy employs an analogous approach in sidecar mode, where balancing is performed directly at the proxy level, which is integrated into the client loop. Server balancing is implemented through a separate component, the load balancer, which independently routes requests to the relevant instance. For this, special tools such as NGINX, HAProxy, Envoy, or cloud solutions such as AWS Elastic Load Balancing (ELB) are used. The advantage of this approach is centralised request management, but there is a risk of a bottleneck if the balancer is not scalable or fails.

Another relevant approach to organising load balancing is the use of a service mash, which is implemented by introducing proxy servers (sidecars) to each microservice. Istio, Linkerd, and Consul Connect provide automatic traffic redirection based on performance metrics, centralised control, tracing, monitoring, authentication, and other service functions without

changing business logic. Istio provides the broadest configuration and integration capabilities with Kubernetes, Linkerd is distinguished by its simplicity and minimal overheads, while Consul Connect focuses on integration with HashiCorp infrastructure and support for hybrid environments. All these solutions increase the reliability and manageability of systems, but require extra computing resources and complex configuration, which should be considered when designing the architecture. To better compare the key load balancing tools, the study analysed their characteristics according to the key criteria: reliability, routing flexibility, configuration complexity, and compatibility with other services. This enables a clearer understanding of the strengths and weaknesses of each solution and helps to choose the best tool depending on the specific needs and features of the microservice architecture. The results of this comparison are presented in Table 2.

Table 2. Comparison of load balancing tools by key characteristics

Tool	Reliability	Routing flexibility	Customisation complexity	Compatibility with other services
NGINX	High, proven over years	Support for basic and advanced routing HyperText Transfer Protocol (HTTP), Transmission Control Protocol (TCP)	Moderate, requires configuration knowledge	Extensive module support, integration with CI/CD
HAProxy	Very high, widely used	High, support for complex balancing rules	Medium, configuration via configuration files	Solid integration with various protocols and monitoring
Envoy	Highly cloud-focused	Very high, support for L3-L7 routing, service mesh	High, requires time for configuration and training	Perfectly integrates with service mash-up platforms
AWS ELB	High, automatic scaling and high availability	Limited compared to NGINX and Envoy, standard methods	Low, manageable configuration via AWS console	Tight integration with AWS ecosystem and services
Istio	High, centralised traffic management, fault-tolerance	Very high, support for intelligent routing based on performance metrics	High, requires knowledge of Kubernetes and configurations	Support for Kubernetes, Prometheus, Jaeger, and other monitoring services
Linkerd	High, simple and lightweight solution	High-level, basic routing for services	Low, easy installation and basic configuration	Solid integration with Kubernetes, Prometheus
Consul Connect	Highly centralised service management	High-performance routing with security policy support	Moderate, requires customisation of the HashiCorp infrastructure	Integration with HashiCorp Vault, Terraform, and other HashiCorp services

Notes: CI/CD – Continuous Integration/Continuous Delivery

Source: compiled by the author of this study

The comparison shows that all of the tools under consideration provide a high level of reliability but differ in terms of routing flexibility and configuration complexity. Envoy is ideal for complex cloud and service-mix environments that require high adaptability, although it is more resource intensive to implement. NGINX and HAProxy offer an optimum balance between functionality and simplicity and are well suited for traditional systems with a medium workload. AWS ELB is best used in the AWS ecosystem, valuing ease

of configuration and deep integration, although with some routing limitations. The choice of a particular tool should be based on infrastructure requirements and project specifics.

Automatic scaling methods are key to ensuring high availability of microservice systems, as they enable the number of service instances to be dynamically adapted to the current load. The key approaches to automatic scaling are horizontal and vertical scaling. Horizontal scaling involves adding or removing instances of

the same service, which allows quickly responding to changes in load and reduces the risk of overloading individual components. Vertical scaling involves changing the CPU, Random Access Memory (RAM) resources of individual instances, which is limited by the physical characteristics of the infrastructure but can be effective for short-term peaks. In cloud environments, automatic scaling is implemented through specialised services that provide an adaptive response to load changes. Specifically, AWS Auto Scaling and Google Cloud Autoscaler are the leading solutions that are widely used to maintain stable operation of microservices.

AWS Auto Scaling works on the principle of continuous monitoring of key metrics such as CPU, number of requests, memory usage, and others. Based on the set thresholds and scaling policies, the service automatically adds or removes instances, maintaining the optimum level of resources. Another major advantage is the tight integration with other AWS services, such as Elastic Load Balancer and CloudWatch, which provides detailed monitoring, alerts, and flexible scaling management. This enables quick response to peak loads and ensures traffic balancing between active instances. The key limitations are a strong dependence on the AWS ecosystem and the complexity of settings when dealing with complex scaling scenarios. Google Cloud Autoscaler supports scaling for both virtual machines and containerised environments in Kubernetes. It monitors CPU, memory, network load, and even custom metrics, enabling precise resource allocation to meet workload demands. Thanks to integration with the Google Kubernetes Engine (GKE), the service provides flexible horizontal container scaling with high response speed. The advantage is support for various types of metrics and fine-tuning of thresholds, which avoids over- or under-scaling. The disadvantage may be a slightly more complex configuration compared to other providers, as well as the need for additional Kubernetes knowledge to fully utilise the features. Both services support integration with health checks and load balancing mechanisms, which allows not only scaling resources but also ensuring business continuity by quickly replacing faulty instances. The choice between these platforms is usually based on the chosen cloud provider, the specifics of the infrastructure, and the scaling needs of the project.

A variety of container orchestration platforms are used to ensure high availability, scalability, and efficient lifecycle management of microservices. Kubernetes is a leading container orchestration system that provides extensive capabilities for automating the deployment, scaling, and management of microservice applications (Zhou *et al.*, 2021). Its architecture provides high availability due to mechanisms for automatically checking the state of containers through liveness and readiness probes, which enable the rapid detection of inoperable service instances. Scaling is performed using Horizontal Pod Autoscaler, which responds to the load by measuring

resource consumption (CPU, memory) or custom metrics. Additionally, support for geo-distributed deployment across multiple availability zones avoids a single point of failure and thus increases system resilience. Kubernetes integrates closely with CI/CD processes (e.g., through GitLab CI or Jenkins), supports Infrastructure as a Code, and can work in conjunction with service mashups such as Istio, which provides advanced traffic, security, and observability management (Mustyala, 2022). Due to its modularity and distributed nature, Kubernetes is considered the most complete and flexible solution for building a scalable, isolated, and self-healing microservice architecture.

OpenShift is an extension of Kubernetes aimed at enterprise use and includes advanced tools for improved management, security, and application lifecycle automation. Apart from the basic features of Kubernetes, OpenShift integrates built-in CI/CD support through OpenShift Pipelines, providing controlled and repeatable application deployment. The platform features enhanced security policies implemented through the Security Context Constraints (SCC) and Role-Based Access Control (RBAC) mechanisms, which enable granular access control to resources. The monitoring and logging system is based on the Prometheus, Grafana, and EFK (Elasticsearch, Fluentd, Kibana) stack, which provides deep observability. Service scaling is supported by built-in controllers, analogous to Kubernetes. Due to its extensive management capabilities, advanced security, and centralised approach to administration, OpenShift is a reasonable choice for large organisations that require not only high availability but also strict control over infrastructure, access policies, and auditing (Al-Harbi & Al-Qahtani, 2024).

Amazon ECS is a managed container orchestration service on the AWS cloud infrastructure focused on simplifying the deployment and management of microservices. ECS provides basic availability mechanisms through health checks, automatic scaling through AWS Auto Scaling, and load balancing through AWS Elastic Load Balancer. All these features are tightly integrated with other AWS services, such as Identity and Access Management for access control, Virtual Private Cloud for network isolation, and Secrets Manager for secure storage of confidential information. Monitoring is implemented through Amazon CloudWatch, which allows creating triggers and alerts based on load metrics, response time, or number of requests. One of the key advantages of ECS is the minimal complexity of setup – the platform hides most of the technical details, which allows deploying a stable and scalable environment with minimal effort. However, this convenience is accompanied by a deep tie to the AWS ecosystem, which can limit flexibility in multi-cloud or hybrid scenarios (Bagai, 2024). Table 3 presents the key features of each platform, such as automatic scaling, load balancing, monitoring, and implementation features that directly affect the resilience and reliability of services.

Table 3. Comparison of container orchestration platforms

Platform	Scaling	Load balancing	Health Checks	Features	Application
Kubernetes	Automatic horizontal scaling	Built-in, supports L4 and L7 levels	Liveness, readiness, startup probes	The most popular, with a developed ecosystem, supports Helm, and has a large community	Suitable for all environments, especially in large distributed systems
OpenShift	Automatic, improved	Built-in, extended	Advanced features	Security, CI/CD, access control	Enterprises, corporate solutions
Amazon ECS	Automatic scaling with AWS Auto Scaling	Built-in, integrated	Supported	Tight integration with AWS infrastructure	Cloud applications, scalable services

Source: compiled by the author of this study based on E. Truyen *et al.* (2019), A.M. Kovalenko (2021)

The integrated use of these platforms allows optimising the deployment of microservices, ensuring their continuity, and responding quickly to failures, which ultimately increases the overall availability and efficiency of the system. Based on this analysis, the following practical recommendations for implementing high availability in microservice architectures can be formulated. It is worth implementing well-known fault tolerance patterns – retry, circuit breaker, and fallback – that allow the system to automatically respond to temporary failures, avoid excessive loads on problematic components, and provide alternative ways to process requests. Correctly configuring the parameters of these patterns is critical to achieving the optimal balance between reliability and performance. One should also use a distributed deployment of microservices across different physical or virtual nodes, as well as in different geographical regions. This approach eliminates a single point of failure, increases resilience to localised failures, reduces the risk of large-scale disruptions, and facilitates faster system recovery. Furthermore, it is vital to implement automatic scaling mechanisms that allow adapting resources to the current load without disrupting services. This ensures system flexibility and saves resources while maintaining stable operation even at peak times. Along with automatic scaling, it is necessary to implement effective load balancing between service instances, which prevents overloading of individual nodes and increases overall system performance.

Monitoring and health checks should be integrated into all levels of the architecture to promptly detect failures and automatically launch recovery processes. This allows promptly responding to problems, minimising downtime, and improving user experience. The use of modern container orchestration platforms such as Kubernetes, Docker Swarm, or OpenShift is a key element of high availability. These platforms support automated deployment, scaling, load balancing, and recovery of microservices, which greatly simplifies the management of complex systems. Finally, the implementation of continuous integration and delivery (CI/CD) practices ensures fast and secure service updates without interruptions, which increases the overall reliability and efficiency of software support.

Comprehensive adherence to these recommendations helps to increase the resilience, reliability, and availability of microservice architectures, minimises the risk of failure and ensures business continuity even in case of a malfunction. To summarise the results, high availability of microservice systems is achieved through a combination of service isolation, fault-tolerance patterns, autoscaling, load balancing, and self-healing mechanisms. The most effective solution in practice is to use Kubernetes with autoscaling, service mesh, and distributed deployment. In contrast, more simplistic approaches such as DNS balancing have limitations in scalability and adaptability.

DISCUSSION

The findings confirmed the significance of using fault tolerance patterns (retry, circuit breaker, fallback), isolation of services in containers, self-healing mechanisms, monitoring and health checks, and load balancing to ensure the resilience of microservice architectures. The implementation of these approaches has helped minimise the impact of failures, localise malfunctions, and ensure system continuity by automating recovery processes. These solutions help to increase system reliability and efficiency, reducing the need for human intervention and improving productivity. The results obtained are consistent with the provisions highlighted in M. Kuppam (2024), where the researcher emphasises the need to develop a resilient architecture through the introduction of self-healing mechanisms that minimise human intervention and automate the response to failures. These research findings have confirmed the practical effectiveness of such approaches when automatic restart or replacement of faulty instances helps to ensure continuous operation of the system. B. Arugula (2024) proposed a model for building resilient cloud-native APIs in event-driven microservice ecosystems that involves autonomous disaster recovery and adaptive request routing. These approaches echo the conclusions of this paper regarding the significance of load balancing between healthy instances, as well as the use of fallback mechanisms that allow the system to continue functioning even in case of individual service failures.

The findings confirmed that automatic scaling is a key mechanism for maintaining the resilience and adaptability of microservice architecture in the cloud environment. Modern services, such as AWS Auto Scaling and Google Cloud Autoscaler, demonstrate high efficiency by monitoring key metrics (CPU, memory, network load) and dynamically managing the number of instances based on predefined policies. Integration with load balancing and health checks ensures fast response to failures and minimises downtime. AWS Auto Scaling is distinguished by its close interaction with other components of the AWS ecosystem (CloudWatch, ELB), which contributes to flexible and centralised scale management. Google Cloud Autoscaler is deeply integrated with GKE and supports custom metrics, providing precise adaptation to load changes in Kubernetes clusters. Both approaches confirm the effectiveness of combining scaling with availability and load control, increasing the overall resilience of microservice infrastructure.

N. Singh *et al.* (2023) reached analogous conclusions, emphasising the role of load balancing mechanisms and service discovery in the Docker Swarm environment for distributed big data systems. The researchers noted the effectiveness of horizontal scaling in reducing the risk of overloading individual containers, which coincides with the conclusions drawn from the study on the benefits of horizontal scaling for flexible adaptation to load changes. S. Rabiou *et al.* (2022) highlighted the problems and challenges of load balancing and auto-scaling in cloud microservices environments, especially the limitations of vertical scaling and the need for dynamic resource management to ensure system resilience. These findings are also consistent with the findings of the present study, which showed that horizontal scaling in combination with service-mashups provides the best level of resilience and adaptability for microservice architectures.

These results confirm that Kubernetes is a leading container orchestration system that provides a prominent level of availability, scalability, and resilience for microservice architectures through deployment automation, horizontal scaling (via Horizontal Pod Autoscaler (HPA)), self-healing mechanisms, and geo-distribution of instances. Built-in health checks (liveness and readiness probes) enable rapid detection of failures, while integration with service mashups such as Istio enables advanced traffic, security, and availability management. The platform tightly integrates with CI/CD processes, supports infrastructure-as-a-service approaches, and is a flexible solution for building an isolated, scalable, and self-healing microservice architecture. In turn, OpenShift, as an enterprise-oriented extension of Kubernetes, offers additional security features (SCC, RBAC), centralised monitoring (EFK, Prometheus, Grafana), and native CI/CD integration through OpenShift Pipelines, making it a viable choice for large organisations with increased access control and audit requirements.

Amazon ECS, as a managed solution within the AWS ecosystem, provides a simplified container management model, automatic scaling, CloudWatch monitoring, and deep integration with other AWS services. ECS enables rapid deployment of resilient services with minimal configuration but is less flexible in multi-cloud environments due to its strong tie to AWS.

The results of this study are in line with the findings of A. Saboor *et al.* (2022), who emphasised the significance of efficient orchestration of containerised microservices to ensure high availability, scalability, and reliability of cloud systems. Specifically, the researchers focused on conceptual approaches to automated deployment and life cycle management of microservices, which is in line with the recommendations identified in the study on the use of orchestration platforms (Kubernetes, etc.) to improve service resilience. R.R. Vangala (2018) confirmed the significance of dynamic orchestration and automatic scaling as key factors in the adaptive resilience of microservice systems. The adaptive framework proposed by the researcher emphasised the need for rapid response to load changes and automatic recovery from failures, which is consistent with the results of the present study on the effectiveness of horizontal scaling and automatic load balancing to maintain system continuity.

Monitoring and health checks have played a key role in ensuring the continuity and stability of microservice systems. By being capable of detecting faults at an early stage and responding to them quickly, these mechanisms considerably reduce the risk of cascading failures and minimise downtime. Most significantly, automated health checks help to isolate problematic components and ensure that traffic is redirected to healthy instances, which increases overall system reliability. It also confirms the significance of distributed deployment as a key factor in the high availability of microservice architecture. Running services on multiple physical or virtual nodes eliminates a single point of failure, providing resource redundancy and load balancing. The geographical location of components in different data centres or regions further increases the system's resilience to local failures, which is especially relevant for critical applications with continuous operation requirements.

The findings are consistent with the approaches described in O.V. Talaver & T.A. Vakaliuk (2023), which emphasises the significance of comprehensive system health management to improve system reliability. The researchers emphasised the need to use automated tools to monitor service performance, which enabled prompt detection of failures and minimised the impact of individual component failures on the functioning of the entire system – a conclusion that is consistent with the emphasis on the significance of health checks in the study. Additionally, the findings of the present study are consistent with the findings of Y. Wang *et al.* (2021), who examined both the benefits and challenges of

microservice architectures. Specifically, authors emphasised that the scalability and resilience of systems directly depend on the proper design of distributed deployment and well-established mechanisms for self-healing services. This coincides with the conclusions about the need to avoid single points of failure by distributing the load across different nodes and geographical areas, which allows the system to stay operational even in case of partial failures.

Overall, the findings of the study confirmed that microservice architecture is a modern and effective approach to creating scalable, flexible, and adaptive software systems that can meet the requirements of a dynamic business environment. The comprehensive implementation of microservice architecture, which includes not only technical tools but also suitable management approaches, enables organisations to create highly reliable, flexible, and scalable software systems that can effectively respond to market changes, reduce support costs, and provide high quality customer service.

CONCLUSIONS

The study found that the stability, reliability, and continuity of microservice systems are achieved through the integrated implementation of a series of architectural approaches, mechanisms, and patterns. Service isolation, automatic recovery of instances, load balancing, and continuous monitoring form the foundation for ensuring stable operation even in case of partial failures or anomalies. Practical cases revealed that the use of fault tolerance patterns – retry, circuit breaker, and fallback – reduces the risk of cascading failures while optimising system response times through adaptive parameter settings. Specifically, retry reduces the probability of failure due to temporary problems. Circuit breaker effectively isolates faulty components, preventing excessive load, while fallback provides alternative request processing scenarios in case of prolonged failures.

Tools such as Consul provided efficient service management, detection, traffic routing, and implementation of fault tolerance patterns (retry, circuit breaker, fallback). In turn, the use of Istio as a service mashup enabled centralised management of inter-service interaction, load balancing, security, and observability. The introduction of distributed deployment of services in

different geographical locations considerably reduced the risk of large-scale failures, ensuring high availability and fast recovery of the system even in case of local failures. AWS and Google Cloud provided built-in services for automatic scaling, load balancing, redundancy, and monitoring, which greatly simplifies the construction of fault-tolerant systems with a prominent level of availability. Automatic scaling based on monitoring of key metrics (CPU, latency, number of requests) allowed dynamically adjusting resources to the current load, maintaining system stability at peak times and saving resources during downturns.

The prominent level of reliability of microservice systems was ensured by container orchestration platforms (Kubernetes, OpenShift, Amazon ECS) that automate application lifecycle management: deployment, scaling, monitoring, and recovery. They support load balancing, container isolation, self-healing, and horizontal scaling, which reduces the risk of failures. Integration with CI/CD and infrastructure as code increases stability and control of changes in the production environment. Comparative analysis revealed that Kubernetes is the most flexible and scalable solution, while OpenShift and Amazon ECS offer complementary enterprise and cloud integrations. The proposed recommendations will help to improve the efficiency, reliability, and resilience of microservice systems. Further research could focus on developing adaptive mechanisms for automatically responding to complex failures in microservice architectures using artificial intelligence and machine learning to predict possible failures and optimise resources. Another major area will be exploring the impact of various security and data protection strategies on the overall availability and resilience of systems, which will ensure both reliability and compliance with modern information security requirements.

ACKNOWLEDGEMENTS

None.

FUNDING

None.

CONFLICT OF INTEREST

None.

REFERENCES

- [1] Al-Harbi, F., & Al-Qahtani, A. (2024). [Software-defined storage \(SDS\): Architecture, benefits, and leading platforms](#). *International Journal of Informatics and Data Science Research*, 1(8), 36-49.
- [2] Arugula, B. (2024). Architecting for resilience: Designing fault-tolerant systems in multi-cloud environments. *International Journal of Emerging Trends in Computer Science and Information Technology*, 5(2), 113-121. [doi: 10.63282/3050-9246.IJETCSIT-V5I2P112](#).
- [3] Bagai, R. (2024). Comparative analysis of AWS model deployment services. *International Journal of Computer Trends and Technology*, 72(5), 102-110. [doi: 10.14445/22312803/ijctt-v72i5p113](#).
- [4] Barua, B., & Kaiser, M.S. (2024). Enhancing resilience and scalability in travel booking systems: A microservices approach to fault tolerance, load balancing, and service discovery. *ArXiv*. [doi: 10.48550/arXiv.2410.19701](#).

- [5] De Souza Miranda, F., dos Santos, D.S., Vilela, R.F., Guez Assunção, W.K., dos Santos, R.C., & Costa Pinto, V.H.S. (2024). A proposed catalog of development patterns for fault-tolerant microservices. In *SBQS '24: Proceedings of the XXIII Brazilian symposium on software quality* (pp. 406-416). New York: Association for Computing Machinery. doi: [10.1145/3701625.3701678](https://doi.org/10.1145/3701625.3701678).
- [6] Foka, M.K. (2024). *Research on the effectiveness and advantages of microservice architecture in Web applications*. (Master's thesis, Zaporizhzhia National University, Zaporizhzhia, Ukraine).
- [7] Kansal, S., & Balasubramaniam, V.S. (2024). Microservices architecture in large-scale distributed systems: Performance and efficiency gains. *Journal of Quantum Science and Technology (JQST)*, 1(4), 633-663. doi: [10.63345/jqst.v1i4.139](https://doi.org/10.63345/jqst.v1i4.139).
- [8] Kovalenko, A.M. (2021). *Methods and tools for auditing the security of the Kubernetes automatic container orchestration system*. (Masters's dissertation, Igor Sikorsky Kyiv Polytechnic Institute, Kyiv, Ukraine).
- [9] Kuppam, M. (2024). The resilient design techniques. In *Enterprise digital reliability* (pp. 87-115). Berkley: Apress. doi: [10.1007/979-8-8688-1032-9_4](https://doi.org/10.1007/979-8-8688-1032-9_4).
- [10] Li, J. (2025). Research on optimization model of high availability and flexibility of blockchain system based on microservice architecture. *Procedia Computer Science*, 261, 207-216. doi: [10.1016/j.procs.2025.04.191](https://doi.org/10.1016/j.procs.2025.04.191).
- [11] Liu, G., Huang, B., Liang, Z., Qin, M., Zhou, H., & Li, Z. (2020). Microservices: Architecture, container, and challenges. In *2020 IEEE 20th international conference on software quality, reliability and security companion (QRS-C)* (pp. 629-635). Macau: IEEE. doi: [10.1109/QRS-C51114.2020.00107](https://doi.org/10.1109/QRS-C51114.2020.00107).
- [12] Márquez, G., Soldani, J., Ponce, F., & Astudillo, H. (2020). *Frameworks and high-availability in microservices: An industrial survey*. In *Proceedings of the XXIII Ibero-American conference on software engineering (CIBSE)* (pp. 57-70). Montevideo: Curran Associates.
- [13] Mustyala, A. (2022). CI/CD pipelines in Kubernetes: Accelerating software development and deployment. *International Journal of Science and Engineering*, 8(3), 1-11. doi: [10.53555/epijise.v8i3.238](https://doi.org/10.53555/epijise.v8i3.238).
- [14] Paz, S., & Bernardino, J. (2018). Web platform assessment tools: An experimental evaluation. In T.A. Majchrzak, P. Traverso, K.-H. Krempels & V. Monfort (Eds.), *Web information systems and technologies* (pp. 45-63). Cham: Springer. doi: [10.1007/978-3-319-93527-0_3](https://doi.org/10.1007/978-3-319-93527-0_3).
- [15] Rabiou, S., Yong, C.H., & Mohamad, S.M.S. (2022). A cloud-based container microservices: A review on load-balancing and auto-scaling issues. *International Journal on Data Science*, 3(2), 80-92. doi: [10.18517/ijods.3.2.80-92.2022](https://doi.org/10.18517/ijods.3.2.80-92.2022).
- [16] Raj, P., & David, G.S.S. (2021). Engineering resilient microservices toward system reliability: The technologies and tools. In R. Achary & P. Raj (Eds.), *Cloud reliability engineering: Technologies and tools* (pp. 77-116). Bora Raton: CRC Press. doi: [10.1201/9781003030973-3](https://doi.org/10.1201/9781003030973-3).
- [17] Roda-Sanchez, L., Garrido-Hidalgo, C., Royo, F., Maté-Gómez, J.L., Olivares, T., & Fernández-Caballero, A. (2023). Cloud-edge microservices architecture and service orchestration: An integral solution for a real-world deployment experience. *Internet of Things*, 22, article number 100777. doi: [10.1016/j.iot.2023.100777](https://doi.org/10.1016/j.iot.2023.100777).
- [18] Rossi, D. (2020). Consistency and availability in microservice architectures. In M.J. Escalona, F.D. Mayo, T.A. Majchrzak & V. Monfort (Eds.), *Web information systems and technologies* (pp. 39-55). Cham: Springer. doi: [10.1007/978-3-030-35330-8_3](https://doi.org/10.1007/978-3-030-35330-8_3).
- [19] Saboor, A., Hassan, M.F., Akbar, R., Shah, S.N.M., Hassan, F., Magsi, S.A., & Siddiqui, M.A. (2022). Containerized microservices orchestration and provisioning in cloud computing: A conceptual framework and future perspectives. *Applied Sciences*, 12(12), article number 5793. doi: [10.3390/app12125793](https://doi.org/10.3390/app12125793).
- [20] Seliviorstrova, T., & Krasnoshapka, N. (2023). Aspects of designing scalable microservices architecture for web services. *Information Technology Computer Science Software Engineering and Cyber Security*, 4, 58-66. doi: [10.32782/it/2023-4-7](https://doi.org/10.32782/it/2023-4-7).
- [21] Singh, N., Hamid, Y., Juneja, S., Srivastava, G., Dhiman, G., Gadekallu, T.R., & Shah, M.A. (2023). Load balancing and service discovery using Docker Swarm for microservice based big data applications. *Journal of Cloud Computing Advances Systems and Applications*, 12, article number 4. doi: [10.1186/s13677-022-00358-7](https://doi.org/10.1186/s13677-022-00358-7).
- [22] Suleiman, N., & Murtaza, Y. (2024). *Scaling microservices for enterprise applications: Comprehensive strategies for achieving high availability, performance optimization, resilience, and seamless integration in large-scale distributed systems and complex cloud environments*. *Applied Research in Artificial Intelligence and Cloud Computing*, 7(6), 46-82.
- [23] Talaver, O.V., & Vakaliuk, T.A. (2023). Reliable distributed systems: Review of modern approaches. *Journal of Edge Computing*, 2(1), 84-101. doi: [10.55056/jec.586](https://doi.org/10.55056/jec.586).
- [24] Truyen, E., van Landuyt, D., Preuveneers, D., Lagaisse, B., & Joosen, W. (2019). A comprehensive feature comparison study of open-source container orchestration frameworks. *Applied Sciences*, 9(5), article number 931. doi: [10.3390/app9050931](https://doi.org/10.3390/app9050931).

- [25] Vangala, R.R. (2018). [Adaptive resilience framework: A comprehensive study on dynamic orchestration and auto-scaling of microservices in cloud-native systems](#). *International Journal of Computer Engineering and Technology*, 9(6), 278-288.
- [26] Wang, Y., Kadiyala, H., & Rubin, J. (2021). Promises and challenges of microservices: An exploratory study. *Empirical Software Engineering*, 26, article number 63. [doi: 10.1007/s10664-020-09910-y](#).
- [27] Waseem, M., Liang, P., Shahin, M., Di Salle, A., & Márquez, G. (2021). Design, monitoring, and testing of microservices systems: The practitioners' perspective. *Journal of Systems and Software*, 182, article number 111061. [doi: 10.1016/j.jss.2021.111061](#).
- [28] Zhou, N., Georgiou, Y., Pospieszny, M., Zhong, L., Zhou, H., Niethammer, C., Pejak, B., Marko, O., & Hoppe, D. (2021). Container orchestration on HPC systems through Kubernetes. *Journal of Cloud Computing Advances Systems and Applications*, 10, article number 16. [doi: 10.1186/s13677-021-00231-z](#).

Висока доступність в мікросервісній архітектурі

Богдан Федоришин

Аспірант

Національний університет «Львівська політехніка»

79000, вул. Степана Бандери, 12, м. Львів, Україна

<https://orcid.org/0009-0005-3779-0186>

Анотація. Метою роботи було дослідження підходів до забезпечення високої доступності мікросервісних систем з акцентом на стійкість до відмов, масштабування і безперервну роботу сервісів. У дослідженні застосовано порівняльно-аналітичний метод, аналізу технічних рішень забезпечення високої доступності, систематизації характеристик платформ оркестрації контейнерів і оцінки інструментів балансування навантаження за критеріями продуктивності, гнучкості, надійності та інтеграційної зручності. Досліджено fault-tolerance патерни – retry, circuit breaker і fallback – які забезпечують гнучке управління помилками, знижують ризик каскадних відмов і підтримують безперервність роботи систем. Встановлено, що поведінка fault-tolerance патернів залежить від конфігурації параметрів виконання, таких як таймаути, ліміти повторних спроб і умови активації fallback-механізмів. Оцінено ефективність таких інструментів, як NGINX, HAProxy, Envoy та Amazon Web Services Elastic Load Balancing, з огляду на їх вплив на масштабованість і стійкість архітектури, а також можливості автоматичного масштабування на прикладі хмарних платформ Amazon Web Services і Google Cloud. Виявлено, що вбудовані сервіси autoscaling забезпечують стабільну роботу сервісів при змінному навантаженні та дозволяють оперативно реагувати на пікові навантаження. Надано огляд оркестраторів контейнерів (Kubernetes, OpenShift, Amazon ECS), серед яких Kubernetes визнано найбільш ефективним завдяки підтримці механізмів самовідновлення, розподіленого розгортання, health checks та інтеграції з CI/CD. Результати дослідження можуть слугувати аналітичною основою для проектування стійких мікросервісних архітектур у хмарному та корпоративному середовищах з метою підвищення надійності, масштабованості та безперервності бізнес-процесів

Ключові слова: платформа оркестрації контейнерів; розподілене розгортання; масштабування; моніторинг; балансування навантаження

ВІСНИК
Черкаського державного технологічного університету

Науковий збірник

Том 30, № 4. 2025

Заснований у 1996 р. Виходить чотири рази на рік

Оригінал-макет видання виготовлено у редакційно-видавничому відділі
Черкаського державного технологічного університету

Відповідальний редактор:

А. Лавданський

Підписано до друку 15 грудня 2025 р.
Формат 60*84/8
Умов. друк. арк. 19,5
Наклад 50 прим.

Адреса видавництва:

Черкаський державний технологічний університет
18006, бульв. Шевченка, 460, м. Черкаси, Україна
Тел.: +380638013283
E-mail: info@bulletin-chstu.com.ua
<https://bulletin-chstu.com.ua/uk>

BULLETIN
of Cherkasy State Technological University

Scientific Collection

Volume 30, No. 4. 2025

Founded in 1996. Published four times per year

The original layout of the publication is made in the editorial and publishing department
of Cherkasy State Technological University

Managing editor:
A. Lavdanskyi

Signed for print of December 15, 2025.
Format 60*84/8
Conventional printed pages 19.5
Circulation 50 copies

Editors Office Address:

Cherkasy State Technological University
18006, 460 Shevchenko Blvd., Cherkasy, Ukraine
Tel.: +380638013283
E-mail: info@bulletin-chstu.com.ua
<https://bulletin-chstu.com.ua/en>